

On the Analysis of Complex Networks: Fundamental Limits, Scalable Algorithms, and Applications

by

Soheil Feizi

B.S., Sharif University of Technology (2008)

S.M., Massachusetts Institute of Technology (2010)

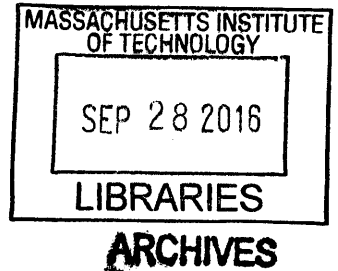
Submitted to the Department of Electrical Engineering and Computer Science

in partial fulfillment of the requirements for the degree of
Doctor of Philosophy

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

September 2016



© Massachusetts Institute of Technology 2016. All rights reserved.

Signature redacted

Author
Department of Electrical Engineering and Computer Science

August 5, 2016

Signature redacted

Certified by.....
Muriel Médard

Cecil H. Green Professor of Electrical Engineering and Computer Science

Signature redacted Thesis Supervisor

Certified by.....
Manolis Kellis

Professor of Electrical Engineering and Computer Science

Signature redacted Thesis Supervisor

Accepted by.....
1 00 Leslie Kolodziejcki

Chair, Department Committee on Graduate Students

On the Analysis of Complex Networks: Fundamental Limits, Scalable Algorithms, and Applications

by

Soheil Feizi

Submitted to the Department of Electrical Engineering and Computer Science
on August 5, 2016, in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy

Abstract

Network models provide a unifying framework for understanding dependencies among variables in data-driven and engineering sciences. Networks can be used to reveal underlying data structures, infer functional modules, and facilitate experiment design. In practice, however, size, uncertainty and complexity of the underlying associations render these applications challenging. In this thesis, we illustrate the use of spectral, combinatorial, and statistical inference techniques in several network science problems.

In Chapters 2-4, we consider network inference challenges. In Chapter 2, we introduce Network Maximal Correlation (NMC) as a multivariate measure of nonlinear association suitable for evaluation on large datasets. We characterize a solution of the NMC optimization using geometric properties of Hilbert spaces for finite discrete and jointly Gaussian random variables. We illustrate an application of NMC and multiple MC in inference of graphical models for bijective, possibly non-monotone, functions of jointly Gaussian variables. As a demonstration of NMC's utility, we infer nonlinear gene association networks and modules in cancer datasets and validate them using survival times of patients. In Chapter 3, we develop a network integration framework to infer gene regulatory networks in human and model organisms fly and worm using diverse and high-throughput datasets. Inferred regulatory interactions have significant overlap with known edges, indicating the robustness and accuracy of the proposed network inference framework. In Chapter 4, we formulate the transitive noise problem in networks as the inverse of matrix transitive closure and introduce an algorithm to solve it efficiently. We demonstrate the effectiveness of our approach in several applications such as regulatory network inference, protein contact map inference and strong collaboration tie inference.

In Chapters 5-8, we consider network analysis challenges. In Chapter 5, we consider the problem of network alignment where the goal is to find a bijective mapping between nodes of two networks to maximize their overlapping edges while minimizing mismatches. This problem is essential in comparative analysis across large datasets and networks. To solve this combinatorial problem, we present a new scalable spec-

tral algorithm which creates an eigenvector relaxation for the underlying optimization. We prove the optimality of the method under certain technical conditions, and show its effectiveness over various synthetic networks as well as in comparative analysis of gene regulatory networks across human, fly and worm species. In Chapter 6, we consider the source inference problem where the goal is to identify the source(s) of propagated signals across biological, social and engineered networks. To solve this problem, we propose a computationally tractable general method based on a path-based network diffusion kernel. We prove mean-field optimality of this method for different scenarios and show its effectiveness over several synthetic networks as well as in identifying sources in a Digg social news network. In Chapter 7, we consider the problem of learning low dimensional structures (such as clusters) in large networks. Here we introduce logistic Random Dot Product Graphs (RDPGs) as a new class of networks which includes most stochastic block models as well as other low dimensional structures. Using this model, we propose a scalable spectral method that solves the maximum likelihood inference problem asymptotically exactly. This leads to a new scalable spectral network clustering algorithm that is robust under different clustering setups. In Chapter 8, we consider the biclustering problem, the analog of clustering on bipartite graphs. This problem has several applications such as inference of co-regulated genes, document classification, and so on. Here we propose an algorithm based on message-passing that closely approximates a general likelihood function and excels at resolving the overlaps between biclusters.

In Chapters 9-12, we consider design challenges of systems and algorithms for engineering networks such as communication networks. In Chapters 9-10, we create a connection between compressive sensing and traditional information theoretic techniques in source, channel and network coding and propose a joint coding scheme over wireless networks based on random projection and restricted eigenvalue principles. Moreover, we characterize fundamental results on the trade-off between the communication rate and the decoding complexity. In Chapters 11-12, we propose an adaptive nonuniform sampling framework, in which time increments between samples are determined as a function of the most recent increments and sample values, obviating the need to track time stamps. We analyze the performance of the proposed method for different stochastic and deterministic signal models and show its effectiveness to enhance measurements of heart ECG signals.

Thesis Supervisor: Muriel Médard

Title: Cecil H. Green Professor of Electrical Engineering and Computer Science

Thesis Supervisor: Manolis Kellis

Title: Professor of Electrical Engineering and Computer Science

Acknowledgments

This thesis would not have been possible without the support of many individuals. First and foremost, I would like to express my special gratitude and thanks to my thesis supervisors Prof. Muriel Médard and Prof. Manolis Kellis. I was very fortunate and lucky to be Muriel's student and a member of her group. Muriel has been an amazing mentor for me, both in academic life and beyond. Her open-mindedness and support provided me the opportunity to work on interdisciplinary research problems. I still remember my first meeting with her when her kindness gave me the confidence and determination to start my research career at MIT. Thank you, Muriel, for all of these years.

I am also indebted to my advisor, Prof. Manolis Kellis. I took Manolis's Computational Biology course at MIT in the beginning of my Ph.D. Having a more theoretical background with an interest for applications, I was really fascinated by the way he was expressing biological problems in the computer science language. That motivated me to pursue interdisciplinary research problems with applications in data sciences and biological sciences. Manolis's excitement towards using theoretical tools and concepts in the area of Computational Biology was one of the main driving forces behind my research. It was a pleasure working with Manolis over the past few years. Thank you so much, Manolis.

I am very grateful to all member of my thesis committee: Prof. Ken R. Duffy, Prof. Ali Jadbabaie, Prof. Brendan Frey and Prof. Robert Nowak. Ken was a great mentor for me during my time in graduate school. He helped me to think about my research problems from a statistical point of view. Our long and enjoyable meetings and his patiently-written emails were always insightful and educational for me. Besides research, Ken was a kind friend. I learned so much from him about different topics from Irish culture to work-out routines. Thank you, Ken.

I was extremely fortunate to have the opportunity to collaborate with Prof. Ali Jadbabaie during my time at MIT. Ali guided me through new topics in optimization and network sciences. His vision had a central role in my research and his insight

can be seen throughout this thesis specially in the network alignment chapter. Ali's rigorous comments and principled ideas have helped me to better understand opportunities and challenges of working in interdisciplinary research areas. Ali, I will always be indebted to you for your fatherly support during my time in the graduate school.

I would also like to express my special gratitude towards Prof. Brendan Frey and Prof. Robert Nowak for being on my thesis committee. I am very thankful for the many enjoyable discussions with them. They have been an inspiration for me during my graduate career. Brendan and Rob, it was an absolute pleasure working with you. I am deeply grateful for your advice and support throughout these years.

I am also grateful to Prof. Alan V. Oppenheim, my academic mentor, for his valuable advice during my time at MIT. His mentorship has led to many interesting and good-spirited discussions relating to different aspects of academic life and beyond. I am also grateful to my colleagues and collaborators Prof. Gerald Quon, Prof Anshul Kundaje, Dr. Daniel Marbach, Prof. Andreas Pfenning, Prof. Jason Ernst, Prof. Sushmita Roy, Prof. Mukul Bansal, and Prof. Daniel Lucani for their friendship and support. Special thanks to my collaborators Prof. Vivek Goyal, Prof. Michelle Effross, Dr. Tarjei S Mikkelsen for their crucial help and wonderful presence.

Over the years, I have had the opportunity to mentor and advise several masters, undergraduates, and high school students, which has contributed immensely to my teaching and mentoring experience. I would like to thank Luke O'Connor, Amy Zhang, Muriel L. Rambeloarison, Aritro Biswas, Eric Mazumdar, Liz Deny, Sebastian Palacios, Orhan Tunc Celiker, Eric Soderstrom, Allison Paul, Michael Colavita and Anna Matilde Tanga. Thank you, all.

Many people have helped and taught me greatly during my time at MIT. I would like to acknowledge all members of both the Network Coding and Reliable Communications and the Computational Biology Groups including Vitaly Abdrashitov, Dave Adams, Ahmad Bairami, Flavio Calmon, Jason Cloud, Jinfeng Du, Ulric Ferner, Kerim Fouli, Bernhard Hupler, Coleen Josephson, Minji Kim, Marie Menshova, Ullrich Monich, Prakash Narayana, Ali ParandehGeibi, Arman Rezaee, Salman Salama-

tian, Shirley Shi, Daniel M Whisman, Weifei Zeng, and specially Michael Lewy from Muriel's group, and Abhishek Sarkar, Wouter Meuleman, Melina Claussnitzer, Jessica Wu, David Hendrix, Chris Bristow, Stefan Washietl, Matt Rasmussen, Matthew Eaton, Irwin Jungreis, Yue Li, Luke Ward, Pouya Kheradpour, Angela Yen, Ben Holmes, Hoon Cho, Rachel Sealfon, Benjamin Iriarte, Luis Barrera, Michal Grzadkowski, Bob Altshuler and Debbie Lehto Goodwin from Manoli's group, among others. In particular, I would like to acknowledge and thank my friends Ali Makhdoumi and Georgios Angelopoulos, with whom I collaborated and co-authored papers on several topics presented in this thesis.

I would like to acknowledge my amazing friends who supported me during this period of my life and made Boston my home. Ali, Azarakhsh, Arman, Donna, Elham, Damoon, Afrooz, Jouya, Sepideh, Omid, Saeed, Sina, Vahid, Afsaneh and Abolhasan, thank you all for being such amazing friends. Last but not least, I would like to thank my family without their unconditional love this journey would have been much more difficult. Words cannot express my gratitude and feeling towards my parents and my brothers whose support and love are the greatest gifts for me. Thank you and love you forever.

Some parts of this thesis have been presented in different journals and conferences or are currently under review:

- **S. Feizi***, A. Makhdoumi*, K. Duffy, M. Kellis, M. Médard, *Network Maximal Correlation*, Presented in RECOMB Conference on Systems Biology and Regulatory Genomic 2015. Available on arXiv: 1606.04789
- **S. Feizi**, G. Quon, M. Recamonde-Mendoza, M. Médard, M. Kellis, Ali Jadbabaie, *Spectral Alignment of Networks*, Presented in Cold Spring Harbor Laboratory Conference, Systems Biology: Networks. Available on arXiv: 1602.04181
- L. O'Connor*, M. Médard and **S. Feizi***, *Clustering over Random Dot Product Graphs*, Available on arXiv:1510.00850
- **S. Feizi**, K. Duffy, M. Kellis, M. Médard, *Network Infusion to Infer Information Sources in Networks*, Presented in RECOMB Conference on Systems Biology

and Regulatory Genomic 2014, and in International Conference on Computational Social Sciences 2015. Available on arXiv:1606.07383

- L. O'Connor* and **S. Feizi***, *Biclustering Using Message Passing*, Advances in Neural Information Processing Systems Foundation (NIPS), 3617-3625, 2014
- **S. Feizi**, G. Angelopoulos, V. K Goyal, M. Médard, *Backward Adaptation for Power Efficient Sampling*, IEEE Transactions on Signal Processing, vol. 62, no. 16, 2014
- **S. Feizi**, D. Marbach , M. Médard, M. Kellis, *Network Deconvolution as a General Method to Distinguish Direct Dependencies in Networks* , Nature Biotechnology, 31, 726-733, 2013
- **S. Feizi**, V. K Goyal, M. Médard, *Time-Stampless Adaptive Nonuniform Sampling for Stochastic Signals*, IEEE Transactions on Signal Processing, vol. 60, no. 10, 2012. Presented in International Conference on Acoustics, Speech, and Signal Processing (ICASSP), Japan, 2012
- **S. Feizi**, M. Médard, *A Power Efficient Sensing/Communication Scheme: Joint Source-Channel-Network Coding by Using Compressive Sensing*, Allerton Conference on Communication, Control, and Computing, 2011
- **S. Feizi**, G. Angelopoulos, V. K Goyal, M. Médard, *Energy-Efficient Time-Stampless Adaptive Nonuniform Sampling*, IEEE sensors, Ireland, 2011
- **S. Feizi**, M. Médard, M. Effros, *Compressive Sensing Over Networks*, Allerton Conference on Communication, Control, and Computing, 2010
- **S. Feizi**, V. K Goyal, M. Médard, *Locally Adaptive Sampling*, Allerton Conference on Communication, Control, and Computing, 2010

Contents

List of Figures	17
List of Tables	35
1 Introduction	37
2 Network Maximal Correlation	45
2.1 Maximal Correlation	51
2.1.1 Maximal Correlation for Finite Categorical and Jointly Gaussian Random Variables	52
2.1.2 Computation of Maximal Correlation	57
2.2 Network Maximal Correlation	58
2.2.1 Network Maximal Correlation: Extensions and Examples	60
2.3 NMC for Finite Discrete Variables	63
2.3.1 Relationship Between NMC and the Maximum Correlation Problem	64
2.3.2 Computation of NMC for Finite Discrete Variables	67
2.3.3 Parallel Computation of NMC	69
2.3.4 A Remark on Computation Challenges of NMC for Continuous Variables	71
2.4 Probabilistic Properties of Network Maximal Correlation for Finite Discrete Random Variables	72
2.4.1 Continuity of Network Maximal Correlation	73

2.4.2	Sample NMC	73
2.5	NMC for Jointly Gaussian Variables	74
2.6	Illustration of NMC's Use	77
2.6.1	Inference of Graphical Models for Functions of Gaussian Variables	77
2.6.2	Performance Assessment of Graphical Model Inference Methods Using Some Examples	80
2.6.3	Inference of Nonlinear Relevance Graphs	83
2.6.4	Illustration of Sample NMC's Use in a Real Data Application	85
2.7	Proofs	89
2.7.1	Proof of Proposition 2.4	89
2.7.2	Proof of Proposition 2.8	90
2.7.3	Proof of Proposition 2.9	91
2.7.4	Proof of Theorem 2.10	92
2.7.5	Proof of Theorem 2.13	92
2.7.6	Proof of Proposition 2.17	95
2.7.7	Proof of Theorem 2.21	96
2.7.8	Proof of Theorem 2.22	96
2.7.9	Proof of Theorem 2.24	101
2.7.10	Proof of Theorem 2.26	103
2.7.11	Proof of Proposition 2.27	106
2.7.12	Proof of Theorem 2.32	106
3	Network Integration	107
3.1	Inference of Feature Specific Regulatory Networks	109
3.2	Inference of Integrated Regulatory Networks	111
4	Network Deconvolution	115
4.1	Network Deconvolution	118
4.1.1	Algorithm	118
4.1.2	Modeling assumptions and intuitions	121
4.1.3	Scaling effects	122

4.1.4	Asymmetric decomposable matrices	124
4.1.5	Computational complexity analysis of network deconvolution	127
4.2	Inferring Gene Regulatory Networks by Network Deconvolution	129
4.2.1	DREAM framework	130
4.2.2	Evaluation Results	133
4.2.3	Effect of Varying ND Parameters	138
4.3	Inferring Protein Structural Contacts	138
4.3.1	Evaluation Results	141
4.4	Application to Co-authorship Collaboration Relationships	146
5	Network Alignment	151
5.1	Network Alignment Problem Setup	156
5.1.1	Problem Formulation and Notation	156
5.1.2	Prior Work	161
5.1.3	Network Alignment and Graph Isomorphism	163
5.2	EigenAlign Algorithm	166
5.2.1	EigenAlign as a Relaxation of Quadratic Assignment	167
5.2.2	Computational Complexity of EigenAlign	171
5.2.3	Performance Characterization of EigenAlign Over Erdős-Rényi Graphs	172
5.3	Low Rank Alignment of Networks	175
5.4	Alignment of Modular Networks	182
5.5	Performance Evaluation Over Synthetic Networks	189
5.5.1	Erdős-Rényi and Power-Law graphs	190
5.5.2	Low Rank Alignment of Networks	195
5.5.3	Modular Network Structures	198
5.6	Performance Evaluation Over Gene Networks	199
5.7	User De-anonymization over twitter follower subgraphs	206
5.8	Proofs	208
5.8.1	Proofs of Lemmas 5.6 and 5.7	208

5.8.2	Proof Of Corollary 5.17	209
5.8.3	Proof of Theorem 5.15	220
5.8.4	Proof Of Theorem 5.18	222
5.8.5	Proof of Theorem 5.25	224
5.8.6	Proof of Theorem 5.26	230
5.8.7	Proof of Theorem 5.27	231
5.8.8	Proof of Theorem 5.28	232
6	Network Infusion	235
6.1	Problem Setup, Computational Difficulties, and Prior Work	240
6.1.1	Source Inference Problem Setup	240
6.1.2	Computational Difficulties of the Source Inference Problem	242
6.1.3	Prior Work	245
6.2	Main Results	248
6.2.1	The path-based Network Diffusion Kernel	248
6.2.2	Global Source Inference	256
6.2.3	Multi-source NI	261
6.2.4	NI for Heterogeneous Network Diffusion	265
6.2.5	Non-parametric NI - Unknown Snapshot Time	268
6.3	NI with Multiple Snapshots	270
6.4	NI Over Synthetic Networks	272
6.5	Inference of News Sources in a Social News Network	280
6.6	Proofs	283
6.6.1	Proof of Proposition 6.8	283
6.6.2	Proof of Proposition 6.11	283
6.6.3	Proof of Proposition 6.16	284
6.6.4	Proof of Theorem 6.17	284
6.6.5	Proof of Theorem 6.20	286
6.6.6	Proof of Theorem 6.26	289
6.6.7	Proof of Proposition 6.27	292

7	Network Clustering	293
7.1	Prior Work	296
7.2	Logistic Random Dot Product Graphs	298
7.2.1	Definitions	298
7.2.2	Relationship between RDPGs and SBMs	301
7.2.3	Maximum-Likelihood Inference of Latent Position Vectors	307
7.3	Performance Evaluation over Synthetic Networks	313
7.4	Proofs	317
8	Network Biclustering	321
8.1	Methods	324
8.1.1	Problem statement	324
8.1.2	BCMP objective function	325
8.1.3	The max-sum algorithm	328
8.1.4	Message updates	330
8.1.5	Biclustering Using Message Passing	331
8.1.6	Learning Model Parameters using Expectation-Maximization	333
8.2	Automatic Parameter Tuning	336
8.2.1	Automatically tuning r_k	336
8.2.2	Automatically tuning δ	336
8.2.3	Automatically detecting the number of clusters	336
8.3	Evaluation Results	337
8.4	Gene Expression Data	340
8.5	Document Classification Example	341
8.6	Proofs	342
8.6.1	Proof of Proposition 1	342
8.6.2	Proof of Proposition 2	343
9	Compressive Sensing over Networks	345
9.1	Technical Background and prior work	347
9.1.1	Compressive Sensing Background	347

9.1.2	Distributed Compression Background	349
9.2	Compressive Sensing and Distributed Source Coding	350
9.2.1	A One-stage Tree with n -Correlated Sources	351
9.2.2	A General Multicast Network with Correlated Sources	355
9.3	Compressive Sensing and Channel Coding	359
10	Joint Source-Channel-Network Coding using Compressive Sensing	367
10.1	Compressive Sensing Background	370
10.2	Main Results	373
10.2.1	Background Results	374
10.2.2	The problem Setup and Notations	374
10.2.3	Joint Source-Channel-Network Coding Scheme	375
10.3	Proofs	381
10.3.1	Proof of Lemma 10.2	381
10.3.2	Proof of Theorem 10.6	382
11	Nonuniform Sampling of Stochastic Signals	385
11.1	TANS Framework	387
11.2	Problem Setup, Signal Models and Background Results	390
11.2.1	Problem Setup	391
11.2.2	Signal Models	391
11.2.3	Generalized Linear Prediction Filter	392
11.3	Greedy TANS	394
11.3.1	Greedy TANS for Autoregressive Signals	395
11.3.2	Greedy TANS for Markovian Signals	396
11.4	Dynamic Programming-Based TANS	403
11.4.1	An Online Source Coding Scheme Based on TANS	404
11.4.2	Approximate Dynamic Programming Methods	406
11.5	Simulation Results	407
11.6	Proofs	411
11.6.1	Proof of Theorem 11.2	411

11.6.2	Proof of Theorem 11.3	412
11.6.3	Proof of Theorem 11.4	413
11.6.4	Proof of Corollary 11.5	414
11.6.5	Proof of Theorem 11.8	415
12	Backward Adaptation for Power Efficient Sampling	417
12.1	Time-Stampless Adaptive Nonuniform Sampling	420
12.2	TANS for Signals with Bounded Derivatives	421
12.2.1	TANS by Polynomial Extrapolation	421
12.2.2	TANS by Incremental Variation	425
12.2.3	Comparison with Synthetic Signals	427
12.3	TANS with Finite Sampling Rates	428
12.4	Performance Evaluation of TANS Schemes on Real Datasets	430
12.5	Energy Analysis of TANS	435
12.5.1	Implementation Considerations of TANS	435
12.5.2	Energy Comparison and Evaluation	437
12.6	Proofs	441
12.6.1	Proof of Theorem 12.1	441
12.6.2	Proof of Theorem 12.2	444
13	Conclusion and Future Directions	445
	Bibliography	457

List of Figures

2-1	Relationships between latent jointly Gaussian variables X_i and observations Y_i for example 1 (panel a), and for example 2 (panel b). . . .	80
2-2	For example 1 (2.52), panel (a) illustrates the covariance matrix of X variables (top left), the covariance matrix of Y variables (top right), the inferred covariance matrix computed by the NMC transformed variables (bottom left), and the inferred covariance matrix computed by the copula method introduced in reference [1]. (b) Covariance estimation error (using an entry-wise l_1 norm metric) for different considered approaches. Using NMC, estimated elements of the covariance matrix are close to true ones.	81
2-3	A similar plot to Figure 2-2 for the case of example 2.	82
2-4	(a) A nonlinear gene module of Glioma cancer (GBMLGG) defined in Section 2.6.4, as a group of genes whose RNA sequence counts are strongly and nonlinearly dependent among each other. (b) Survival time curves for the corresponding nonlinear cancer module. Survival times of cancer patients are significantly associated with average expression of inferred gene module.	87

2-5	An example illustrating the framework considered in Section 2.6.4 for Glioma (GBMLGG). Using normalized RNA sequence counts from the TCGA data portal, we construct both linear and nonlinear complete dependency graphs between genes. We then select top 5% interactions among genes in the NMC network with the largest nonlinear association increases compared to their linear association strengths. Then, using spectral clustering [2] along with a heuristic step to merge small overlapping clusters, we identify dense and large clusters over the thresholded network. In this example, we illustrate the inferred nonlinear gene module over networks with density parameters 1%, 5% and 10%.	88
3-1	(a) The proposed framework to infer integrative regulatory networks. (b) The proposed framework to infer physical feature-specific regulatory networks.	111
3-2	Contributions of input feature-specific networks in integrated interactions.	113
3-3	AUROC and AUPR scores of feature-specific and integrated regulatory networks in human, fly and worm species.	114
4-1	Direct edges in a network (solid blue arrows) can lead to indirect relationships (dashed red arrows) due to transitive information flow. These indirect contributions can be of length two (e.g. $1 \rightarrow 2 \rightarrow 3$), three (e.g. $1 \rightarrow 2 \rightarrow 3 \rightarrow 5$) or higher, and can combine both direct and indirect effects (e.g. $2 \rightarrow 4$ and $2 \rightarrow 3 \rightarrow 4$), and multiple indirect effects along varying paths (e.g. $2 \rightarrow 3 \rightarrow 5$, $2 \rightarrow 4 \rightarrow 5$, $2 \rightarrow 3 \rightarrow 4 \rightarrow 5$). Network deconvolution seeks to reverse the effect of transitive information flow across all indirect paths, in order to recover the true direct network (blue edges) based on the observed network (combined blue and red edges).	118

4-2 Panel (a) shows an example of how eigenvalues of the observed matrix change by application of network deconvolution. In general, ND decreases magnitudes of high positive eigenvalues that were inflated due to transitivity effects over the network. Panel (b) shows the network deconvolution nonlinear mapping function that is applied to eigenvalues of the unscaled observed dependency matrix to compute the eigenvalues of the direct dependency matrix. 121

4-3 This figure illustrates the effect of β in the performance of network deconvolution in various considered applications: Panel (a) demonstrates the overall score of network deconvolution as a gene regulatory network inference method, over three considered networks *E. coli*, *in silico*, and *S.cerevisiae*. This overall score is computed by combining AUPR and AUROC scores as explained in Section 4.2. Panel (b) demonstrates the effect of β in the performance of network deconvolution in inferring protein evolutionary constraints. Y-axis shows the overall contact map prediction rate gains of top 100 predictions over all 15 tested proteins as explained in Section 4.3. Panel (c) illustrates the effect of β in the performance of network deconvolution in inferring weak and string social ties. Y-axis shows the strong tie prediction ratio over the considered co-authorship network as explained in Section 4.4. 125

4-4 Panel (a) shows CPU time usage of performing network deconvolution on networks with 1000 to 50000 nodes and various eigen sparsity levels. As it is illustrated in this figure, network deconvolution can be performed efficiently even on very large networks if the network is low rank. Panel (b) shows the CPU time usage of ND for various network sizes averaged over networks whose eigen sparsity is less than or equal to 100. This figure demonstrates that unlike the baseline network deconvolution algorithm which has a complexity of order $\mathcal{O}(n^3)$, complexity of performing network deconvolution on low rank networks is approximately linear with respect to the number of nodes in the network. (c) Parallelization of network deconvolution is important to apply it effectively on very large networks. This panel demonstrates the total CPU time usage and also CPU time usage per core (per subgraph) for different subgraph sizes of a network with 5000 nodes. . . . 128

4-5 DREAM5 participants were solicited to infer gene regulatory interactions from three different expression compendia. For each network, a list of directed, unsigned edges are predicted, ordered according to confidence scores. To score each of these predicted networks, organism specific gold standards containing the known transcription factor to target gene interactions (true positives) are used as compiled in the DREAM project [3]. These gold standards are not available in the inference part. For the evaluation, all transcription factor-target gene pairs that are not part of the gold standards are considered as negatives.130

4-6 (a) ND applied to the inferred networks of top-scoring methods from DREAM5 leads to consistent improvements for mutual information (MI) and correlation based methods. ND also improves other top-scoring methods, including the best performing method of the DREAM5 challenge (GENIE3). Moreover, the community network obtained by integrating network predictions from individual methods (1-9) before ND is outperformed by the community network based on deconvolved networks. (b) Network motif analysis showing the relative performance of inference methods for cascades (casc.) and feed-forward loops (FFL) before and after ND. Red/blue corresponds to increased/decreased prediction accuracy of the two motif types relative to the overall performance of the method before ND (measured by the area under the ROC curve, AUROC). The original methods (before ND, left side) have different relative performances for cascades and FFLs, e.g., mutual information-based network inference (MI) tends to include feed-forward edges (red arrow), resulting in higher accuracy for FFLs but lower accuracy for cascades, while the opposite is true for the Inferelator and ANOverence. The deconvolved networks (after ND, right side) show significantly higher accuracy for true cascade network motifs for all methods, and moderately improved accuracy for FFLs on average. 134

4-7 A cross validation analysis to choose the eigenvalue scaling parameter β in the application of ND on DREAM5 networks. We train this parameter on 1/3 of input networks and evaluate the performance on the remaining 2/3 of networks. Experiments have been repeated 100 times. 139

- 4-8 Overall total and non-redundant contact discovery rates. (a) Recovery of truly interacting amino acid residues with increasing numbers of top-scoring pairs for mutual information (MI), direct information (DI), and network deconvolution (ND). Recovery after ND is consistently higher than MI for each of tested proteins. ND leads to a small but consistent improvement across the top predictions of direct information (DI). In the strongest-scoring predictions, ranked between 0 and 50 for each method, ND outperforms DI, both independently and as a post-processing step. (b) Discovery of non-redundant contacts, that cannot be explained by previously discovered ones. ND shows improvement over both MI and DI specially in top predictions. (c),(d) We show here an alternate definition of redundant contacts, using a window-based approach (predictions that are within a distance window of previously-discovered contacts are considered redundant). In this metric, ND shows consistent performance increase at all cut-offs. 142
- 4-9 Contact recovery rate after ND is consistently higher than MI for each of 15 tested proteins. ND leads to a small but consistent improvement across the top predictions of direct information (DI). Remarkably, in the strongest-scoring predictions, ranked between 0 and 50 for each method, ND outperforms DI, both independently and as a post-processing step. We excluded protein 1hzx (highlighted in the figure) from overall performance comparisons as all methods had very low performance on it possibly due to its poor alignment. 143
- 4-10 Focusing on non-redundant contacts, that cannot be explained by previously discovered ones, ND shows a significant improvement over the performance of MI. The improvement of ND over DI is even more pronounced when considering non-redundant interactions. 144

4-11	Performance trend of ND is robust to the choice of proximity threshold to define contacts. However, choosing a small threshold leads to sparse contact maps over which significant statistical results cannot be inferred, while considering large thresholds may introduce false contacts. We choose threshold 7 Angstrom and also report the results for both smaller (5 Angstrom, part a) and larger (8 Angstrom, part b) thresholds.	145
4-12	(a) Use of network deconvolution to distinguishing strong ties (red) and weak ties (green) in the largest connected component of a co-authorship network containing 379 authors. True collaboration strengths were computed by summing the number of co-authored papers and down-weighting each paper by the number of additional co-authors. ND only had access to unweighted co-authorship edges, but exploiting transitive relationships to weigh down weak ties resulting in accurate predictions (solid lines), demonstrating that this information lies within the network edges, and that ND is well-suited for discovering it. (b) Beyond the binary classification of strong and weak ties, we found a strong correlation across all 2742 edges connecting 1589 authors, between the weights assigned by ND (y-axis) and the true collaboration strengths (x-axis) obtained using additional publication details.	147
5-1	(a) An illustration of matched, mismatched, and neutral mappings, for undirected graphs. (b) An illustration of inconsistent mappings, for directed graphs, where they are matches in one direction, and mismatches in the other direction.	157
5-2	Framework of EigenAlign algorithm 5.9.	166
5-3	(a) Example modular network structures of Definition 5.23. (b) An illustration of the module-network bipartite graph whose nodes correspond to modules in the original networks.	183

5-4	Eigenvector scores of within and across module mappings of the expected alignment network with parameters $n = 20$, $m = 4$, $p = 0.3$, $q = 0.01$, and (a) $\alpha = 100$, (b) $\alpha = 10$. Parameters of panel (a) satisfy conditions of Theorem 5.25.	185
5-5	Examples of (a) Erdős-Rényi, and (b) power law networks used in Section 5.5.1.	189
5-6	EigenAlign (EA) performance over power law graphs (a) with different alignment scores, (b) with alignment network scaling, and (c) with the greedy assignment method. At each point, simulations have been repeated 20 times.	190
5-7	Performance evaluation of alignment methods over Erdős-Rényi graphs. EigenAlign outperforms IsoRank, NetAlign and Klau optimization in low and high network densities, in all considered noise levels, in both noise models. At each point, simulations have been repeated 10 times.	191
5-8	Performance evaluation of alignment methods over power law graphs. EigenAlign outperforms IsoRank, NetAlign, and Klau optimization in low and high network densities, in all considered noise levels, in both noise models. At each point, simulations have been repeated 10 times.	192
5-9	Running time of network alignment methods over Erdős-Rényi graphs.	194
5-10	Performance evaluation of network alignment methods over (a) Erdős-Rényi ($p = 0.2$), and power law ($\theta = 3$) graphs, with $n = 50$. The computationally expensive SDP-based method outperforms other methods, while EigenAlign has the second best performance. At each point, simulations have been repeated 20 times.	195
5-11	Running time of network alignment methods over Erdős-Rényi graphs with $n = 50$ and $p = 0.2$. The SDP-based method has significantly higher computational complexity compared to other network alignment techniques, which prohibits its use over large graphs. At each point, simulations have been repeated 20 times.	196

5-12	Performance evaluation of different network alignment methods including LowRank Align over Erdős-Rényi graphs with $n = 50$ nodes and edge density of $p = 0.2$. Experiments have been repeated 100 times for all methods except for the SDP-based alignment. For the SDP-based alignment, experiments have been repeated five times owing to its high computational complexity.	197
5-13	Performance evaluation of different network alignment methods including LowRank Align over regular graphs with $n = 50$ nodes and edge density of $p = 0.2$. Experiments have been repeated 100 times for all methods except for the SDP-based alignment. For the SDP-based alignment, experiments have been repeated five times owing to its high computational complexity.	198
5-14	Performance evaluation of network alignment methods over modular network structures with $n = 20$, $m = 2$, $p = 0.2$, and (a) $q = 0.05$, (b) $q = 0.01$. The recovery rate of the EigenAlign+SDP method is close to the one of the SDP-based method, while it has significantly lower computational complexity.	199
5-15	Average running time of network alignment methods over modular network structures with $n = 20$, $m = 2$, $p = 0.2$, and $q = 0.05$. The average running time of the SDP+EigenAlign algorithm is significantly lower than the one of the SDP-based method. Simulations have been repeated 20 times.	200
5-16	(a) A comparative analysis framework across human, fly and worm regulatory networks. (b) An application example of a network alignment method (EigenAlign) over three gene families across human-fly.	201

5-17	The number of matched and mismatched interactions inferred using different network alignment methods across (a) human-fly and (b) human-worm networks. (c) An average match-mismatch ratio for different network alignment methods, in both human-fly and human-worm comparisons. (d) This panel shows average match-mismatch ratios over pairs of biological processes across both human-fly and human-worm networks. EigenAlign mappings preserve regulatory patterns of biological processes across both human-fly and human-worm networks significantly better than the ones of NetAlign and IsoRank methods. .	202
5-18	Examples of conserved biological processes (p -values < 0.05) across human and fly inferred by EigenAlign.	203
5-19	(a) A suitable alignment should map central nodes in one network to central nodes in the other network. Correlation of (b) degree centrality, (c) eigenvector centrality, and (d) page-rank centrality measures of aligned genes using different network alignment methods, in both human-fly and human-worm comparisons.	204
5-20	Performance evaluation of different network alignment methods including LowRank Align over co-regulatory networks. EigenAlign and LowRank Align both lead to similar solutions.	205
5-21	De-anonymization of user IDs over twitter subgraphs. We consider small subgraphs of the twitter follower network sampled in 2008 and 2009 years, with 53 and 88 nodes, and 58 and 91 edges, respectively. Panel (a) illustrates two large components of these networks. Panel (b) illustrates the fraction of correctly de-anonymized users using different network alignment methods. The EigenAlign algorithm is more efficient than both NetAlign and IsoRank methods in user de-anonymization over these graphs.	207

5-22	Expected alignment scores of different mapping combinations. Ellipses with the same color represent corresponding modules across two networks. Nodes with the same color represent true mappings across two networks.	222
6-1	Network Infusion (NI) aims to identify source node(s) by reversing information propagation in the network. NI is based on a path-based network diffusion process that closely approximates the observed diffusion pattern, while leading to a tractable source inference method for large complex networks. The displayed network and infection pattern are parts of the Digg social news network.	237
6-2	Likelihood scores based on the SI diffusion model and the path-based network diffusion kernel for a depicted line network. Although the underlying diffusion model is SI, NI optimization leads to a tight approximation of the exact solution for a wide range of parameter t (except the range between two vertical dashed lines). Unlike SI, NI computation can be efficiently extended to large complex networks, owing to decoupling of its likelihood terms.	243
6-3	(a) A graph considered in Remark 6.5. (b) An example graph with overlapping shortest paths between nodes 0 and 9.	246
6-4	Tightness of the first order approximation of the path-based network diffusion kernel over example networks.	253
6-5	The frequency of having overlapping shortest paths between two randomly selected nodes over an Erdős-Rényi graph with parameter p . In sparse graph ($p \leq \frac{1}{n}$), the overlap frequency is small. Experiments are repeated 20,000 times for each case.	255
6-6	(a) A Markov chain of a hypo-exponential distribution. (b) A Markov chain of a mixed hypo-exponential distribution.	265
6-7	Examples of (a) Erdős-Rényi, and (b) power law networks, with 50 nodes.	273

6-8	Performance of source inference methods over (a,b) Erdős-Rényi, and (c,d) power law networks graphs with 250 nodes and an average density of 1%. Error bars represent \pm standard deviations of empirical results for each case. Experiments have been repeated 100 times so that error margins are small.	274
6-9	(a) An example of a grid network. Performance of source inference methods over grid networks with 250 nodes with (a) one and (b) five independent samples. Error bars represent \pm standard deviations of empirical results for each case. Experiments have been repeated 100 times so that error margins are small.	275
6-10	(a) An asymmetric grid network. Performance of source inference methods over asymmetric grid networks with 250 nodes with (a) one and (b) five independent samples. Node with red color is the source node. Error bars represent \pm standard deviations of empirical results for each case. Experiments have been repeated 100 times so that error margins are small.	276
6-11	(a) An asymmetric grid network. Performance of source inference methods over asymmetric grid networks with 250 nodes with (a) one and (b) five independent samples. Node with red color is the source node. Error bars represent \pm standard deviations of empirical results for each case. Experiments have been repeated 100 times so that error margins are small.	277

6-12 (a) An integrative source inference framework which combines prediction ranks of NI-ME, NI-ML and distance centrality methods. (b) The cumulative number of votes for a particular story over time. (c) Median prediction ranks of true sources for over 3,500 Digg stories inferred by individual and integrative source inference methods. NI-ME outperforms other individual methods with a median prediction rank of 5, while the integrative source inference leads to the best performance with a median prediction rank of 3. (d) Histogram of prediction ranks of sources inferred by different source inference methods. NI-ME and integrative source inference algorithms infer news sources optimally (in their top prediction) for approximately 30% and 38% of stories, respectively. 280

6-13 We consider news spread over the Digg social news networks for over 3,500 news stories. In approximately 65% of cases, nodes that have received the news at a time t do not have any neighbors who had already received the news by that time, violating the basic conditional independency assumption of the SI model. Furthermore, the empirical distribution of remaining news propagation times over network edges cannot be approximated closely by a single distribution of a homogeneous SI diffusion model, even by fitting a general Weibull distribution to the observed data (Section 6.5). 281

6-14 An example of a regular tree used in the proof of Theorem 6.20. . . . 286

6-15 An illustrative figure of the proof of Theorem 6.26. 289

7-1	Normalized mean squared error (one minus squared correlation) of inferred latent positions for two SBMs (a-b) and a non-SBM logistic RDPG (c). The top eigenvectors of the adjacency matrix A and the modularity matrix M do not characterize the community structure in panel (a) and in panel (b), respectively. In panel (c), the top eigenvector of A does not recover the latent structure. In contrast, our method successfully recovers the underlying latent position vectors in all cases.	313
7-2	Performance comparison of our method (logistic RDPG) against spectral clustering methods in different clustering setups. Panels (a)-(e) illustrate networks that can be characterized by SBM, while panel (f) illustrates a non-SBM network model. The scale of the x axis is different in panel (f) than the rest of the panels. Our proposed method performs consistently well, while other methods exhibit sensitive and inconsistent performance in different network clustering setups.	314
7-3	Performance comparison of our method (logistic RDPG) against semidefinite programming-based clustering methods.	315
8-1	Biclustering is the analogue of clustering on a bipartite graph. (a) Biclustering allows nodes to be reordered in a manner that reveals modular structures in the bipartite graph. (b) The rows and columns of an adjacency matrix are similarly biclustered and reordered.	325
8-2	A graphical model of the BCMP objective function.	330
8-3	Performance comparison of the proposed method (BCMP) with ISA and LAS, for Bernoulli and Gaussian models, and for overlapping and non-overlapping biclusters. On the y axis is the total number of misclassified row-column pairs. Either the noise level or the amount of overlap is on the x axis.	337
8-4	Total size and average density of biclusters reported by BCMP, ISA and LAS for three gene expression datasets. a) <i>In Silico</i> . b) <i>E. coli</i> . c) <i>S. cerevesiae</i> .	340

8-5	Running time for BCMP. In Proposition 8.8, we showed that the complexity of a round of message updates for BCMP is nearly linear, i.e., $\mathcal{O}(K(N + \log M)(M + \log N))$. This figure illustrates the running time of BCMP on a personal computer, with $N = M$	342
9-1	A depth one multicast tree with n sources.	349
9-2	A general multicast network with n sources.	355
9-3	A modularized decoder for a) compressive sensing with source coding, (b) compressive sensing with channel coding.	358
9-4	A Point-to-Point Gaussian Channel	360
10-1	The considered problem setup where sources have temporal and spatial dependencies, and receivers desire to reconstruct sources within an allowed distortion level.	369
10-2	A schematic view of our proposed joint source-channel-network coding scheme. Each source projects its high dimensional signal to a lower dimensional space by using a random projection matrix followed by a pre-coding on-off module. Intermediate nodes perform linear network coding. Each receiver uses compressive sensing decoders (LASSO) to reconstruct source signals.	370
10-3	To control the reconstruction error, the loss function should have high curvature around the optimal solution of LASSO (equation 10.1) along the directions of $\mathcal{C}(\mathcal{S}; \alpha)$	372
10-4	An example of a network topology where the network coding matrix can be designed to satisfy the RE condition. Hence, the network can transform sources with spatial dependencies which do not satisfy the RE condition, to a proper lower dimensional space.	378

11-1	A schematic view of the TANS framework: sampling times are determined by a function of m most recently taken samples. Hence, it is not necessary to save sampling times (time stamps) for use in the reconstruction process.	387
11-2	Demonstrating the behavior of different parts of sampling state cost $c(S_i, T_i) = d(S_i, T_i) + a(S_i, T_i)$, where $d(S_i, T_i)$ is the reconstruction distortion and $a(S_i, T_i)$ is the rate penalty function.	390
11-3	A hidden Markov chain considered in Markovian signal model of equation (11.4).	392
11-4	Illustration of Corollary 11.5 by simulation. The red curve is the solution of (11.11). The blue curve is optimal sampling rate T^* . Note that, their difference is always less than 1.	408
11-5	Rate-distortion curves of Theorem 11.8 for a Markovian signal of equation (11.4) with the signal parameters $\alpha_0 = 0.01$ and $\alpha_1 = 0.99$. P_e represents error probability of the estimation process. $P_e = 0$ curve corresponds to the so-called genie-aided scheme.	409
11-6	Average sampling rate versus average reconstruction distortion for a Markovian signal with parameters $\alpha_0 = 0.01$ and $\alpha_1 = 0.99$ for methods: (i) uniform sampling with causal line-connecting (CLC) reconstruction, (ii) uniform sampling with non-causal line-connecting (NCLC) reconstruction, (iii) uniform sampling with generalized linear prediction (GLP) filtering (iv) greedy TANS with generalized linear prediction (GLP) filtering, and (v) analytical lower bound for greedy TANS based on Theorem 11.8.	410

11-7	Average sampling rate versus average reconstruction distortion for a Markovian signal with parameters $\alpha_0 = 0.7$ and $\alpha_1 = 0.97$ for methods: (i) uniform sampling with causal line-connecting (CLC) reconstruction, (ii) uniform sampling with non-causal line-connecting (NCLC) reconstruction, (iii) uniform sampling with generalized linear prediction (GLP) filtering (iv) greedy TANS with generalized linear prediction (GLP) filtering, and (v) analytical lower bound for greedy TANS based on Theorem 11.8.	411
11-8	A rate-distortion plot of a dynamic programming-based TANS for an online source coding application explained in Section 11.4.1.	412
11-9	Comparison of a dynamic programming-based TANS with greedy TANS for a Markovian signal model.	412
12-1	TANS setup	419
12-2	An underlying hidden Markov chain used for the signal creation. . . .	423
12-3	(a) A sample signal and its MC state transitions used in Section (12.2.1). (b) Sampling rate versus time for this signal with TPE method. . . .	426
12-4	A rate-distortion comparison among TPE, TIV and uniform sampling.	427
12-5	An ECG signal modeled as a P-QRS-T complex.	429
12-6	(a) A 6-sec ECG signal from MIT-BIH database, with the DC component removed. (b-h) Reconstructed ECG signals with an effective compression ratio of 6, using different sampling-compression schemes: (b) Uniform Downsampling, (c) Uniform Sampling followed by a DCT-based compression scheme. (d) Level-crossing Nonuniform Sampling. (e) Compressed sensing. (f) TANS by polynomial extrapolation (TPE). (g) TANS by incremental variations (TIV). (h) TANS with finite sampling rates (TFR).	431
12-7	PRD in terms of the transmitted number of samples for each technique.	435
12-8	MSE in terms of transmitted number of samples for each technique. .	436

12-9	High-level block diagram of a system implementing TANS framework. A nonuniform ADC is controlled by a timing control logic, implementing the sampling function depending on the signal and sampling model.	436
12-10(a)	A system using uniform sampling followed by transformation-based compression. (b) A level-crossing nonuniform sampling system. (c) A compressed sensing system.	439

List of Tables

3.1	Number of genes and TFs covered by gene expression data.	113
3.2	Number of TFs covered by evolutionary conserved motifs and TF binding datasets.	114
4.1	Effect of varying both density (α) and eigenvalue scaling (β) parameters on the DREAM5 networks.	138
4.2	Tested protein names and PDB IDs	140
9.1	Comparison of Distributed Compression Methods of Correlated Sources	353
12.1	A energy comparison of 3 different schemes for sampling, compressing and transmitting a 6-sec ECG signal.	441

Chapter 1

Introduction

Network science has been widely adopted in recent years in diverse settings, including molecular and cell biology [4–6], social science [7], communication and engineering networks [8–10]. As datasets become larger, networks provide an efficient system-level approach for understanding dependencies between variables and revealing functionality of their relationships. In data-driven applications, network modeling provides a unifying framework to succinctly represent data, to analyze it, and to design targeted experiments.

For instance, gene regulatory networks [3,4] are used in molecular and cell biology to describe regulatory relationships between transcription factors and their target genes while protein-protein interaction networks [11] capture physical interactions among proteins. Inferring these networks and their system-level functions play crucial roles in elucidating the impact of molecular variations on complex traits and human diseases [12–15]. Beyond biological applications, networks have been used widely in social science to represent social ties and relationships among individuals [16]. Given the recent surge of social networks like Facebook or ResearchGate, recognizing weak and strong ties is increasingly important for recommending friends or colleagues, recognizing conflicts of interest, or evaluating an author’s contribution to a team. Further, several network science tools developed for biological and social networks can be used for engineered networks such as communication networks and vice versa [8].

Although network science has reached to a diverse set of applications, it is plagued by several general challenges. The focus of this thesis is to develop and apply methods for better understanding some of these key challenges. Chapters 2-4 of this thesis consider network inference challenges where we aim to learn dependencies among variables represented by a network topology. In Chapters 5-8 of this thesis, we consider network analysis challenges including problems of network alignment, source inference in networks, network clustering and biclustering. Finally in Chapters 9-12, we consider the design challenges of systems and algorithms for engineering networks using Compressive Sensing, Signal Processing and Information Theory. Below we provide more detail about each of these challenges and our contributions.

In Chapter 2, we introduce Network Maximal Correlation (NMC) as a multivari-

ate measure of nonlinear association among random variables. NMC is defined via an optimization that infers (non-trivial) transformations of variables by maximizing aggregate inner products between transformed variables. We characterize a solution of the NMC optimization using geometric properties of Hilbert spaces for finite discrete and jointly Gaussian random variables. For finite discrete variables, we propose an algorithm based on alternating conditional expectation to determine NMC. We also show that empirically computed NMC converges to NMC exponentially fast in sample size. For jointly Gaussian variables, we show that under some conditions the NMC optimization is an instance of the Max-Cut problem. We then illustrate an application of NMC and multiple MC in inference of graphical model for bijective, possibly non-monotone, functions of jointly Gaussian variables generalizing the copula setup developed by Liu et al. [1, 17]. Finally, we illustrate NMC's utility in a real data application of learning nonlinear dependencies among genes in a cancer dataset.

In Chapter 3, we leverage genome-wide functional genomics datasets to infer gene regulatory networks in human, fly, and worm. We use a rank-based integration approach to combine feature-specific input networks which are constructed using genome-wide transcription factor binding profiles, conserved sequence motif instances and gene expression levels for multiple cell types. We find that inferred regulatory interactions have significant overlap with known interactions in TRANSFAC, REDfly and EdgeDB benchmarks, for human, fly and worm species, indicating the robustness and accuracy of the used network inference pipeline.

Recognizing direct relationships between variables connected in a network is a pervasive problem in biological, social and information sciences, as correlation-based networks contain numerous indirect edges due to higher-order relationships. In Chapter 4, we present a heuristic method for inferring direct effects from an observed correlation matrix containing both direct and indirect effects. We formulate the problem as the inverse of network convolution, and introduce an algorithm that removes the combined effect of all indirect paths of arbitrary length in a closed-form solution by exploiting eigen-decomposition and infinite-series sums. We demonstrate the effectiveness of our approach in several network applications: distinguishing direct targets

in gene expression regulatory networks; recognizing directly-interacting amino-acid residues for protein structure prediction from sequence alignments; and distinguishing strong collaborations in co-authorship social networks using connectivity information alone.

In Chapter 5, we study the network alignment problem which refers to the problem of finding a bijective mapping across vertices of two graphs to maximize the number of overlapping edges and/or to minimize the number of mismatched interactions across networks. This problem arises in many fields such as computational biology, social sciences and computer vision and is often cast as an expensive quadratic assignment problem (QAP). Although spectral methods have received significant attention in different network science problems such as network clustering, the use of spectral techniques in the network alignment problem has been limited partially owing to the lack of principled connections between spectral methods and relaxations of the network alignment optimization. In this chapter, we propose a network alignment framework that uses an orthogonal relaxation of the underlying QAP in a maximum weight bipartite matching optimization. Our method takes into account the ellipsoidal level sets of the quadratic objective function by exploiting eigenvalues and eigenvectors of (transformations of) adjacency graphs. Our framework not only can be employed to provide a theoretical justification for existing heuristic spectral network alignment methods, but it also leads to a new scalable network alignment algorithm which outperforms existing ones over various synthetic and real networks. Moreover, we generalize the objective function of the network alignment problem to consider both matched and mismatched interactions in a standard QAP formulation. This can be critical in applications where networks have low similarity and therefore we expect more mismatches than matches. We assess the effectiveness of our proposed method theoretically for certain classes of networks, through simulations over various synthetic network models, and in two real-data applications; in comparative analysis of gene regulatory networks across human, fly and worm, and in user de-anonymization over twitter follower subgraphs.

Several significant models have been developed that enable the study of diffusion

of signals across biological, social and engineered networks. Within these established frameworks, the inverse problem of identifying the source of the propagated signal is challenging, owing to the numerous alternative possibilities for signal progression through the network. In real world networks, the challenge of determining sources is compounded as the true propagation dynamics are typically unknown, and when they have been directly measured, they rarely conform to the assumptions of any of the well-studied models. In Chapter 6 we introduce a method called Network Infusion (NI) that has been designed to circumvent these issues, making source inference practical for large, complex real world networks. The key idea is that to infer the source node in the network, full characterization of diffusion dynamics, in many cases, may not be necessary. This objective is achieved by creating a diffusion kernel that well-approximates standard diffusion models, but lends itself to inversion, by design, via likelihood maximization or error minimization. We apply NI for both single-source and multi-source diffusion, for both single-snapshot and multi-snapshot observations, and for both homogeneous and heterogeneous diffusion setups. We prove the mean-field optimality of NI for different scenarios, and demonstrate its effectiveness over several synthetic networks. Moreover, we apply NI to a real-data application, identifying news sources in the Digg social network, and demonstrate the effectiveness of NI compared to existing methods. Finally, we propose an integrative source inference framework that combines NI with a distance centrality-based method, which leads to a robust performance in cases where the underlying dynamics are unknown.

Inference of clusters over networks is a central problem in machine learning. Commonly, it is formulated as a discrete optimization, and a continuous relaxation is used to obtain a spectral algorithm. An alternative problem formulation arises by considering a latent space model, in which edge probabilities are determined by continuous latent positions. A model of particular interest is the Random Dot Product Graph (RDPG), which can be fit using an efficient spectral method; however, this method is based on a heuristic that can fail, even in simple cases. In Chapter 7, we consider a closely related latent space model, the Logistic RDPG, which uses a logistic link function to map from latent position inner products to edge likelihoods. Over this

model, we show that asymptotically exact maximum likelihood inference of the latent position vectors can be achieved using a spectral method. Our method involves computing the top eigenvectors of a normalized adjacency matrix and scaling the eigenvectors using a regression step. Through simulations, we show that this method is more accurate and more robust than existing spectral and semidefinite network clustering methods. In particular, the novel regression scaling step is essential to the performance gain of the proposed method.

Biclustering is the analog of clustering on a bipartite graph. Existing methods infer biclusters through local search strategies that find one cluster at a time; a common technique is to update the row memberships based on the current column memberships, and vice versa. In Chapter 8, we propose a biclustering algorithm that maximizes a global objective function using message passing. Our objective function closely approximates a general likelihood function, separating a cluster size penalty term into row- and column-count penalties. Because we use a global optimization framework, our approach excels at resolving the overlaps between biclusters, which are important features of biclusters in practice. Moreover, Expectation-Maximization can be used to learn the model parameters if they are unknown. In simulations, we find that our method outperforms two of the best existing biclustering algorithms, ISA and LAS, when the planted clusters overlap. Applied to three gene expression datasets, our method finds coregulated gene clusters that have high quality in terms of cluster size and density.

In Chapters 9-12, we consider design challenges of systems and algorithms for engineering networks such as communication networks. In Chapter 9, we demonstrate some applications of compressive sensing over networks. We make a connection between compressive sensing and traditional information theoretic techniques in source coding and channel coding. Our results provide an explicit trade-off between the rate and the decoding complexity. Although optimal decoders to recover the original signal, compressed by source coding have high complexity, the compressive sensing decoder is a linear or convex optimization. First, we investigate applications of compressive sensing on distributed compression of correlated sources. Here, by using

compressive sensing, we propose a compression scheme for a family of correlated sources with a modularized decoder, providing a trade-off between the compression rate and the decoding complexity. We call this scheme *Sparse Distributed Compression*. We use this compression scheme for a general multicast network with correlated sources. Here, we first decode some of the sources by a network decoding technique and then, we use a compressive sensing decoder to obtain the whole sources. Then, we investigate applications of compressive sensing on channel coding. We propose a coding scheme that combines compressive sensing and random channel coding for a high-SNR point-to-point Gaussian channel. We call this scheme *Sparse Channel Coding*. We propose a modularized decoder providing a trade-off between the capacity loss and the decoding complexity. At the receiver side, first, we use a compressive sensing decoder on a noisy signal to obtain a noisy estimate of the original signal and then, we apply a traditional channel coding decoder to find the original signal.

In Chapter 10, we propose a joint source-channel-network coding scheme, based on compressive sensing principles, for wireless networks with AWGN channels (that may include multiple access and broadcast), with sources exhibiting temporal and spatial dependencies. Our goal is to provide a reconstruction of sources within an allowed distortion level at each receiver. We perform joint source-channel coding at each source by randomly projecting source values to a lower dimensional space. We consider sources that satisfy the restricted eigenvalue (RE) condition as well as more general sources for which the randomness of the network allows a mapping to lower dimensional spaces. Our approach relies on using analog random linear network coding. The receiver uses compressive sensing decoders to reconstruct sources. Our key insight is the fact that, compressive sensing and analog network coding both preserve the source characteristics required for compressive sensing decoding.

In Chapter 11, we introduce a Time-stampless Adaptive Nonuniform Sampling (TANS) framework, in which time increments between samples are determined by a function of the m most recent increments and sample values. Since only past samples are used in computing time increments, it is not necessary to save sampling times (time stamps) for use in the reconstruction process. We focus on two TANS

schemes for discrete-time stochastic signals: a greedy method, and a method based on dynamic programming. We analyze the performances of these schemes by computing (or bounding) their trade-offs between sampling rate and expected reconstruction distortion for autoregressive and Markovian signals. Simulation results support the analysis of the sampling schemes. We show that, by opportunistically adapting to local signal characteristics, TANS may lead to improved power efficiency in some applications.

Advances in sampling and coding theory have contributed significantly towards lowering power consumption of resource-constrained devices, *e.g.* battery-operated sensor nodes, enabling them to operate for extended periods of time. In Chapter 12, rate and energy efficiency of a proposed adaptive nonuniform sampling framework in Chapter 11 is examined and compared against state-of-the-art methods. TANS addresses one of the main limitations of nonuniform sampling schemes: sampling times do not need to be stored/transmitted since they can be computed using a function of previously taken samples. The sampling rate is adapted continuously with the aim of reducing the rate and therefore the energy consumption of the sampling process when the signal is varying slowly. Three TANS methods are proposed for different signal models and sampling requirements: (i) TANS by polynomial extrapolation, which only assumes the third derivative of the signal is bounded but requires no other specific knowledge of the signal; (ii) TANS by incremental variation, where the sampling time intervals are chosen from a lattice; and (iii) TANS constrained to a finite set of sampling rates. Practical implementation details of TANS are discussed, and its rate and energy performance are compared with uniform sampling followed by a transformation-based compression, nonuniform sampling and compressed sensing. Our results demonstrate that TANS provides significant improvements in terms of both the rate-distortion performance and the energy consumption compared against the other approaches.

In the following chapters, we explain these problems and our contributions with more detail ¹.

¹The notation is chapter-specific and is explained at each chapter of the thesis separately.

Chapter 2

Network Maximal Correlation

Identifying relationships among variables in large datasets is an increasingly important task in systems biology [18], social sciences [19], finance [20], and other fields. For independent observations of bivariate data, several measures exist that characterize the strength of the association based on a linear fit (e.g., Pearson’s correlation [21], canonical correlation [22, 23]), rank statistics (e.g., Spearman’s correlation [24]), and information content (e.g., mutual information [25, 26]). Some of these measures have been extended to the multivariate setting. For instance, Chow and Liu [27] has used mutual information in the inference of tree graphical models, while Liu et al. [1, 17] introduced a copula setup based on rank statistics such as Spearman’s [24] correlation coefficient to characterize graphical models for some nonlinear functions of jointly Gaussian variables. Another method to capture a nonlinear association between two variables is the randomized dependence coefficient [28] where they fix a set of nonlinear functions and then use a randomized rank statistics to compute association between nonlinear transformations of variables.

A classical measure of nonlinear relationships between two random elements X_1 and X_2 is *Maximal Correlation* (MC), introduced by Hirschfeld [29], Gebelein [30], Sarmanov [31] and Rényi [32], having also appeared in the work of Witsenhausen [33], Ahlswede and Gács [34], and Lancaster [35]. MC determines the possibly nonlinear transformations of the two variables, subject to zero mean and unit variance, that maximizes their Pearson’s correlation. MC not only computes an association strength between the variables, but it also characterizes a functional relationship between them.

Definition 2.1 (Maximal Correlation). Let (Ω, \mathcal{F}, P) be a probability triple and for $k \in \{1, 2\}$ let (\mathcal{X}_k, β_k) be a measurable space with $X_k : \Omega \rightarrow \mathcal{X}_k$ a random element. Maximal correlation between the two (not necessarily real-valued) random elements X_1 and X_2 is defined as

$$\rho(X_1, X_2) \triangleq \sup_{\phi_1, \phi_2} \mathbb{E}[\phi_1(X_1) \phi_2(X_2)], \quad (2.1)$$

such that $\phi_k : \mathcal{X}_k \rightarrow \mathbb{R}$ is Borel measurable, $\mathbb{E}[\phi_k(X_k)] = 0$, and $\mathbb{E}[\phi_k(X_k)^2] = 1$, for $k = 1, 2$.

MC can be computed efficiently for both discrete [36] and continuous real-valued [35] random variables (Section 2.1). For discrete random variables, under mild conditions, MC is equal to the second largest singular value of an scaled joint probability distribution matrix and the optimal transformations of the variables can be characterized using right and left singular vectors of the scaled probability distribution matrix. An algorithm to compute MC termed Alternating Conditional Expectation (ACE) is introduced by Breiman and Friedman [36], and further analyzed by Buja [37]. Recently, MC has been used in many applications. It is related to strong data processing inequalities and contraction coefficients, being recently investigated by Anantharam et al. [38], Polyanskiy [39], Raginsky [40]. Further studies include applications in information theory [41, 42], information-theoretic security and privacy [43, 44], data processing [38, 40] and other fields [45–49]. In particular, [45] has introduced multipartite maximal correlation and showed its connection with the multipartite hypercontractivity ribbon. Multipartite maximal correlation aims to find a correlation matrix C with the largest r such that $C - rI$ is positive semidefinite, where I is the identity matrix. This objective function is different than the one of the network maximal correlation (2.3) discussed in the present paper.

Many modern datasets consist of independent observations of high-dimensional multivariate random elements X_1, \dots, X_n (i.e., $n\Gamma^2$). One approach to characterize relationships among these elements would be to determine the bivariate MC for each pair of them. In this approach one would solve the following optimization for all pairs (i, i') , $1 \leq i, i' \leq n$:

$$\sup_{\phi_{i,i'}, \phi_{i',i}} \mathbb{E} [\phi_{i,i'}(X_i) \phi_{i',i}(X_{i'})], \quad (2.2)$$

such that $\phi_{i,i'} : \mathcal{X}_i \rightarrow \mathbb{R}$ is Borel measurable, $\mathbb{E}[\phi_{i,i'}(X_i)] = 0$, and $\mathbb{E}[\phi_{i,i'}(X_i)^2] = 1$, for $1 \leq i, i' \leq n$.

Should they exist, let $\phi_{i,i'}^*$ for $1 \leq i, i' \leq n$ denote the optimal solution of this optimization. By this approach, each element X_i is assigned to $n-1$ transformation functions $\{\phi_{i,i'}^* : 1 \leq i' \leq n, i' \neq i\}$. In some applications, this can cause interpretability

and over-fitting issues. To circumvent these issues, we consider an alternate multivariate extension of the MC optimization (2.1) with two types of constraints. We seek to formulate an optimization that (i) assigns a single transformation function for each element; and (ii) its objective function can consider a subset of variable pairs. The latter is described a graph and is motivated by several distinct considerations: (a) there may be *a priori* known structure that indicates some elements are necessarily unrelated; (b) one may not care about the association of certain elements; or (c) the restriction may serve as a computation reduction technique.

Definition 2.2 (Network Maximal Correlation). Let $G = (V, E)$ be a graph with vertices $V = \{1, 2, \dots, n\}$ and edges $E \subseteq \{(i, i') : i, i' \in V, i \neq i'\}$. The Multivariate Maximal Correlation of X_1, \dots, X_n given G is defined to be

$$\rho_G(X_1, \dots, X_n) \triangleq \sup_{\phi_1, \dots, \phi_n} \sum_{(i, i') \in E} \mathbb{E}[\phi_i(X_i) \phi_{i'}(X_{i'})], \quad (2.3)$$

such that $\phi_i : \mathcal{X}_i \rightarrow \mathbb{R}$ is Borel measurable, $\mathbb{E}[\phi_i(X_i)] = 0$, and $\mathbb{E}[\phi_i(X_i)^2] = 1$, for all $1 \leq i \leq n$.

When no confusion arises, we use ρ_G to refer to the NMC. The NMC optimization maximizes the aggregate inner products between transformed variables. NMC naturally generalizes the bivariate MC to the multivariate setting. For example, the bivariate MC is equivalent to the NMC when the graph $G = (V, E)$ has two nodes connected by an edge. That is, $V = \{1, 2\}$ and $E = \{(1, 2)\}$. As another example, if $V = \{1, 2, 3\}$ and $E = \{(1, 2), (1, 3)\}$, the objective function of the NMC optimization (2.3) aims to find transformation functions $\{\phi_i(X_i)\}_{i=1}^3$ which maximize

$$\mathbb{E}[\phi_1(X_1) \phi_2(X_2) + \phi_1(X_1) \phi_3(X_3)].$$

Note that in this example the transformation function $\phi_1(X_1)$ appears in two terms in the objective function causing additional coupling constraints in optimization (2.3) compared to multiple bivariate MC optimizations (2.2). We investigate this optimization for a general graph $G = (V, E)$. In different applications, we highlight

appropriate selections for G in the NMC optimization ¹.

Since $\mathbb{E}[\phi_i(X_i)^2] = 1$ for any $1 \leq i \leq n$, we have

$$\rho_G = |E| - \inf_{\phi_1, \dots, \phi_n} \frac{1}{2} \sum_{(i, i') \in E} \mathbb{E}[(\phi_i(X_i) - \phi_{i'}(X_{i'}))^2],$$

which means that the NMC optimization (2.3) is equivalent to finding functions of random variables that minimizes the Mean Squared Error (MSE) among all neighboring random variables. This can be useful in different applications such as fitting a nonlinear regression model (e.g., $\phi_i(X_i) = \sum_{i'=\{1, \dots, n\} \setminus \{i\}} \phi_{i'}(X_{i'}) + Z$).

The techniques we have developed to characterize local and global optimums of the aforementioned NMC optimization can be utilized in other related formulations as well. For example, we also introduce absolute NMC, which maximizes the total absolute pairwise correlations among variables (Section 2.2.1). Absolute NMC is appropriate for applications where the strength of an association does not depend on the sign of the correlation coefficient. Moreover, the NMC optimization can be regularized to consider fewer nonlinear transformations or to restrict the set of possible transformations, further avoiding over-fitting issues (Section 2.2.1).

In Section 2.2, for finite discrete random variables, we characterize the solution of the NMC optimization using a natural basis expansion in an appropriate Hilbert space. For finite discrete random variables, we show that the NMC optimization is an instance of the Maximum Correlation Problem (MCP), which is NP-hard [50–53].

In Section 2.3, using the results from the Multivariate Eigenvector Problem (MEP) [50], we characterize necessary conditions that are satisfied at the global optimum of the NMC optimization. We propose an efficient algorithm based on Alternating Conditional Expectation (ACE) [36] that converges to a local optimum of the NMC optimization. We also provide guidelines for choosing appropriate starting points of the ACE algorithm to avoid local maximizers. The proposed algorithm does not require the formation of joint distribution matrices which could be expensive for variables with large number of categories. We develop a distributed version of the ACE

¹We use the terminology *graph* and *network* interchangeably.

algorithm to compute an approximation of the NMC value for large and dense graphs using graph partitioning. We characterize a bound on the expected performance of the proposed algorithm for poly-growth graphs which are often used in analyzing distributed graph algorithms (see e.g., [54–56] and references therein).

In some applications, often only noisy samples from the joint distributions are available. In Section 2.4, we prove a finite sample bound and error guarantees for NMC. In particular, under general conditions we prove that NMC of finite discrete random variables is continuous with respect to their joint probability distributions. That is, a small perturbation in the joint distribution results in a small change in the NMC value. Moreover, we show that the probability of discrepancy between NMC and sample NMC (NMC computed using empirical distributions) decays exponentially fast as the sample size grows.

In Section 2.5, for jointly Gaussian variables, we use Hermite-Chebyshev polynomials as the basis for the functions of Gaussian random variables and characterize the optimal solution of the NMC optimization. Under some conditions, we show that the NMC optimization is equivalent to the Max-Cut problem, which is NP-complete [57]. However, there exist algorithms to approximate its solution using Semidefinite Programming (SDP) within certain approximation factors [58]. Moreover, we provide sufficient conditions under which a solution of the NMC optimization can be characterized exactly.

In Section 2.6, we illustrate some applications of NMC. First we illustrate an application of NMC in inference of graphical model for bijective, possibly non-monotone, functions of jointly Gaussian variables. Graphical models provide a useful framework for characterizing dependencies among variables and for efficient computation of marginal and the mode of distributions [59]. Moreover, Gaussian graphical models play an important role in some applications such as linear regression [60], partial correlation [61], maximum likelihood estimation [62], etc. Here we consider a setup in which observed variables are related to latent jointly Gaussian variables through some link functions. These link functions are unknown, bijective, and can be linear or nonlinear. We prove that under some conditions the inverse covariance matrix of

the transformed variables computed by the NMC optimization (as well as multiple MC optimizations) characterizes the underlying graphical model. Our setup generalizes the *copula* model developed by Liu et al. [1, 17] where the link functions are monotone, all increasing or decreasing. In our setup, link functions are only assumed to be bijective and can be non-monotone. We provide examples which illustrate circumstances where our approach identifies graphical models that the method in [1] could not. This generalization is significant, as there are many practical situations where the relationships between observed and latent variables are non-monotone [63]. Then we use the NMC optimization to characterize a nonlinear global relevance graph and propose a greedy algorithm for its approximate inference. Finally, we provide an example to demonstrate NMC's utility in a real data application where we apply sample NMC to cancer datasets [64] and infer nonlinear gene association networks. Over these inferred networks, we deduce gene modules that are significantly associated with survival times of individuals while these modules are not detected using linear association measures.

2.1 Maximal Correlation

Recall the definition of maximal correlation 2.1. For $k = i, i'$, we let $\phi_k^*(X_k)$ denote an optimal solution of optimization (2.1) should it exist. The existence of a solution $(\phi_i^*(X_i), \phi_{i'}^*(X_{i'}))$ for the MC optimization (2.1) and the uniqueness of such a solution have been investigated in [36]. Maximal correlation $\rho(X_i, X_{i'})$ is always between 0 and 1, where a high MC value indicates a strong association between two variables [30]. Unlike Pearson's linear correlation [21], MC only depends on the joint distribution of the variables and not on their sample spaces \mathcal{X}_i . Several works have investigated different aspects of optimization (2.1) for both discrete [36] and continuous [35] random variables. For discrete variables, under mild conditions, MC is equal to the second largest singular value of the following Q -matrix [36]:

Definition 2.3. Let $P_{X_i, X_{i'}}$ be the joint distribution of discrete variables X_i and $X_{i'}$ with finite alphabet sizes $|\mathcal{X}_i|$ and $|\mathcal{X}_{i'}|$, respectively. Define a matrix $Q_{i, i'} \in \mathbb{R}^{|\mathcal{X}_i| \times |\mathcal{X}_{i'}|}$

whose (j, j') element is

$$Q_{i,i'}(j, j') \triangleq \frac{P_{X_i, X_{i'}}(j, j')}{\sqrt{P_{X_i}(j)P_{X_{i'}}(j')}}. \quad (2.4)$$

The matrix is called the Q -matrix of the distribution $P_{X_i, X_{i'}}$.

In this case, optimal transformations of variables can be characterized using right and left singular vectors of the normalized probability distribution matrix. For Gaussian variables, in [35], Lancaster introduced projections over Hilbert spaces to compute MC. Interestingly in this case the maximal correlation and Pearson's linear correlation is equivalent. Below we extend the Hilbert space projection approach to derive existing MC results for finite discrete variables as this approach will prove useful for the NMC generalization.

2.1.1 Maximal Correlation for Finite Categorical and Jointly Gaussian Random Variables

For $i = 1, 2, \dots, n$, we define a Hilbert space H_i as

$$H_i = \{ \phi_i \circ X_i \mid \phi_i : \mathcal{X}_i \rightarrow \mathbb{R}, \phi_i \text{ is Borel measurable, } \mathbb{E}[\phi_i \circ X_i] < \infty, \mathbb{E}[(\phi_i \circ X_i)^2] < \infty \}, \quad (2.5)$$

where the inner product is defined as $\mathbb{E}[\phi_i(X_i) \phi_i'(X_i)]$, and the norm is defined as $\sqrt{\mathbb{E}[(\phi_i \circ X_i)^2]}$. We use the notation $\phi_i \circ X_i$ and $\phi_i(X_i)$ interchangeably. The assumption that the Hilbert spaces H_i , for $1 \leq i \leq n$, is compact holds when X_i is a finite discrete or a Gaussian random variable amongst other cases [35]. Since every Hilbert space has an orthonormal basis (Theorem 2.4, [65]), we let $\{\psi_{i,j}\}_{j=1}^{\infty}$ represent an orthonormal basis for H_i .

Consider the following optimization:

$$\begin{aligned}
 \sup_{\{a_{i,j}\}_{j=1}^{\infty}, \{a_{i',j'}\}_{j'=1}^{\infty}} & \sum_{j,j'} a_{i,j} a_{i',j'} \rho_{i,i'}^{j,j'} & (2.6) \\
 & \sum_{j=1}^{\infty} a_{k,j}^2 = 1, \quad k = i, i', \\
 & \sum_{j=1}^{\infty} a_{k,j} \mathbb{E}[\psi_{k,j}(X_k)] = 0, \quad k = i, i',
 \end{aligned}$$

where

$$\rho_{i,i'}^{j,j'} \triangleq \mathbb{E}[\psi_{i,j}(X_i) \psi_{i',j'}(X_{i'})]. \quad (2.7)$$

Proposition 2.4. For $k = i, i'$, let $\{a_{k,j}^*\}_{j=1}^{\infty}$ be a solution of optimization (2.6), should it exists. Then

$$\phi_k^*(X_k) = \sum_{j=1}^{\infty} a_{k,j}^* \psi_{k,j}(X_k) \quad (2.8)$$

is a solution ² the optimization (2.1).

Proof. A proof is presented in Section 2.7.1. □

This proposition follows directly from Theorem 2.10. Proposition 2.4 provides an alternative optimization (2.6) to determine maximal correlation. As illustrated in Examples 2.5 and 2.6, selecting appropriate orthonormal bases for Hilbert spaces H_i and $H_{i'}$ is critical to obtaining a tractable optimization problem (2.6). In the following, we use Proposition 2.4 to solve the maximal correlation optimization for finite discrete variables as well as for jointly Gaussian variables.

Example 2.5 (MC for Finite Discrete Random Variables). Maximal correlation between two discrete (finite alphabets) random variables X_i and $X_{i'}$ can be written as the second largest eigen value of the matrix $Q_{i,i'}$ of Definition 2.3 [38, 39, 49]. This

²This means convergence in the L_2 sense, i.e., $\lim_{n \rightarrow \infty} \|\phi_i^*(X_i) - \sum_{j=1}^n a_{i,j}^* \psi_{i,j}(X_i)\| = 0$. Throughout the paper since we only work with the inner product of functions, without loss of generality, we replace $\phi_i^*(X_i)$ by $\sum_{j=1}^{\infty} a_{i,j}^* \psi_{i,j}(X_i)$. Also note that for the case of finite discrete random variables, the convergence in L_2 is in fact pointwise. For further details, see [35, 66].

example explains why this characterization arises naturally from the Hilbert space representation of this quantity. For two discrete random variables X_i and $X_{i'}$ with a joint distribution $P_{X_i, X_{i'}}(\cdot, \cdot)$, let $\{1, \dots, |\mathcal{X}_i|\}$ and $\{1, \dots, |\mathcal{X}_{i'}|\}$ be alphabets of random variables X_i and $X_{i'}$, respectively. We choose the following orthonormal bases for H_i and $H_{i'}$:

$$\psi_{i,j}(X_i) = \mathbf{1}\{X_i = j\} \frac{1}{\sqrt{P_{X_i}(j)}}, \quad \psi_{i',j}(X_{i'}) = \mathbf{1}\{X_{i'} = j\} \frac{1}{\sqrt{P_{X_{i'}}(j)}}.$$

By these selections of bases, we have

$$\rho_{i,i'}^{j,j'} = \mathbb{E}[\psi_{i,j}(X_i) \psi_{i',j'}(X_{i'})] = \frac{P_{X_i, X_{i'}}(j, j')}{\sqrt{P_{X_i}(j)P_{X_{i'}}(j')}}.$$

Moreover, we have

$$\mathbb{E}[\psi_{k,j}(X_k)] = \sqrt{P_{X_k}(j)}, \quad k = i, i'.$$

Thus, optimization (2.6) is simplified to the following optimization:

$$\begin{aligned} \max \quad & \sum_{j,j'} a_{i,j} a_{i',j'} \frac{P_{X_i, X_{i'}}(j, j')}{\sqrt{P_{X_i}(j)P_{X_{i'}}(j')}} \\ & \sum_{j=1}^{|\mathcal{X}_k|} (a_{k,j})^2 = 1, \quad k = i, i', \\ & \sum_{j=1}^{|\mathcal{X}_k|} a_{k,j} \sqrt{P_{X_k}(j)} = 0, \quad k = i, i'. \end{aligned} \tag{2.9}$$

According to Proposition 2.4, to solve MC optimization (2.1) it is sufficient to find the solution to optimization (2.9). In the following, we show that solution of the optimization (2.9) can be computed in a closed form using matrix spectral decomposition.

For $i = 1, \dots, n$, let

$$\mathbf{a}_i \triangleq (a_{i,1}, a_{i,2}, \dots, a_{i,|\mathcal{X}_i|})^T, \quad \sqrt{\mathbf{p}_i} \triangleq \left(\sqrt{P_{X_i}(1)}, \sqrt{P_{X_i}(2)}, \dots, \sqrt{P_{X_i}(|\mathcal{X}_i|)} \right)^T. \quad (2.10)$$

Optimization (2.9) can be re-written as follows:

$$\begin{aligned} \max \quad & \mathbf{a}_i^T Q_{i,i'} \mathbf{a}_{i'} & (2.11) \\ & \|\mathbf{a}_k\|_2 = 1, \quad k = i, i', \\ & \mathbf{a}_k \perp \sqrt{\mathbf{p}_k}, \quad k = i, i', \end{aligned}$$

where $Q_{i,i'}$ is defined in Definition 2.3, and \mathbf{a}_k and $\sqrt{\mathbf{p}_k}$ are defined in (2.26). Since the maximum singular value of $Q_{i,i'}$ is one, with the right and left singular vectors being $\sqrt{\mathbf{p}_i}$ and $\sqrt{\mathbf{p}_{i'}}$, the feasible set of optimization (2.11) includes unit-norm vectors orthogonal to the leading singular vectors of $Q_{i,i'}$. Thus, the optimal value of optimization (2.11) is equal to the second largest singular value of $Q_{i,i'}$ and optimal vectors \mathbf{a}_i^* and $\mathbf{a}_{i'}^*$ are left and right singular vectors corresponding to the second largest singular value of $Q_{i,i'}$.

Example 2.6 (MC for Jointly Gaussian Random Variables). Using Hermite-Chebyshev polynomials as a basis for Jointly Gaussian random variables in the study of maximal correlation is introduced by Lancaster [35] and further studied by Goodman and Johnson [67]. This example briefly explains this approach. In Section 2.5, we use a similar approach to characterize network maximal correlation for more general multivariate jointly Gaussian variables.

Let $(X_i, X_{i'})$ be jointly Gaussian variables with the correlation coefficient $\rho_{i,i'}$. The k -th Hermite-Chebyshev polynomial is defined as

$$\Psi_k(x) \triangleq (-1)^k e^{x^2} \frac{d^k}{dx^k} e^{-x^2}. \quad (2.12)$$

These polynomials form an orthonormal basis with respect to Gaussian distributions

[35]. That is,

$$\int_{-\infty}^{\infty} \Psi_j(x_i) \Psi_{j'}(x_{i'}) p(x_i, x_{i'}) dx_i dx_{i'} = (\rho_{i,i'})^j \mathbf{1}\{j = j'\}, \quad (2.13)$$

where $p(x_i, x_{i'})$ is the joint density function of Gaussian variables X_i and $X_{i'}$ with correlation $\rho_{i,i'}$. Let $\psi_{k,j}$ to be the j -th Hermite-Chebyshev polynomial, for $k = i, i'$.

Thus,

$$\mathbb{E}[\psi_j(X_k)] = \mathbf{1}\{j = 0\}, \quad k = i, i'. \quad (2.14)$$

Therefore optimization (2.6) can be written as

$$\begin{aligned} \sup \quad & \sum_{j=0}^{\infty} a_{i,j} a_{i',j} (\rho_{i,i'})^j \\ & \sum_{j=0}^{\infty} (a_{k,j})^2 = 1, \quad k = i, i', \\ & a_{k,0} = 0, \quad k = i, i'. \end{aligned} \quad (2.15)$$

Since $|\rho_{i,i'}| \leq 1$, an optimal solution of optimization (2.15) is obtained when $|a_{i,1}| = 1$, $|a_{i',1}| = 1$, while other coefficients are equal to zero. The signs of $a_{k,1}$ for $k = i, i'$ are determined so that $a_{i,1} a_{i',1} \rho_{i,i'} = |\rho_{i,i'}|$. This leads to the maximal correlation $|\rho_{i,i'}|$ between two variables that is equal to the absolute value of the correlation coefficient between them when the two random variables are jointly Gaussian. Moreover, optimal transformation functions are either identity or negative identity, i.e.,

$$\phi_k(X_k) = a_{k,1} \psi_1(X_k) = \pm X_k, \quad k = i, i',$$

where signs of variables are selected so that $a_{i,1} a_{i',1} \rho_{i,i'} = |\rho_{i,i'}|$.

Algorithm 1 Alternating Conditional Expectation (ACE)

Initialization: $\phi_i^{(0)}(X_i), \phi_{i'}^{(0)}(X_{i'})$ with mean zero.
for $k=\{0,1,\dots\}$...
 $\phi_i^{(k+1)}(X_i) = \mathbb{E}[\phi_{i'}^{(k)}(X_{i'})|X_i]$.
update: $\phi_i^{(k+1)}(X_i) = \frac{\phi_i^{(k+1)}(X_i)}{\sqrt{\mathbb{E}[(\phi_i^{(k+1)}(X_i))^2]}}$
 $\phi_{i'}^{(k+1)}(X_{i'}) = \mathbb{E}[\phi_i^{(k)}(X_i)|X_{i'}]$.
update: $\phi_{i'}^{(k+1)}(X_{i'}) = \frac{\phi_{i'}^{(k+1)}(X_{i'})}{\sqrt{\mathbb{E}[(\phi_{i'}^{(k+1)}(X_{i'}))^2]}}$
update: $\rho_G^{(k+1)} = \mathbb{E}[\phi_i^{(k+1)}(X_i)\phi_{i'}^{(k+1)}(X_{i'})]$
end

2.1.2 Computation of Maximal Correlation

For discrete random variables, a solution of the optimization (2.1) can be characterized by the second largest singular value of the normalized joint distribution matrix (Section 2.1.1). Given samples of variables (i.e., $\{x_i^{(s)}, x_{i'}^{(s)}\}_{s=1}^m$), one can compute MC using the empirical joint distribution of variables $P_m(\cdot, \cdot)$. Robustness of the MC computation using empirical distributions is discussed in Section 2.4.

An iterative approach to compute MC without forming the joint distribution function is based on *Alternating Conditional Expectation (ACE)* [36]. Briefly, at iteration k , the ACE algorithm computes optimal transformation functions $\{\phi_i^{(k)}(X_i)\}$ using conditional expectations in a Gauss-Seidel manner [68]. If the correlation value at iteration k denoted by $\rho_G^{(k)}$ does not increase by a certain value, the algorithm terminates. We describe the steps of this algorithm in Algorithm 1.

Proposition 2.7 ([36]). The sequence $\rho_G^{(k)}$ generated by Algorithm 1 converges to a local optimum of the optimization (2.1). Moreover if the starting points of Algorithm 1 are such that vectors

$$\left(\phi_i^{(0)}(1)\sqrt{P_{X_i}(1)}, \dots, \phi_i^{(0)}(|\mathcal{X}_i|)\sqrt{P_{X_i}(|\mathcal{X}_i|)}\right)$$

and

$$\left(\phi_{i'}^{(0)}(1)\sqrt{P_{X_{i'}}(1)}, \dots, \phi_{i'}^{(0)}(|\mathcal{X}_{i'}|)\sqrt{P_{X_{i'}}(|\mathcal{X}_{i'}|)}\right)$$

are not orthogonal to the span of the left and right singular vectors corresponding to the second largest singular value of $Q_{i,i'}$, then the ACE algorithm 1 converges to the global optimum. In addition, if the $Q_{i,i'}$ matrix has unique singular vectors (left and right) corresponding to the second largest singular value, optimal transformation functions are unique maximizers of the optimization (2.1).

2.2 Network Maximal Correlation

Recall the definition of the Network Maximal Correlation (Definition 2.2). NMC infers nonlinear transformation functions assigned to each node variable so that the aggregate pairwise correlations over the graph $G = (V, E)$ is maximized ³.

Let $\phi_i^*(\cdot)$ be a solution of the NMC optimization (2.3) should it exist. Then, an edge maximal correlation between variables i and i' is defined as

$$\rho_G^{(i,i')}(X_1, \dots, X_n) \triangleq |\mathbb{E}[\phi_i^*(X_i) \phi_{i'}^*(X_{i'})]|,$$

where $(i, i') \in E$. Unlike the bivariate MC formulation of (2.1), the edge maximal correlation is a function of joint distribution of all variables. Thus, an edge maximal correlation between variables X_i and $X_{i'}$ is always smaller than or equal to their bivariate maximal correlation, i.e.,

$$\rho_G^{(i,i')}(X_1, \dots, X_n) \leq \rho(X_i, X_{i'}).$$

We next provide a framework to study NMC and its properties. The following proposition shows the existence of optimal transformations of the NMC optimization (2.3):

Proposition 2.8. Under the assumption that all the Hilbert spaces $\{H_i : 1 \leq i \leq n\}$ are compact, there exist functions ϕ_i^* such that $\mathbb{E}[\phi_i^*(X_i)] = 0$ and $\mathbb{E}[\phi_i^*(X_i)^2] = 1$ for $1 \leq i \leq n$, that achieve the optimal value of optimization (2.3).

Proof. A proof is presented in Section 2.7.2. □

³Some other extensions of the NMC optimization have been investigated in Section 2.2.1.

Define \mathcal{P}_i as a conditional expectation operation that maps elements of the space $H_{i'}$ (for any $i' \neq i$) onto elements of H_i , $1 \leq i \leq n$, i.e.,

$$\mathcal{P}_i : H_{i'} \rightarrow H_i, \quad \mathcal{P}_i(\phi_{i'}(X_{i'})) = \mathbb{E}[\phi_{i'}(X_{i'}) | X_i] \in H_i.$$

The following proposition characterizes $\{\phi_i^*\}_{i=1}^n$ using the \mathcal{P}_i operators:

Proposition 2.9. There exist optimal transformation functions of NMC optimization

(2.3) $\{\phi_i^*, 1 \leq i \leq n\}$ satisfying

$$\phi_i^*(X_i) = \frac{\sum_{i' \in \mathcal{N}(i)} \mathcal{P}_i(\phi_{i'}^*(X_{i'}))(X_i)}{\text{var}\left(\sum_{i' \in \mathcal{N}(i)} \mathcal{P}_i(\phi_{i'}^*(X_{i'}))(X_i)\right)} \quad (2.16)$$

where $\mathcal{N}(i)$ represents neighbors of node i in the graph $G = (V, E)$, and $\text{var}(\cdot)$ is the variance of a random variable.

Proof. A proof is presented in Section 2.7.3. □

Proposition 2.9 characterizes a property of optimal transformations of NMC (2.3) using conditional expectation operations without explicitly computing them. In the following, we use orthonormal representations of the Hilbert spaces H_i to produce a constructive approach to solve the NMC optimization. Recall that $\{\psi_{i,j}\}_{j=1}^\infty$ represents an orthonormal basis for H_i .

Theorem 2.10. Consider the following optimization:

$$\begin{aligned} \sup_{\{\{a_{i,j}\}_{i=1}^n\}_{j=1}^\infty} & \sum_{(i,i') \in E} \sum_{j,j'} a_{i,j} a_{i',j'} \rho_{i,i'}^{j,j'} \\ & \sum_{j=1}^\infty a_{i,j}^2 = 1, \quad 1 \leq i \leq n, \\ & \sum_{j=1}^\infty a_{i,j} \mathbb{E}[\psi_{i,j}(X_i)] = 0, \quad 1 \leq i \leq n, \end{aligned} \quad (2.17)$$

where

$$\rho_{i,i'}^{j,j'} \triangleq \mathbb{E}[\psi_{i,j}(X_i) \psi_{i',j'}(X_{i'})]. \quad (2.18)$$

Let $\phi_i^*(\cdot)$ and $\{a_{i,j}^*\}_{j=1}^\infty$ for $1 \leq i \leq n$ be solutions of optimizations (2.3) and (2.17), respectively. Then,

$$\phi_i^*(X_i) = \sum_{j=1}^\infty a_{i,j}^* \psi_{i,j}(X_i).$$

is a solution of the NMC optimization (2.3). Moreover, $\rho_{i,i'}^{j,j'}$ for $j, j' \geq 1$ are coefficients of the basis expansion of $P_{X_i, X_{i'}}$ with respect to the basis $\{\psi_{i,j} \psi_{i',j'}\}_{j,j'}$, i.e.,

$$P_{X_i, X_{i'}}(x_i, x_{i'}) = \sum_{j,j'} \rho_{i,i'}^{j,j'} \psi_{i,j}(x_i) \psi_{i',j'}(x_{i'}).$$

Proof. A proof is presented in Section 2.7.4. □

Selecting appropriate Hilbert spaces H_i is critical to have a tractable optimization problem (2.17). In the following, we consider the NMC optimization for finite discrete variables, while the Gaussian case is discussed in Section 2.5.

2.2.1 Network Maximal Correlation: Extensions and Examples

Absolute NMC

The optimization (2.3) maximizes aggregate pairwise correlations over the network. In some applications, the strength of an association does not depend on the sign of the correlation coefficient. In those cases, one can re-write the NMC optimization (2.3) to maximize the total absolute pairwise correlations over the graph as follows:

Definition 2.11 (Absolute Network Maximal Correlation). Consider the following

optimization:

$$\rho_G^A(X_1, \dots, X_n) \triangleq \sup_{\phi_1, \dots, \phi_n} \sum_{(i, i') \in E} |\mathbb{E}[\phi_i(X_i) \phi_{i'}(X_{i'})]|, \quad (2.19)$$

such that $\phi_i : \mathcal{X}_i \rightarrow \mathbb{R}$ is Borel measurable, $\mathbb{E}[\phi_i(X_i)] = 0$, and $\mathbb{E}[\phi_i(X_i)^2] = 1$, for all $1 \leq i \leq n$. $G = (V, E)$ is a graph with vertices $V = \{1, 2, \dots, n\}$ and edges $E = \{(i, i') : i, i' \in V, i \neq i'\}$. We refer to this optimization as an absolute NMC optimization.

A similar approach can be used to characterize the absolute NMC optimization (2.19) by introducing extra variables $s_{i, i'}$ to represent correlation signs of edges (i, i') :

$$\begin{aligned} \sup_{\phi_1, \dots, \phi_n} \quad & \sum_{(i, i') \in E} s_{i, i'} \mathbb{E}[\phi_i(X_i) \phi_{i'}(X_{i'})] \\ & \mathbb{E}[\phi_i(X_i)] = 0, \quad \mathbb{E}[\phi_i^2(X_i)] = 1, \quad 1 \leq i \leq n \\ & s_{i, i'} \in \{-1, 1\}, \quad 1 \leq i, i' \leq n. \end{aligned} \quad (2.20)$$

In this case, similar to Proposition (2.9), there exists an optimal solution $\{\phi_i^*(X_i)\}_{i=1}^n$ of optimization (2.19) such that

$$\phi_i^*(X_i) = \frac{\sum_{i' \in \mathcal{N}(i)} s_{i, i'}^* \mathcal{P}_i(\phi_{i'}^*(X_{i'}))}{\text{var} \left(\sum_{i' \in \mathcal{N}(i)} s_{i, i'}^* \mathcal{P}_i(\phi_{i'}^*(X_{i'})) \right)},$$

where

$$s_{i, i'}^* = \begin{cases} 1 & \text{if } \mathbb{E}[\phi_i^*(X_i) \phi_{i'}^*(X_{i'})] \geq 0 \\ -1 & \text{if } \mathbb{E}[\phi_i^*(X_i) \phi_{i'}^*(X_{i'})] < 0 \end{cases}$$

Regularized NMC

The NMC optimization (2.3) results in n possibly nonlinear transformation functions $\phi_i^*(X_i)$ whose correlation with the original variables can be small. In some

applications, one may wish to restrict the set of possible transformations of the NMC optimization to control the correlation between transformed and original variables. This can be done by introducing a regularization term in the definition of NMC as follows.

Definition 2.12 (Regularized NMC). The regularized NMC of real-valued X_1, \dots, X_n connected by a graph $G = (V, E)$ is defined as the solution of the following optimization:

$$\rho_G^R(X_1, \dots, X_n) = \sup_{\phi_1, \dots, \phi_n} (1 - \lambda) \sum_{(i, i') \in E} \mathbb{E} [\phi_i(X_i) \phi_{i'}(X_{i'})] + \lambda \sum_{i \in V} \mathbb{E} [\phi_i(X_i) (X_i - \mathbb{E}[X_i])], \quad (2.21)$$

such that $\phi_i : \mathcal{X}_i \rightarrow \mathbb{R}$ is Borel measurable, $\mathbb{E}[\phi_i(X_i)] = 0$, and $\mathbb{E}[\phi_i(X_i)^2] = 1$, for all $1 \leq i \leq n$. $G = (V, E)$ is a graph with vertices $V = \{1, 2, \dots, n\}$ and edges $E = \{(i, i') : i, i' \in V, i \neq i'\}$. $0 \leq \lambda \leq 1$ is the regularization parameter.

Unlike MC and NMC, which only depend on the joint distributions of variables, the regularized NMC depends on both joint distributions and support of variables because of the regularization term. Moreover, one can define regularized absolute NMC similarly to the optimization (2.19).

Let the optimal transformation functions computed by optimization (2.21) be $\phi_{i,\lambda}^*(X_i)$. If $\lambda = 0$, $\phi_{i,\lambda}^*(X_i) = \phi_i^*(X_i)$, while if $\lambda = 1$, $\phi_{i,\lambda}^*(X_i) = X_i$. By varying λ between 0 and 1, transformation functions $\phi_{i,\lambda}^*(X_i)$ vary from $\phi_i^*(X_i)$ to X_i . Define

$$R_{G,\lambda}(X_1, \dots, X_n) \triangleq \sum_{(i, i') \in E} \mathbb{E} [\phi_{i,\lambda}^*(X_i) \phi_{i',\lambda}^*(X_{i'})].$$

Therefore, $R_{G,0} = \rho_G$ while $R_{G,1}$ is the total linear correlations over the network. By the definition of NMC, $R_{G,0} \geq R_{G,1}$. Note that we can use a similar algorithm to Algorithm 3 to compute regularized NMC of Definition 2.12. The objective function

of the regularized NMC optimization (2.21) can be written as follows:

$$\sum_{i \in V} \mathbb{E} \left[\phi_i(X_i) \left(\frac{1-\lambda}{2} \sum_{j \in \mathcal{N}(i)} \phi_{i'}(X_{i'}) + \lambda(X_i - \mathbb{E}[X_i]) \right) \right], \quad (2.22)$$

where $\mathcal{N}(i)$ represents the set of neighbors of node i in the graph $G = (V, E)$. To compute the regularized NMC, one can use an algorithm similar to the ACE Algorithm 3 with the following updates for transformation functions:

$$\phi_i^*(X_i) = \mathbb{E} \left[\frac{1-\lambda}{2} \sum_{j \in \mathcal{N}(i)} \phi_{i'}(X_{i'}) + \lambda(X_i - \mathbb{E}[X_i]) \mid X_i \right]. \quad (2.23)$$

2.3 NMC for Finite Discrete Variables

In this section, we analyze the NMC optimization (2.3) for finite discrete random variables, and then introduce an efficient algorithm to compute NMC. We then propose a parallelizable version of the NMC algorithm based on network partitioning and characterize a bound for its expected performance. Computation of absolute NMC, regularized NMC, and NMC for continuous variables are explained in Sections 2.2.1, 2.2.1 and 2.3.4, respectively.

First we introduce some notation. For any vector $\mathbf{v} = (v_1, \dots, v_d) \in \mathbb{R}^d$ and $p \geq 1$, we let $\|\mathbf{v}\|_p$ represent the standard p -norm of the vector \mathbf{v} defined as

$$\|\mathbf{v}\|_p = \left(\sum_{i=1}^d v_i^p \right)^{\frac{1}{p}}.$$

For $p = 2$, we drop the subscript, i.e., $\|\mathbf{v}\| = \|\mathbf{v}\|_2$. The infinite norm of a vector is defined as

$$\|\mathbf{v}\|_\infty = \max_{1 \leq i \leq d} v_i.$$

The inner product between two vectors \mathbf{v} and \mathbf{w} is defined as

$$\langle \mathbf{v}, \mathbf{w} \rangle = \sum_{i=1}^d v_i w_i.$$

For a matrix $V \in \mathbb{R}^{d_1} \times \mathbb{R}^{d_1}$, the matrix norm is a vector norm on $\mathbb{R}^{d_1 \times d_1}$. I.e.,

$$\|V\|_p = \sup_{\mathbf{w} \neq 0} \frac{\|V\mathbf{w}\|_p}{\|\mathbf{w}\|_p}.$$

For $p = 2$, we drop the subscript.

2.3.1 Relationship Between NMC and the Maximum Correlation Problem

Let X_i be a discrete random variable with alphabet $\{1, \dots, |\mathcal{X}_i|\}$. Let

$$\psi_{i,j}(x) = \mathbf{1}\{x = j\} \frac{1}{\sqrt{P_{X_i}(j)}} \quad (2.24)$$

be an orthonormal basis for H_i , where $\mathbf{1}$ is the indicator function. We assume that all the elements of the alphabet $x_i \in \mathcal{X}_i$ have positive probabilities as otherwise they can be neglected without loss of generality. We can write

$$\rho_{i,i'}^{j,j'} = \mathbb{E}[\psi_{i,j}(X_i) \psi_{i',j'}(X_{i'})] = \frac{P_{X_i, X_{i'}}(j, j')}{\sqrt{P_{X_i}(j) P_{X_{i'}}(j')}}.$$

Therefore, the optimization (2.17) is simplified to the following:

$$\begin{aligned} \max_{\{a_{i,j}\}_{i=1}^n}_{j=1}^{\infty} \quad & \sum_{(i,i') \in E} \sum_{j,j'} a_{i,j} a_{i',j'} \frac{P_{X_i, X_{i'}}(j, j')}{\sqrt{P_{X_i}(j) P_{X_{i'}}(j')}} \\ & \sum_{j=1}^{|\mathcal{X}_i|} (a_{i,j})^2 = 1, \quad 1 \leq i \leq n, \\ & \sum_{j=1}^{|\mathcal{X}_i|} a_{i,j} \sqrt{P_{X_i}(j)} = 0, \quad 1 \leq i \leq n. \end{aligned} \quad (2.25)$$

For $i = 1, \dots, n$, let

$$\begin{aligned} \mathbf{a}_i &\triangleq (a_{i,1}, a_{i,2}, \dots, a_{i,|\mathcal{X}_i|})^T \\ \sqrt{\mathbf{p}_i} &\triangleq \left(\sqrt{P_{X_i}(1)}, \sqrt{P_{X_i}(2)}, \dots, \sqrt{P_{X_i}(|\mathcal{X}_i|)} \right)^T. \end{aligned} \quad (2.26)$$

Therefore, the optimization (2.25) can be re-written as follows:

$$\begin{aligned} \max_{\{\mathbf{a}_i\}_{i=1}^n} \quad & \sum_{(i,i') \in E} \mathbf{a}_i^T Q_{i,i'} \mathbf{a}_{i'} \\ & \|\mathbf{a}_i\|_2 = 1, \quad 1 \leq i \leq n, \\ & \mathbf{a}_i \perp \sqrt{\mathbf{p}_i}, \quad 1 \leq i \leq n, \end{aligned} \quad (2.27)$$

where $Q_{i,i'}$ is defined in Definition 2.3 and \perp represents orthogonality between two vectors.

The optimization (2.27) is not convex nor concave in general. Below, we show that this optimization is an instance of the Maximum Correlation Problem (MCP) proposed by Hotelling [51,52]. By making this connection, we use established techniques via the Multivariate Eigenvalue Problem (MEP) to solve optimization (2.27).

For each i , since $I_{|\mathcal{X}_i|} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_i}^T$ is positive semidefinite, we take its square root⁴ and write

$$B_i \triangleq \sqrt{I_{|\mathcal{X}_i|} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_i}^T}, \quad (2.28)$$

where $I_{|\mathcal{X}_i|}$ is a $|\mathcal{X}_i| \times |\mathcal{X}_i|$ identity matrix. Let $\mathbf{b}_i = B_i \mathbf{a}_i$. Let $U_i \Sigma_i U_i^T$ be the singular value decomposition of B_i where $U_i^{(j)}$ is the j -th column of U_i and $\sigma_i^{(j)}$ is the j -th singular value of B_i . We will show that only one singular value of B_i is zero which corresponds to the singular vector $\sqrt{\mathbf{p}_i}$. Without loss of generality, suppose $\sigma_i^{(1)} = 0$,

⁴The square root of a symmetric positive semidefinite matrix A is defined as $\sqrt{A} = U \Sigma^{1/2} U^T$ where $A = U \Sigma U^T$.

for all i . Define A_i , a $|\mathcal{X}_i| \times |\mathcal{X}_i|$ matrix, as follows:

$$A_i \triangleq \left(\left[U_i^{(2)}, \dots, U_i^{(|\mathcal{X}_i|)} \right] \text{diag} \left(\frac{1}{\sigma_i^{(2)}}, \dots, \frac{1}{\sigma_i^{(|\mathcal{X}_i|)}} \right) \left[U_i^{(2)}, \dots, U_i^{(|\mathcal{X}_i|)} \right]^T \right). \quad (2.29)$$

As we show in the proof of Theorem 2.13, $\sigma_i^{(j)} \neq 0$, for all $1 \leq i \leq n$, and $j \geq 2$, A_i is well-defined according to (2.29).

Theorem 2.13. The NMC optimization (2.27) can be re-written as follows:

$$\begin{aligned} \max_{\mathbf{b}_1, \dots, \mathbf{b}_n} \quad & \sum_{(i,i') \in E} \mathbf{b}_i^T A_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) A_{i'} \mathbf{b}_{i'} \\ \text{s.t.} \quad & \|\mathbf{b}_i\|_2 = 1 \quad 1 \leq i \leq n. \end{aligned} \quad (2.30)$$

Proof. A proof is presented in Section 2.7.5. □

Let C be a matrix consisting of submatrices $C_{i,i'} \in \mathbb{R}^{|\mathcal{X}_i| \times |\mathcal{X}_{i'}|}$ where if $(i, i') \in E$,

$$C_{i,i'} \triangleq A_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) A_{i'}, \quad (2.31)$$

otherwise $C_{i,i'}$ is an all zero matrix of size $|\mathcal{X}_i| \times |\mathcal{X}_{i'}|$. Note that since the graph $G = (V, E)$ does not have self-loops, $C_{i,i}$ is a zero matrix for $1 \leq i \leq n$.

Let $\mathbf{b} \triangleq (\mathbf{b}_1^T, \dots, \mathbf{b}_n^T)^T \in \mathbb{R}^{|\mathcal{X}|}$, where $|\mathcal{X}| = \sum_{i=1}^n |\mathcal{X}_i|$ and $\mathbf{b}_i \in \mathbb{R}^{|\mathcal{X}_i|}$.

Corollary 2.14. The NMC optimization (2.30) can be written as follows:

$$\begin{aligned} \max_{\mathbf{b}} \quad & \mathbf{b}^T C \mathbf{b} \\ & \|\mathbf{b}_i\|_2 = 1, \quad 1 \leq i \leq n. \end{aligned} \quad (2.32)$$

The optimization (2.32) is in the standard form of the MCP problem, proposed by Hotelling [51, 52], to find the linear combination of one set of variables that correlates maximally with the linear combination of another set of variables.

Definition 2.15 (Multivariate Eigenvalue Problem [50]). The first-order optimality

condition for optimization (2.32) is the existence of real-valued scalars, namely, Lagrange multipliers $\lambda_1, \dots, \lambda_n$, such that the following system of equations is satisfied:

$$\begin{aligned} \sum_{i'=1}^n C_{i,i'} \mathbf{b}_{i'} &= \lambda_i \mathbf{b}_i, \quad 1 \leq i \leq n \\ \|\mathbf{b}_i\|_2 &= 1, \quad 1 \leq i \leq n. \end{aligned} \tag{2.33}$$

This system of equations is called Multivariate Eigenvalue Problem (MEP).

2.3.2 Computation of NMC for Finite Discrete Variables

In Section 2.3.1, we established a connection among the NMC optimization (2.3), the Maximum Correlation Problem (MCP) and the Multivariate Eigenvalue Problem (MEP) (see e.g., [50–53]). After showing that the NMC optimization can be reformulated as the MCP, we use existing techniques from the literature to solve it. Several local maximizers exist for cases where finding a global optimum of optimization (2.32) is computationally difficult [50, 69]. For example, an aggregated power method that iterates on blocks of C was proposed by Horst [70] as a general technique for solving the MEP (Definition 2.15) numerically.

Below, we summarize general algorithmic ideas to solve MCP:

- (1) First, an efficient algorithm is used to solve MEP which is the necessary first order condition for MCP. This step is studied in references [50, 70].
- (2) Next, a strategy is used to choose starting points of the algorithm or to jump out of the local minima of optimization (2.32). We adopt this step from [53, 71].

An efficient algorithm to solve MEP: Algorithm 2 proposed by [50] is a Gauss-Seidel algorithm [68] to solve MEP. This algorithm is a variant of the classical power iteration method (see e.g. [72]). Let

Algorithm 2 Gauss-Seidel Algorithm for MEP

Input: $C \in \mathbb{R}^{|\mathcal{X}|} \times \mathbb{R}^{|\mathcal{X}|}$.
Initialization: $\mathbf{b}^{(0)} \in \mathbb{R}^{|\mathcal{X}|}$.
for $k = 0, 1, \dots$
 for $i = 1, \dots, n$
 $\tilde{\mathbf{b}}_i^{(k)} = \sum_{i'=1}^{i-1} C_{i,i'} \mathbf{b}_{i'}^{(k+1)} + \sum_{i'=i+1}^n C_{i,i'} \mathbf{b}_{i'}^{(k)}$.
 $\mathbf{b}_i^{(k+1)} = \frac{\tilde{\mathbf{b}}_i^{(k)}}{\|\tilde{\mathbf{b}}_i^{(k)}\|_2}$
 end
end

$$r(\mathbf{b}) = \mathbf{b}^T C \mathbf{b} = \sum_{i,i'=1}^n \beta_i^T C_{i,i'} \beta_{i'} \quad (2.34)$$

$$\lambda_i(\mathbf{b}) = \mathbf{b}_i^T [C_{i,1}, \dots, C_{i,n}] \mathbf{b} = \sum_{i'=1}^n \beta_i^T C_{i,i'} \beta_{i'} \quad (2.35)$$

Algorithm 2 is an iterative algorithm. Let $\beta^{(k)} = (\beta_1^{(k)}, \dots, \beta_n^{(k)})$ be the update of Algorithm 2 at iteration k .

Theorem 2.16 (Theorem 5.1 [53]). The sequence $\{r(\mathbf{b}^{(k)})\}$ generated by Algorithm 2 is monotonically increasing and convergent.

Theorem 2.16 indicates that the Algorithm 2 finds a local optimum of the optimization (2.32). According to Proposition 2.14, this solution provides a local optimum for the NMC optimization (2.3). In the following we introduce a strategy that uses a local optimum of the optimization (2.32) and constructs another solution for the optimization (2.32) with strictly higher objective function value.

A strategy for avoiding local optima of MCP:

Proposition 2.17. Let $\bar{\beta}$ be a solution of the MEP (2.33). Suppose that there exists an $1 \leq i \leq n$ such that $\lambda_i(\bar{\beta}) < 0$. Let $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_n)$ be defined as: $\hat{\beta}_i = -\bar{\beta}_i$, for any i such that $\lambda_i(\bar{\beta}) < 0$, and $\hat{\beta}_{i'} = \bar{\beta}_{i'}$, otherwise. Then, we have $r(\hat{\beta}) > r(\bar{\beta})$.

Proof. A proof is presented in Section 2.7.6. □

Algorithm 2 finds a local optimum $\bar{\beta}$ for the optimization (2.32). This solution can

Algorithm 3 Network ACE to compute NMC

Input: G, X_1, \dots, X_n ,
Initialization: $\phi_1^{(0)}(X_1), \dots, \phi_n^{(0)}(X_n)$ with mean zero and unit variance.
for $k = 0, 1, \dots$
 for $i = 1, \dots, n$
 $\tilde{\phi}_i^{(k)}(X_i) = \mathbb{E} \left[\sum_{i'=1}^{i-1} \phi_{i'}^{(k+1)}(X_{i'}) | X_i \right] + \mathbb{E} \left[\sum_{i'=i+1}^n \phi_{i'}^{(k)}(X_{i'}) | X_i \right]$, for $i' \in \mathcal{N}(i)$.
 update: $\phi_i^{(k+1)}(X_i) = \frac{\tilde{\phi}_i^{(k)}(X_i)}{\sqrt{\mathbb{E} \left[\tilde{\phi}_i^{(k)}(X_i)^2 \right]}}$
 end
 $\rho_G^{(k+1)} = \sum_{(i,i') \in E} \mathbb{E} \left[\phi_i^{(k+1)}(X_i) \phi_{i'}^{(k+1)}(X_{i'}) \right]$
end

be translated to a solution of the NMC optimization (2.3) according to Theorem 2.10, equations (2.24), (2.28), and using $\mathbf{b}_i = B_i \mathbf{a}_i$. This leads to a direct algorithm to find a solution for the NMC optimization (2.3) based on alternating conditional expectation (Algorithm 3). Similarly to Algorithm 2, Algorithm 3 converges to the local optimum of the NMC optimization (2.3). We then use a strategy similarly to the one of Proposition 2.17 to jump out of local maximizers. At each iteration of the algorithm, we update transformation functions as follows: Suppose at iteration k , transformation functions are $\{\phi_i^{(k)}\}_{i=1}^n$. If we fix all variables except the transformation function of node i , an optimal solution of $\phi_i^{(k+1)}$ can be written as the normalized conditional expectation of functions of its neighbors. In each update, the objective function of the NMC optimization increases or remains the same. Note that for the bivariate case (i.e., $n = 2$), Algorithm 3 is simplified to the ACE algorithm [36] for the MC computation (Section 2.1). If the number of nodes n is large, then computation of NMC may be expensive. In the following, we propose an approach to compute NMC using parallel computation.

2.3.3 Parallel Computation of NMC

For large and dense networks, exact computation of NMC may become computationally expensive. For those cases, we propose a parallelizable algorithm which approximates NMC using network partitioning. The idea can be described as follow.

For a given graph $G = (V, E)$,

- (1) Partition the graph into small disjoint sets.
- (2) Estimate NMC for each partition independently.
- (3) Combine NMC solutions over sub-graphs to form an approximation of NMC for the original graph.

Definition 2.18 ([73]). An (ϵ, k) -partitioning of graph $G = (V, E)$ is a distribution on finite partitions of V so that for any partition $\{V_1, \dots, V_M\}$ with non zero probability, $|V_m| \geq k$, for all $1 \leq m \leq M$, where M is the size of the partition. Moreover, the probability that an edge falls across partitions is bounded by ϵ : for any $e \in E$, $\mathbb{P}[e \in E^c] \leq \epsilon$, where $E^c = E \setminus \bigcup_{m=1}^M (V_m \times V_m)$ is the set of cut edges. The probability is with respect to the distribution on partitions.

Definition 2.19. A graph G is poly-growth if there exists $r > 0$ and $C > 0$, such that for any vertex v in the graph,

$$|N_v(d)| \leq Cd^r,$$

where $N_v(d)$ is the number of nodes within distance d of v in G (distance is defined as the shortest path length on the graph).

Reference [73] describes the following procedure for generating an (ϵ, k) -partitioning on a graph:

1. Given $G = (V, E)$, k , and $\epsilon > 0$, we define the truncated geometric distribution as follows:

$$\mathbb{P}[x = l] = \begin{cases} \epsilon(1 - \epsilon)^{l-1}, & l < k, \\ (1 - \epsilon)^{k-1}, & l = k. \end{cases} \quad (2.36)$$

2. We then order nodes arbitrarily $1, \dots, N$. For node i , we sample R_i according to distribution (2.36) and assign all nodes within that distance from node i to color i . If a node has already colored, we re-color it with color i .

3. All nodes with the same color form a partition.

Proposition 2.20 ([73]). If G is a poly-growth graph, then by selecting $k = \Theta(\frac{r}{\epsilon} \log \frac{r}{\epsilon})$, the above procedure results in an (ϵ, Ck^r) partition.

Next, we use an (ϵ, k) -graph partitioning to approximate NMC over large graphs using parallel computations. Consider the following approach:

- (1) Given an (ϵ, k) - partitioning of G , we sample a partition $\{V_1, \dots, V_M\}$ of V .
- (2) For each partition $1 \leq m \leq M$, we compute NMC restricted to $G_m = (V_m, E \cap (V_m \times V_m))$, denoted by $\hat{\rho}_{G_m}$.
- (3) Let $\hat{\rho}_G = \sum_{m=1}^M \hat{\rho}_{G_m}$ be an approximation of ρ_G .

In the following, we bound the approximation error by bounding boundary effects:

Theorem 2.21. Consider an (ϵ, k) - partitioning of the graph G . We have,

$$\mathbb{E}[\hat{\rho}_G] \geq (1 - \epsilon)\rho_G, \quad (2.37)$$

where the expectation is over (ϵ, k) - partitions of graph G .

Proof. A proof is presented in Section 2.7.7. □

2.3.4 A Remark on Computation Challenges of NMC for Continuous Variables

If variables are continuous and we only observe samples from their joint distributions, the empirical computation of conditional expectations in Algorithm 3 may be challenging owing to the lack of sufficient samples. one approach to compute empirical conditional expectations at point $x_i \in \mathbb{R}$ is to use all samples in its B_i neighborhood. With m samples $\{(x_1^{(s)}, \dots, x_n^{(s)})\}_{s=1}^m$, for any i, i' , this leads to

$$P_{X_i, X_{i'}}^{(m)}(x_i, x_{i'}) = \frac{1}{m} \sum_{s=1}^m \mathbf{1} \left\{ x_i^s \in \left[x_i - \frac{B_i}{2}, x_i + \frac{B_i}{2} \right], x_{i'}^s \in \left[x_{i'} - \frac{B_{i'}}{2}, x_{i'} + \frac{B_{i'}}{2} \right] \right\}.$$

In our ACE implementations to compute NMC for continuous variables we have considered both fixed and variable window sizes (i.e., B_i 's). In the variable window size case, we choose B_i 's so that the number of samples in different bins are the same. The default number of considered bins is 10.

2.4 Probabilistic Properties of Network Maximal Correlation for Finite Discrete Random Variables

In many applications, often only noisy measurements or samples of joint distributions are observed. In this section, we prove a finite sample generalization bound, and error guarantees for Network Maximal Correlation of finite discrete random variables. Specifically, under general conditions, we prove that Network Maximal Correlation is a continuous measure with respect to the joint probability distribution. That is, a small perturbation in the distribution results in a small change in the NMC value. Moreover we prove that the probability of discrepancy between NMC and sample NMC (NMC computed using empirical distributions) greater than any given level decays exponentially fast as the sample size m grows.

Throughout this subsection we only consider finite discrete random variables. Moreover to simplify notation we let P be the matrix representation of probability distribution P_{X_1, \dots, X_n} . We assume that all the elements of the alphabet $x_i \in \mathcal{X}_i$ have positive probabilities as otherwise they can be neglected without loss of generality, and define

$$\delta_{X_i}(P) \triangleq \arg \min_{x_i \in \mathcal{X}_i} P_{X_i}(x_i) > 0, \quad 1 \leq i \leq n. \quad (2.38)$$

The empirical distribution of these variables using m observed samples is defined as $P^{(m)}(x_1, \dots, x_n) = \frac{1}{m} \sum_{s=1}^m \mathbf{1}\{x_1^{(s)} = x_1, \dots, x_n^{(s)} = x_n\}$, where $\{(x_1^{(s)}, \dots, x_n^{(s)})\}_{s=1}^m$ are i.i.d. samples drawn according to the distribution P_{X_1, \dots, X_n} . The vector of observed samples of variable X_i is denoted by $\mathbf{x}_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(m)})$.

2.4.1 Continuity of Network Maximal Correlation

Let P and \tilde{P} be two distributions over alphabets $(\mathcal{X}_1, \dots, \mathcal{X}_n)$. Let $K = \max_{1 \leq i \leq n} |\mathcal{X}_i|$. Thus, $\mathcal{K} \triangleq K^n$ is an upper bound on the alphabet size of the joint distribution. In the following, we show that if the infinity norm distance between P and \tilde{P} is small (i.e., $\|P - \tilde{P}\|_\infty = \max_{x_1 \in \mathcal{X}_1, \dots, x_n \in \mathcal{X}_n} |P(x_1, \dots, x_n) - \tilde{P}(x_1, \dots, x_n)| \leq \gamma$), their corresponding NMC values (denoted by ρ_G and $\tilde{\rho}_G$ respectively) are close to each other.

Theorem 2.22. Network Maximal Correlation is a continuous function of the joint probability distribution P . Let $\|P - \tilde{P}\|_\infty \leq \gamma$, for $\gamma \leq \delta^{3/2} \mathcal{K}^{-1}$. Then, we have

$$|\rho_G - \tilde{\rho}_G| \leq \gamma \mathcal{K} |E| \frac{8}{\delta^2}, \quad (2.39)$$

where $\delta = \min_{1 \leq i \leq n} \left(\min\{\delta_{X_i}(P), \delta_{X_i}(\tilde{P})\} \right)$.

Proof. A proof is presented in Section 2.7.8. □

Corollary 2.23. Let ρ and $\tilde{\rho}$ be bivariate MCs with respect to distributions P and \tilde{P} , respectively. Let $\|P - \tilde{P}\|_\infty \leq \gamma$. For any $\gamma \leq \delta^{3/2} \mathcal{K}^{-1}$, we have

$$|\rho - \tilde{\rho}| \leq \gamma \mathcal{K} \frac{8}{\delta^2}, \quad (2.40)$$

where \mathcal{K} and δ are defined according to Theorem 2.22.

2.4.2 Sample NMC

Let $\{(x_1^{(s)}, \dots, x_n^{(s)})\}_{s=1}^m$ be i.i.d. samples drawn according to a joint distribution P_{X_1, \dots, X_n} . Let $P^{(m)}$ denote the empirical distribution obtained from these samples. Network Maximal Correlation computed using this empirical probability distribution is called *Sample Network Maximal Correlation* and is denoted by $\rho_G^{(m)}(\{(x_1^{(s)}, \dots, x_n^{(s)})\}_{s=1}^m)$. For simplicity and when no confusion arises, we refer to the sample NMC by $\rho_G^{(m)}$. In the following, we show that the probability of discrepancy greater than any given value between $\rho_G^{(m)}$ and ρ_G decays exponentially fast as the sample size m grows.

Theorem 2.24. For any $\eta, \epsilon > 0$, if

$$m \geq 2 \left(\frac{32|E|\mathcal{K}}{\epsilon\delta'^2} \right)^2 \log \left(\frac{8n}{\eta} \right), \quad (2.41)$$

then we have

$$\mathbb{P}[|\rho_G^{(m)} - \rho_G| > \epsilon] \leq \eta, \quad (2.42)$$

where $\delta' = \min_{1 \leq i \leq n} \delta_{X_i}(P)$.

Proof. A proof is presented in Section 2.7.9. □

Note that the number of required samples to learn the joint probability distribution reliably is a function of the alphabet size. This is reflected in the right hand side of the bound provided in Theorem 2.24 through the term \mathcal{K} .

Corollary 2.25. Let ρ and $\rho^{(m)}$ be bivariate MC and sample MC, respectively. For

$$m \geq 2 \left(\frac{32\mathcal{K}}{\epsilon\delta'^2} \right)^2 \log \left(\frac{16}{\eta} \right). \quad (2.43)$$

we have

$$\mathbb{P}[|\rho^{(m)} - \rho| > \epsilon] \leq \eta, \quad (2.44)$$

where \mathcal{K} and δ' are defined according to Theorem 2.24.

2.5 NMC for Jointly Gaussian Variables

Suppose that (X_1, \dots, X_n) are jointly Gaussian variables with zero means and unit variances. Let $\rho_{i,i'}$ be the correlation coefficient between variables X_i and $X_{i'}$. We assume that variables are distinct: $|\rho_{i,i'}| \neq 1$ if $i \neq i'$. Let $G_c = (V_c, E_c)$ be the covariance graph corresponding to these variables where $V_c = \{1, 2, \dots, n\}$, and $(i, i') \in E_c$ iff $\rho_{i,i'} \neq 0$. Let $\psi_{i,j}$ to be the j -th Hermite-Chebyshev polynomial, for

$1 \leq i \leq n$. We have

$$\rho_{i,i'}^{j,j'} = \mathbb{E}[\psi_{i,j}(X_i) \psi_{i',j'}(X_{i'})] = (\rho_{i,i'})^j \mathbf{1}\{j = j'\}.$$

Moreover, using the definition of Hermite-Chebyshev polynomials (2.12), we have

$$\mathbb{E}[\psi_{i,j}(X_i)] = \mathbf{1}\{j = 0\}, \quad 1 \leq i \leq n. \quad (2.45)$$

because all of these functions for $j \geq 1$ have zero means over a Gaussian distribution.

Therefore, optimization (2.17) can be written as

$$\begin{aligned} \max_{\{\{a_{i,j}\}_{i=1}^n\}_{j=1}^\infty} \quad & \sum_{(i,i') \in E} \sum_{j=1}^\infty a_{i,j} a_{i',j} (\rho_{i,i'})^j \\ & \sum_{j=1}^\infty (a_{i,j})^2 = 1, \quad 1 \leq i \leq n. \end{aligned} \quad (2.46)$$

We establish that solving the optimization (2.46) is NP-complete by identifying one instance of this problem that reduces to the max-cut problem, which is NP-complete [57].

Theorem 2.26. Let $s_i \in \{-1, 1\}$ for $1 \leq i \leq n$. Suppose

$$\begin{aligned} \sum_{i' \neq i} (1 - s_i s_{i'}) \rho_{i,i'} &\geq 0, \quad \forall 1 \leq i \leq n, \\ \sum_{i' \neq i} s_i s_{i'} \rho_{i,i'} &\geq \sum_{i' \neq i} (\rho_{i,i'})^2, \quad \forall 1 \leq i \leq n. \end{aligned}$$

Then, $\mathbf{a}_i^* = (0, s_i, 0, \dots, 0)$, for $1 \leq i \leq n$ is a global maximizer of the NMC optimization (2.46) over a complete graph without self-loops.

Proof. A proof is presented in Section 2.7.10. □

Proposition 2.27. Under assumptions of Theorem 2.26, the NMC optimization (2.3)

is simplified to the following max-cut optimization:

$$\begin{aligned} \max_{s_i} \quad & \sum_{i \neq i'} s_i s_{i'} \rho_{i,i'} \\ & s_i \in \{-1, 1\}, \quad 1 \leq i \leq n. \end{aligned} \tag{2.47}$$

Moreover, for all $1 \leq i \leq n$, we have $\phi_i^*(X_i) = s_i^* X_i$, where ϕ_i^* and s_i^* are optimal solutions of optimizations (3) and (2.47), respectively.

Proof. A proof is presented in Section 2.7.11. □

Under conditions of Theorem 2.26, one can compute the strength of the nonlinear relationships among variables by solving multiple bivariate MC optimization (2.1). However, one needs to solve optimization (2.47) to compute the signs of covariance coefficients.

For bivariate jointly Gaussian variables, conditions of Proposition 2.27 are always satisfied. This provides an alternative proof that in the bivariate setup $\phi_i^*(X_i) = \pm X_i$ (as we discussed in Section 2.1). For multivariate jointly Gaussian variables however an optimal NMC solution $\phi_i^*(X_i)$ can be different than $\pm X_i$. Proposition 2.27 provides conditions under which $\phi_i^*(X_i) = \pm X_i$ for the multivariate setup.

In general, the Max-Cut optimization (2.47) is NP-complete [57]. However, there exist algorithms to approximate its solution using SemiDefinite Programming (SDP) [58]. Moreover, for planar graphs the Max-Cut problem is dual to the route inspection problem which may be solved in polynomial time [74].

Assumption 2.28. Let (X_1, \dots, X_n) be jointly Gaussian variables, with zero means and unit variances, such that $\rho_{i,i'} \neq 1$ if $i \neq i'$, and

$$\sum_{i' \neq i} \rho_{i,i'} \geq \sum_{i' \neq i} (\rho_{i,i'})^2, \quad \forall 1 \leq i \leq n, \tag{2.48}$$

where $\rho_{i,i'}$ is the correlation coefficient between X_i and $X_{i'}$.

If all correlation coefficients are non-negative (i.e., $\rho_{i,i'} \geq 0$ for $1 \leq i, i' \leq n$),

Assumption 2.28 is immediately satisfied. However, Assumption 2.28 is more general as some correlation coefficients can be negative as well.

Corollary 2.29. Under Assumption 2.28, $\phi_i^*(X_i) = X_i$ is an optimal solution of the NMC optimization (2.46).

2.6 Illustration of NMC's Use

In this section, we illustrate some applications of NMC. First we show that multiple MC and NMC optimizations can be employed to infer graphical models for functions of jointly Gaussian variables, under some conditions. Next, over synthetic data we compare the performance of the proposed graphical model inference algorithm with the copula method proposed by Liu et al. [1]. Finally, we provide an example to demonstrate NMC's utility in a real data application.

2.6.1 Inference of Graphical Models for Functions of Gaussian Variables

Suppose that (X_1, \dots, X_n) are jointly Gaussian variables with the covariance matrix Λ_X . Without loss of generality, we assume all variables have zero means and unit variances. I.e., $\mathbb{E}[X_i] = 0$ and $\mathbb{E}[X_i^2] = 1$, for all $1 \leq i \leq n$. Let J_X be the information (precision) matrix [59] of these variables where $J_X = \Lambda_X^{-1}$. Define $G_X = (V_X, E_X)$ such that, $(i, j) \in E_X$ if and only if $J_X(i, j) \neq 0$.

Theorem 2.30 ([59]). If $(i, i') \notin E_X$, then

$$X_i \perp\!\!\!\perp X_{i'} | \{X_k, k \neq i, i'\} \quad (2.49)$$

where $\perp\!\!\!\perp$ represents independence between variables.

This Theorem has been stated in other references as well (e.g., [75]). Theorem 2.30 represents a way to model explicitly the joint distribution of Gaussian variables using

a graphical model $G_X = (V_X, E_X)$. This result is critical in several applications involved with Gaussian variables which requires computation of marginal distributions, or computation of the mode of the distribution. These computations can be performed efficiently over the graphical model using belief propagation approaches [59]. Moreover, Gaussian graphical models play an important role in many applications such as linear regression [60], partial correlation [61], maximum likelihood estimation [62], etc. In many applications, even if variables are not jointly Gaussian, a Gaussian approximation is used often, partially owing to the efficient inference of their graphical models.

In the following, under some conditions, we use the multiple MC (2.2) and NMC (2.3) optimizations to characterize graphical models for functions of latent jointly Gaussian variables. These functions are unknown, bijective, and can be linear or nonlinear. More precisely, let $Y_i = f_i(X_i)$, where $f_i : \mathbb{R} \rightarrow \mathbb{R}$ is a bijective function. Our goal is to characterize a graphical model for the variables (Y_1, Y_2, \dots, Y_n) without the knowledge of $f_i(\cdot)$ functions.

Consider the following optimizations:

$$\sup_{g_{i,i'}, g_{i',i}} \mathbb{E}[g_{i,i'}(Y_i) g_{i',i}(Y_{i'})], \quad (2.50)$$

such that $g_{i,i'} : \mathcal{Y}_i \rightarrow \mathbb{R}$ is Borel measurable, $\mathbb{E}[g_{i,i'}(Y_i)] = 0$, and $\mathbb{E}[g_{i,i'}(Y_i)^2] = 1$, for $1 \leq i, i' \leq n$, and

$$\max_{g_i : \mathbb{R} \rightarrow \mathbb{R}} \sum_{(i,i')} \mathbb{E}[g_i(Y_i) g_{i'}(Y_{i'})], \quad (2.51)$$

$$\mathbb{E}[g_i(Y_i)] = 0, \quad 1 \leq i \leq n,$$

$$\mathbb{E}[g_i^2(Y_i)] = 1, \quad 1 \leq i \leq n,$$

such that $g_i : \mathcal{Y}_i \rightarrow \mathbb{R}$ is Borel measurable, $\mathbb{E}[g_i(Y_i)] = 0$, and $\mathbb{E}[g_i(Y_i)^2] = 1$, for $1 \leq i \leq n$. Optimization (2.50) solves multiple MC between all pairs of variables while Optimization (2.51) solves NMC considering a fully connected graph. Suppose

$g_{i,i'}^*(Y_i)$ and $g_i^*(Y_i)$ represent optimal solutions of (2.50) and (2.51), respectively.

Theorem 2.31. $g_{i,i'}^*(Y_i) = \pm X_i$ for all $1 \leq i, i' \leq n$ is an optimal solution of (2.50). Moreover, if Assumption 2.28 holds, $g_i^*(Y_i) = X_i$ for all $1 \leq i \leq n$ is an optimal solution of (2.51).

Proof. A proof is presented in Section 2.7.12. □

In the multiple MC optimization (2.50), each variable Y_i is assigned to $n - 1$ transformation functions $\{g_{i,i'}^*(Y_i) : 1 \leq i' \neq i \leq n\}$. However, when variables X_1, \dots, X_n are jointly Gaussian, all these functions are equal to $\pm X_i$. In general this is not true for non-Gaussian distributions. On the other hand, when the Assumption 2.28 holds, a solution of the NMC optimization (2.51) recovers the sign of latent variables as well.

Define matrices $\Lambda_{MultiMC}$ and Λ_{NMC} such that

$$\begin{aligned}\Lambda_{MultiMC}(i, i') &= \mathbb{E}[g_{i,i'}^*(Y_i)g_{i',i}^*(Y_{i'})] \\ \Lambda_{NMC}(i, i') &= \mathbb{E}[g_i^*(Y_i)g_{i'}^*(Y_{i'})].\end{aligned}$$

Moreover, let $J_{MultiMC} = \Lambda_{MultiMC}^{-1}$ and $J_{NMC} = \Lambda_{NMC}^{-1}$. Define $G_{NMC} = (V_{NMC}, E_{NMC})$ such that, $(i, j) \in E_{NMC}$ iff $J_{NMC}(i, j) \neq 0$. Let \mathcal{G}_{NMC} be the set of all possible G_{NMC} since the optimal solution of the optimization (2.51) may not be unique. Similarly define $G_{MultiMC}$ and $\mathcal{G}_{MultiMC}$.

Corollary 2.32. Let G_Y be a graphical model of variables $Y_i = f_i(X_i)$ where $f_i : \mathbb{R} \rightarrow \mathbb{R}$ is a bijective function, for $1 \leq i \leq n$. Then $G_Y \in \mathcal{G}_{MultiMC}$. Moreover, if X_i for $1 \leq i \leq n$ satisfy Assumption 2.28, $G_Y \in \mathcal{G}_{NMC}$.

Corollary 2.32 characterizes the graphical model of variables Y_i 's that are related to latent jointly Gaussian variables X_i 's through unknown bijective functions f_i 's. The family of distributions considered in this corollary is broad and includes many Gaussian distributions as well as distributions whose variables are bijective functions of Gaussian variables. Graphical models characterized in Corollary 2.32 can be used

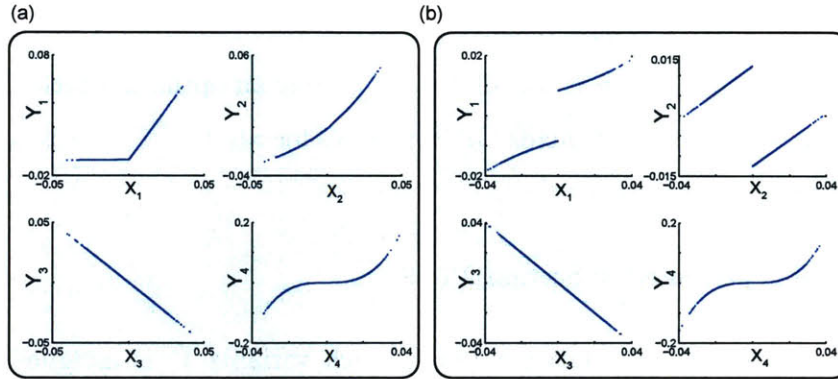


Figure 2-1: Relationships between latent jointly Gaussian variables X_i and observations Y_i for example 1 (panel a), and for example 2 (panel b).

in computation of marginal distributions, computation of the mode of the joint distribution, and in other applications of estimation and prediction similarly to the case of Gaussian graphical models. Note that Corollary 2.32 only considers f_i functions that are bijective.

2.6.2 Performance Assessment of Graphical Model Inference Methods Using Some Examples

In this section, we provide two examples to highlight differences between the setup considered in Section 2.6.1 and the one in reference [1]. Reference [1] considers the graphical model inference problem for functions of jointly Gaussian variables when the underlying link functions are all increasing (or decreasing) functions. In [1] to characterize the underlying graphical model, rank statistics between variables and connections between Spearman's and Pearson's correlation coefficients are employed. However, in the more general setup considered in Section 2.6.1, link functions are bijective and can be non-monotone.

Consider four jointly Gaussian variables X_1, \dots, X_4 , each with unit variance and

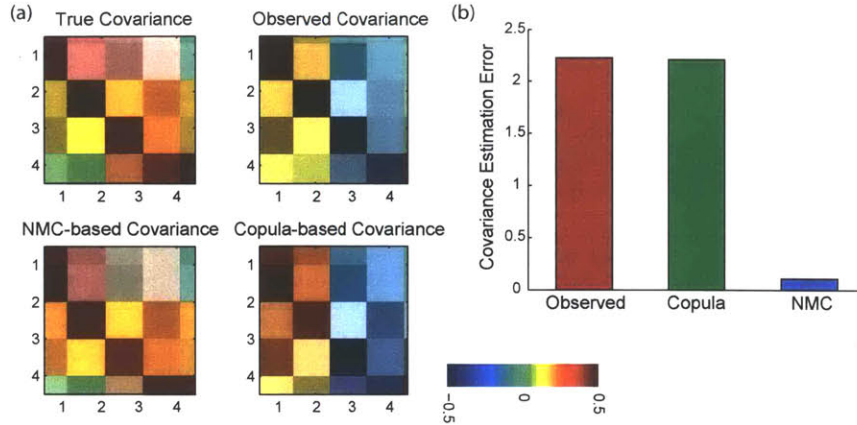


Figure 2-2: For example 1 (2.52), panel (a) illustrates the covariance matrix of X variables (top left), the covariance matrix of Y variables (top right), the inferred covariance matrix computed by the NMC transformed variables (bottom left), and the inferred covariance matrix computed by the copula method introduced in reference [1]. (b) Covariance estimation error (using an entry-wise l_1 norm metric) for different considered approaches. Using NMC, estimated elements of the covariance matrix are close to true ones.

zero mean. We observe samples from variables $Y_i = f_i(X_i)$. In example 1, we have

$$Y_1 = f_1(X_1) = \begin{cases} 10X_1, & \text{if } X_1 \geq 0, \\ \frac{1}{10}X_1, & \text{otherwise,} \end{cases}, \quad Y_2 = f_2(X_2) = e^{20X_2}, \quad (2.52)$$

$$Y_3 = f_3(X_3) = -X_3, \quad Y_4 = f_4(X_4) = X_4^3.$$

In example 2, we consider

$$Y_1 = f_1(X_1) = \begin{cases} e^{20X_1}, & \text{if } X_1 \geq 0, \\ -e^{-20X_1}, & \text{otherwise,} \end{cases} \quad (2.53)$$

$$Y_2 = f_2(X_2) = \begin{cases} \frac{1}{\max(X_2)}X_2 - 1, & \text{if } X_2 \geq 0, \\ \frac{-1}{\min(X_2)}X_2 + 1, & \text{otherwise,} \end{cases}$$

$$Y_3 = f_3(X_3) = -X_3, \quad Y_4 = f_4(X_4) = X_4^3.$$

Relationships between samples of these variables are illustrated in Figure 2-1. We

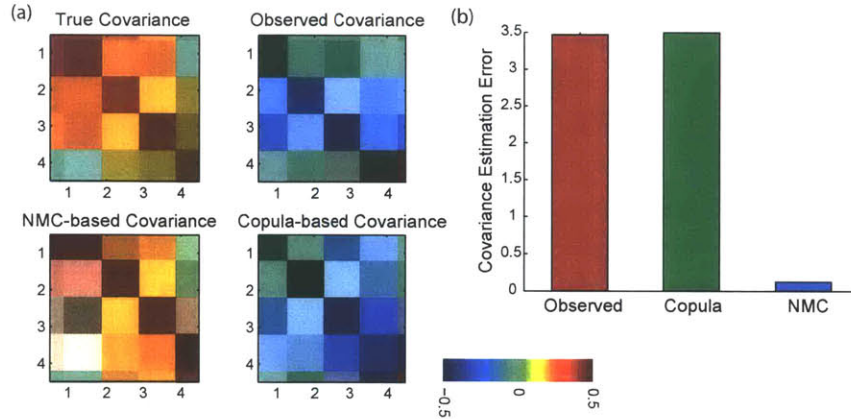


Figure 2-3: A similar plot to Figure 2-2 for the case of example 2.

have drawn 10,000 samples of these random variables. Note that all variables are normalized to have zero means and unit variances. In example 1 (equations (2.52)) all functions are monotone, however some of them are increasing and some are decreasing which does not satisfy the assumptions required in the setup of reference [1]. In example 2 (equations (2.53)) the function $f_2(\cdot)$ is neither increasing nor decreasing which again does not satisfy the assumptions of the setup considered in reference [1].

The functions $f_i(\cdot)$ remain unknown for the inference part. As per Theorem 2.32, we use transformed variables to infer the underlying covariance and precision matrices. As is illustrated in Figures 2-2 and 2-3, the covariance matrices computed using observed variables Y_i show significant errors compared to the true covariance matrices owing to the existence of nonlinear relationships. Moreover, the covariance matrices computed using the copula approach introduced in reference [1] show similarly high errors since the considered setups violate their assumptions. However, the inferred covariance matrices using transformed variables computed by the NMC optimization closely approximate the true covariance coefficients. To compute NMC we used Algorithm 3 with $B_i = 10$ for all i (Section 2.3.4). The Small errors in covariance coefficient estimation in these examples are owing to the computation of the NMC solution using empirical distributions according to the ACE algorithm 3.

2.6.3 Inference of Nonlinear Relevance Graphs

Relevance graphs (RG's) play an important role in many applications including systems biology, social and economic sciences as they characterize variable pairs with highest observed similarities [76]. Let X_1, \dots, X_n be n real-valued random variables with zero means and unit variances. If the similarity measure between variables X_i and X_j denoted by $S(X_i, X_j) \in \mathbb{R}$ is independent of the other variables, the resulting graph with $S(X_i, X_j)$ as the weight on the edge between i and j is called a *pairwise relevance graph* (PRG).

Definition 2.33. Consider the following optimization:

$$G_S^* = \arg \max_{\substack{G=(V,E) \\ |E|=k}} \sum_{(i,j) \in E} S(X_i, X_j). \quad (2.54)$$

Then, $G_S^* = (V, E^*)$ is called a pairwise relevance graph (PRG) of variables X_1, \dots, X_n , with k edges, corresponding to the similarity measure $S(\cdot, \cdot)$. k controls the edge density of the inferred relevance graph.

PRG's can be inferred efficiently in practice. In the following, we highlight two examples of such graphs:

Example 2.34. If $S(\cdot, \cdot)$ is a correlation-based similarity measure (e.g., $S(X_i, X_j) = |\mathbb{E}[X_i X_j]|$), the optimization (2.54) results in a correlation-based PRG. This graph only captures the top k linear associations among variables, ignoring nonlinear ones.

Example 2.35. Mutual information between two discrete variables X_i and X_j is defined as follows [25]:

$$I(X_i; X_j) \triangleq \sum_{x_i \in \mathcal{X}_i} \sum_{x_j \in \mathcal{X}_j} P_{X_i, X_j}(x_i, x_j) \log \left(\frac{P_{X_i, X_j}(x_i, x_j)}{P_{X_i}(x_i) P_{X_j}(x_j)} \right)$$

If $S(\cdot, \cdot)$ is a mutual information-based similarity measure (e.g., $S(X_i, X_j) = I(X_i; X_j)$), optimization (2.54) results in an MI-based PRG.

If the similarity measure depends on all variables, the resulting graph is called a *global relevance graph* (GRG). A global relevance graph can capture system level

properties of observed dependencies among variables. However, inference of a GRG is computationally challenging in general. Below, we introduce an NMC-based global relevance graph that captures observed nonlinear associations among variables and also provides explicit nonlinear transformation functions through which variables are associated.

Definition 2.36. The NMC-based GRG of variables X_1, \dots, X_n , with k edges is defined as

$$G_{NMC}^* = \arg \max_{\substack{G=(V,E) \\ |E|=k}} \rho_G(X_1, \dots, X_n), \quad (2.55)$$

where ρ_G is defined in (2.3).

The optimization (2.55) is combinatorial since it requires computation of NMC over all graphs with k edges. However, unlike mutual information-based GRG's, the required sample size to have a reliable computation of G_{NMC}^* does not increase with the size of the network.

Below, we propose a greedy algorithm to find an approximate solution for optimization (2.55). Suppose at iteration r , the inferred network is $G = (V, E)$ whose nodes are assigned to transformations ϕ_i^r . At this iteration, for each non-connected pair of nodes, we wish to add the corresponding edge to the network, compute NMC, and finally select the edge with the highest NMC increase. However, this is computationally expensive as it requires NMC computation multiple times. Instead, at this iteration, we add an edge to the network with the highest correlation of its node transformations inferred in the last iteration. Then, we update all transformation functions by applying the ACE algorithm 3. We repeat this procedure explained in Algorithm 4 till the inferred network has a certain number of edges. Algorithm 4 starts from a graph with no edge and gradually populates it until the inferred graph has k edges. Alternatively, one can start with a complete graph and remove interactions with lowest nonlinear correlations iteratively. Because NMC is more efficient to determine over sparse graphs, the former has lower computational complexity than

Algorithm 4 Inference of a global relevance graph using NMC

Input: X_1, \dots, X_n
Initialization: $E = \{\}, \mathcal{F}_1^* = \mathcal{F}_1, \dots, \mathcal{F}_n^* = \mathcal{F}_n$
for $r = 1 : k$ **do**
 let $\{i, j, \phi_i^*, \phi_j^*\} = \operatorname{argmax} \mathbb{E}[\phi_{i'}(X_{i'})\phi_{j'}(X_{j'})]$
 subject to: $\{i', j'\} \in G \times G \setminus E, \phi_{i'} \in \mathcal{F}_{i'}^*, \phi_{j'} \in \mathcal{F}_{j'}^*$
 update graph: $E = E \cup \{(i, j)\}$
 update: run Algorithm 3 on G
 For $i \in E$, let ϕ_i^* be the output function
 update function sets: For $i \in E$ let $\mathcal{F}_i^* = \{\phi_i^*\}$
end for
run Algorithm 3 on G to obtain $\rho_{G^*}(X_1, \dots, X_n)$
Output: $E, \rho_{G^*}(X_1, \dots, X_n), \phi_i^*(X_i), \dots, \phi_n^*(X_n)$

the later.

2.6.4 Illustration of Sample NMC's Use in a Real Data Application

Having illustrated applications of the NMC optimization in learning nonlinear dependencies among variables, here we provide an example to demonstrate NMC's utility in a real data application. It is important to note that this real-data example is based on several heuristic steps and assumptions about the underlying data which we explain in Section 2.6.4.

For the data, we use normalized RNA sequence counts from the TCGA data portal for the Glioma cancer (GBMLGG) at the gene level. Other cancer types can be considered in our framework as well. We use the processed data provided in reference [64]. Let the variable X_i denote the RNA sequence counts of the gene i , $1 \leq i \leq n$. We have samples from this variable in m patients denoted by $\{x_i^j\}_{j=1}^m$. Let $\phi_i^*(X_i)$ for $1 \leq i \leq n$ be a solution for the NMC optimization (2.3) over a complete graph without self-loops. Let $A_{NMC} \in \mathbb{R}^{n \times n}$ be a symmetric adjacency matrix where $A_{NMC}[i, i'] = \mathbb{E}[\phi_i^*(X_i)\phi_{i'}^*(X_{i'})]$ for $1 \leq i, i' \leq n$. Moreover, let A_{Lin} be a symmetric linear covariance matrix where $A_{Lin}[i, i'] = \mathbb{E}[X_i X_{i'}]$ for $1 \leq i, i' \leq n$. Figure 2-4-a illustrates top 5% of elements of the $|A_{NMC}| - |A_{Lin}|$ matrix for the

Glioma Cancer (GBMLGG). The non-zero elements of this matrix represent gene pairs whose RNA sequence counts are strongly and nonlinearly associated to each other. Using a heuristic spectral clustering method, we find a subset of gene pairs that are strongly and nonlinearly associated to each other (so-called a nonlinear gene module). An example of such a nonlinear gene module is illustrated by a yellow box in Figure 2-4-a.

Next we partition individuals to two equal-size groups based on their average gene expression ranks in that module and perform a standard survival time analysis for each module based on Kaplan-Meier procedure to estimate survival function [77]. We compute its associated log-rank p -value to determine its association with individual survival times in the considered cancer type [78]. We find that this gene module is significantly associated with survival times of cancer patients (Figure 2-4-b) while it is not detected using linear association measures. Several references [79–83] have hypothesized that complex nonlinear relationships among genes may play important roles in cancer pathways. The proposed NMC algorithm and inferred nonlinear gene modules can be used in discovering such complex nonlinear relationships in different cancer types. To substantiate these inferences, further experiments should be performed to determine the involvement of these nonlinear gene interactions and modules in different cancers, which is beyond the scope of the present paper.

Further Details on Data and Inference Assumptions

Data: Cancer is a complex disease involving abnormal cell growth with the potential to invade or spread to other parts of the body [84]. Different studies have shown associations of micro RNA patterns in different human cancers [63,64]. In this section, we use normalized RNA sequence counts from the TCGA data portal for the Glioma cancer (GBMLGG) at the gene level. Other cancer types can be considered in our framework as well. We use the processed data provided in reference [64]. Let the variable X_i denote the RNA sequence counts of the gene i . We have samples from this variable in m patients denoted by $\{x_i^j\}_{j=1}^m$. In the following we explain inference assumptions and steps we have taken in this application.

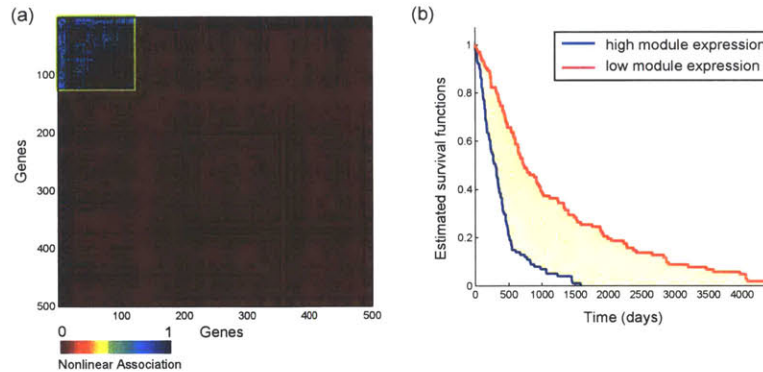


Figure 2-4: (a) A nonlinear gene module of Glioma cancer (GBMLGG) defined in Section 2.6.4, as a group of genes whose RNA sequence counts are strongly and non-linearly dependent among each other. (b) Survival time curves for the corresponding nonlinear cancer module. Survival times of cancer patients are significantly associated with average expression of inferred gene module.

Assumptions: We assume that different RNA sequence count samples for a gene (i.e., $\{x_i^j\}_{j=1}^m$) are independent. That is, there is no between-patient correlations among these samples [85]. Also we assume that conditions of Theorem 2.26 holds. That is, input data comes from, possibly nonlinear, functions of some latent jointly Gaussian variables satisfying conditions of Theorem 2.26. These functions are unknown and bijective. This assumption is less restrictive than the one of the standard covariance analysis where input variables are assumed to be jointly Gaussian.

Inferring Nonlinear Gene-Gene Dependency Networks: For each cancer type, first we select the top 500 highly-variant genes based on their normalized variances (i.e., $n = 500$) [86]. Then, we compute both linear and nonlinear sample associations among genes. In this application, we wish to infer the strength of associations among genes. To compute the nonlinear association network (i.e., the NMC network), according to Section 2.6.1, we compute multiple sample MC's among gene pairs using Algorithm 1. Moreover, we infer a linear association network by computing all pairwise correlations among gene expressions. We then select 5% of interactions among genes in the NMC network with largest nonlinear association increases compared to their linear association strengths. In Figure 2-5, we consider some other density parameters as well (e.g., 1% and 10%).

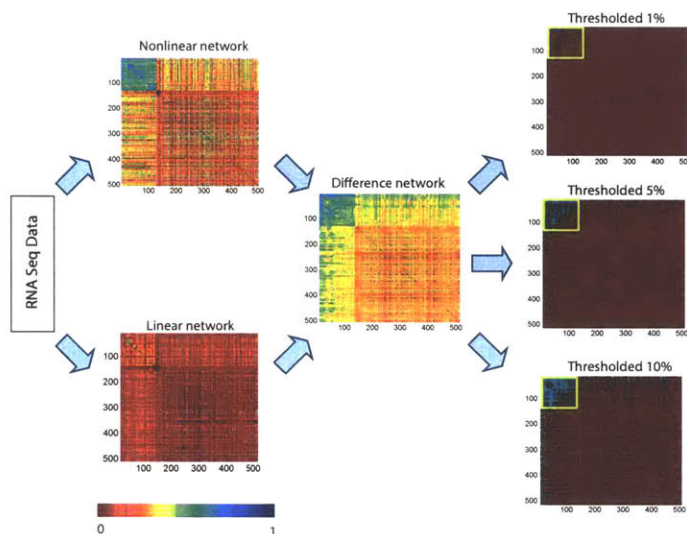


Figure 2-5: An example illustrating the framework considered in Section 2.6.4 for Glioma (GBMLGG). Using normalized RNA sequence counts from the TCGA data portal, we construct both linear and nonlinear complete dependency graphs between genes. We then select top 5% interactions among genes in the NMC network with the largest nonlinear association increases compared to their linear association strengths. Then, using spectral clustering [2] along with a heuristic step to merge small overlapping clusters, we identify dense and large clusters over the thresholded network. In this example, we illustrate the inferred nonlinear gene module over networks with density parameters 1%, 5% and 10%.

Inferring Nonlinear Gene Modules: Next, we partition each network to k groups using a standard spectral clustering algorithm based on the modularity transformation [2]. We consider different values of k (between 1 and 20) to obtain dense and large clusters. Then, in a heuristic manner, we merge small and heavily overlapping modules to form large and dense ones. We define a *gene module* as a group of genes that are densely connected to each other in the network. A gene module is called nonlinear if it is present in the NMC network but not in the linear one. We use a permutation test [87] to compute a p -value for each gene module in the network by permuting the network structure and comparing the density of the module in the original network with the ones in permuted ones. We only consider gene modules with p -values less than 0.05.

Validations using Survival Time Analysis: For each inferred nonlinear mod-

ule, we partition individuals to two equal-size groups based on their average gene expression ranks in that module. We then perform a standard survival time analysis for each module based on Kaplan-Meier procedure to estimate survival function [77], and compute its associated log-rank p -value to determine its association with individual survival times in the considered cancer type [78]. We perform Benjamini and Hochberg multiple hypothesis correction [88] over the computed p -values of different nonlinear modules. Figure 2-4 illustrates an example nonlinear gene module for Glioma Cancer (GBMLGG). This gene module is significantly associated with survival times of cancer patients while it is not detected using linear association measures. To the best of our knowledge, this gene module has not been previously identified.

2.7 Proofs

2.7.1 Proof of Proposition 2.4

Proof. Recall that $\{\psi_{i,j}\}_{j=1}^{\infty}$ is the orthonormal basis of H_i for $1 \leq i \leq n$. We can represent functions ϕ_i and $\phi_{i'}$ in terms of the basis functions, i.e., for any $x_i \in \mathcal{X}_i$, $x_{i'} \in \mathcal{X}_{i'}$,

$$\phi_i(x_i) = \sum_{j=1}^{\infty} a_{i,j} \psi_{i,j}(x_i), \quad \phi_{i'}(x_{i'}) = \sum_{j=1}^{\infty} a_{i',j} \psi_{i',j}(x_{i'}),$$

for two sequences of coefficients $\{a_{i,j}\}_{j=1}^{\infty}$ and $\{a_{i',j}\}_{j=1}^{\infty}$. Thus, the constraint $\mathbb{E}[\phi_k(X_k)^2] = 1$ in optimization (2.1) would be translated into $\sum_{j=1}^{\infty} a_{k,j}^2 = 1$ and the constraint $\mathbb{E}[\phi_k(X_k)] = 0$ is simplified to $\sum_{j=1}^{\infty} a_{k,j} \mathbb{E}[\psi_{k,j}(X_k)] = 0$, for $k = i, i'$. Moreover, we have

$$\mathbb{E}[\phi_i(X_i)\phi_{i'}(X_{i'})] = \sum_{j,j'=1}^{\infty} a_{i,j} a_{i',j'} \mathbb{E}[\psi_{i,j}(X_i)\psi_{i',j'}(X_{i'})].$$

Thus, maximal correlation optimization (2.1) can be re-written as follows:

$$\begin{aligned} \sup_{\{a_{i,j}\}_{j=1}^{\infty}, \{a_{i',j}\}_{j=1}^{\infty}} & \sum_{j,j'} a_{i,j} a_{i',j'} \mathbb{E} [\psi_{i,j}(X_i) \psi_{i',j'}(X_{i'})] \\ & \sum_{j=1}^{\infty} a_{k,j}^2 = 1, \quad k = i, i', \\ & \sum_{j=1}^{\infty} a_{k,j} \mathbb{E} [\psi_{i,j}(X_i)] = 0, \quad k = i, i'. \end{aligned} \quad (2.56)$$

Moreover, $\{\psi_{i,j} \psi_{i',j'}\}_{j,j'}$ forms a basis for L_2 functions of X_i and $X_{i'}$, i.e., $\{\phi_{i,i'} : \mathcal{X}_i \times \mathcal{X}_{i'} \rightarrow \mathbb{R} : \mathbb{E} [\phi_{i,i'}^2(X_i, X_{i'})] < \infty\}$. Since $P_{X_i, X_{i'}}(\cdot, \cdot) \in$ belongs to this set of functions, we can write

$$P_{X_i, X_{i'}}(x_i, x_{i'}) = \sum_{j,j'} \rho_{i,i'}^{j,j'} \psi_{i,j}(x_i) \psi_{i',j'}(x_{i'}). \quad (2.57)$$

Using equation (2.57), we have

$$\mathbb{E} [\psi_{i,j}(X_i) \psi_{i',j'}(X_{i'})] = \rho_{i,i'}^{j,j'}. \quad (2.58)$$

Using (2.58) in optimization (2.56) leads to optimization (2.6). This completes the proof. \square

2.7.2 Proof of Proposition 2.8

Proof. Let

$$F_i = \{\phi_i(X_i) \in H_i : \mathbb{E} [\phi_i(X_i)] = 0, \mathbb{E} [\phi_i(X_i)^2] = 1\}.$$

Now consider the function defined on $F_1 \times \dots \times F_n$, as

$$T(\phi_1, \dots, \phi_n) = \sum_{(i,i') \in E} \mathbb{E} [\phi_i(X_i) \phi_{i'}(X_{i'})].$$

Below, we show that T is continuous and $F_1 \times \cdots \times F_n$ is a compact subset of a finite dimensional space (in the product topology). We then employ Weierstrass Extreme Value Theorem, that says every continuous real-valued function on a compact space attains its extreme values, to conclude the proof. For any $1 \leq i \leq n$, first we show F_i is compact. Since F_i is a subset of compact space H_i , it suffices to show that it is closed. This follows from the definition of F_i and the norm on H_i . In particular, if $\{\phi_i^N(X_i)\}_{N=1}^\infty$ is a sequence in F_i that is converging to $\phi_i(X_i)$, then we have

$$\begin{aligned} 0 &= \lim_{N \rightarrow \infty} \mathbb{E} [(\phi_i^N(X_i) - \phi_i(X_i))^2] = \lim_{N \rightarrow \infty} \mathbb{E} [(\phi_i^N(X_i))^2] + \mathbb{E} [\phi_i(X_i)^2] - 2\mathbb{E} [\phi_i^N(X_i)\phi_i(X_i)] \\ &\geq 1 + \mathbb{E} [\phi_i(X_i)^2] - 2\sqrt{\mathbb{E} [\phi_i(X_i)^2]} = (1 - \sqrt{\mathbb{E} [\phi_i(X_i)^2]})^2, \end{aligned}$$

which leads to $\mathbb{E} [\phi_i(X_i)^2] = 1$, where we used Cauchy-Schwartz inequality along with $\mathbb{E} [(\phi_i^N(X_i))^2] = 1$ for all N . Since $F_1 \times \cdots \times F_n$ is the product of finitely many compact sets, therefore, it is compact as well. Next we show T is continuous. Since T has finitely many arguments, it suffices to show that it is continuous with respect to each argument. Since we have

$$|\mathbb{E} [\phi_i(X_i)\phi_{i'}(X_{i'})] - \mathbb{E} [\tilde{\phi}_i(X_i)\phi_{i'}(X_{i'})]| \leq \sqrt{\mathbb{E} [(\phi_i(X_i) - \tilde{\phi}_i(X_i))^2]} \sqrt{\mathbb{E} [\phi_{i'}(X_{i'})^2]},$$

each term $\mathbb{E} [\phi_i(X_i)\phi_{i'}(X_{i'})]$ is continuous in every $\phi_k(X_k)$ for $k = i, i'$. □

2.7.3 Proof of Proposition 2.9

Proof. Using Proposition 2.8, an optimal solution for the Optimization 2.3 exists which we denote it by $\{\phi_i(X_i)^*\}_{i=1}^n$. Next, using this optimal solution we can construct another set of optimal functions satisfying the condition of Proposition 2.9. For a

given $1 \leq i_0 \leq n$ with neighbors $\mathcal{N}(i_0)$, we have

$$\begin{aligned}
 \sum_{(i,i') \in E} \mathbb{E} [\phi_i(X_i) \phi_{i'}(X_{i'})] &= \sum_{i' \in \mathcal{N}(i_0)} \mathbb{E} [\phi_{i_0}(X_{i_0}) \phi_{i'}(X_{i'})] \\
 &+ \sum_{(i,i') \in E, i \neq i_0, i' \neq i_0} \mathbb{E} [\phi_i(X_i) \phi_{i'}(X_{i'})] \\
 &= \mathbb{E} \left[\phi_{i_0}(X_{i_0}) \sum_{i' \in \mathcal{N}(i_0)} \mathcal{P}_{i_0}(\phi_{i'}(X_{i'}))(X_{i_0}) \right] \\
 &+ \sum_{(i,i') \in E, i \neq i_0, i' \neq i_0} \mathbb{E} [\phi_i(X_i) \phi_{i'}(X_{i'})] \\
 &\leq \sqrt{\mathbb{E} \left[\left(\sum_{i' \in \mathcal{N}(i_0)} \mathcal{P}_{i_0}(\phi_{i'}(X_{i'}))(X_{i_0}) \right)^2 \right]} \\
 &+ \sum_{(i,i') \in E, i \neq i_0, i' \neq i_0} \mathbb{E} [\phi_i(X_i) \phi_{i'}(X_{i'})],
 \end{aligned}$$

where the last inequality follows from the Cauchy-Schwartz inequality and becomes an equality if $\phi_{i_0}^*(X_{i_0}) = c \sum_{i' \in \mathcal{N}(i_0)} \mathcal{P}_{i_0}(\phi_{i'}^*(X_{i'}))(X_{i_0})$, for some constant $c > 0$. Since we have $\mathbb{E}[\phi_{i_0}^*(X_{i_0})] = 0$ and $\text{var}(\phi_{i_0}^*(X_{i_0})) = 1$, we obtain

$$\phi_{i_0}^*(X_{i_0}) = \frac{\sum_{i' \in \mathcal{N}(i_0)} \mathcal{P}_{i_0}(\phi_{i'}^*(X_{i'}))(X_{i_0})}{\text{var} \left(\sum_{i' \in \mathcal{N}(i_0)} \mathcal{P}_{i_0}(\phi_{i'}^*(X_{i'}))(X_{i_0}) \right)}.$$

□

2.7.4 Proof of Theorem 2.10

The proof of this theorem is a natural extension of the proof of Proposition 2.4.

2.7.5 Proof of Theorem 2.13

Proof. To prove Theorem 2.13, first we show that the constraints $\mathbb{E}[\phi_i(X_i)] = 0$ and $\mathbb{E}[\phi_i(X_i)^2] = 1$ lead to the same optimal solution as the constraints $\text{var}(\phi_i(X_i)) = 1$.

We then show that the constraint $\mathbf{a}_i \perp \sqrt{\mathbf{p}_i}$ can be incorporated into the objective

function, without changing the optimal solution.

Lemma 2.37. The NMC optimization (2.3) can be written as follows:

$$\begin{aligned} \max_{\phi_1, \dots, \phi_n} \quad & \sum_{(i, i') \in E} \mathbb{E}[(\phi_i(X_i) - \bar{\phi}_i)(\phi_{i'}(X_{i'}) - \bar{\phi}_{i'})] \\ & \text{var}(\phi_i(X_i)) = 1, \quad 1 \leq i \leq n, \end{aligned} \quad (2.59)$$

where $\bar{\phi}_i$ and $\text{var}(\phi_i(X_i))$ represent the mean and the variance of the random variable $\phi_i(X_i)$.

Proof. In this proof, we denote the random variable $\phi_i(X_i)$ by ϕ_i . We let the optimal objective value of optimization (2.59) be $\tilde{\rho}_G$. We also let ϕ_i^* be an optimal solution of (2.3). The set of functions ϕ_i^* for $i = 1, \dots, n$ is feasible for optimization (2.59) and therefore we have $\rho_G \leq \tilde{\rho}_G$. On the other hand, let ϕ_i^{**} be an optimal solution of optimization (2.59). Let $\tilde{\phi}_i = \phi_i^{**} - \bar{\phi}_i^{**}$. The set of functions $\tilde{\phi}_i$ for $i = 1, \dots, n$ is feasible for optimization (2.3). Thus, we have $\rho_G \geq \tilde{\rho}_G$. Therefore, we have that $\rho_G = \tilde{\rho}_G$. \square

Let $\mathbf{a}_i(j) = \phi_i(j)/\sqrt{\mathbf{p}_i(j)}$ for $1 \leq j \leq |\mathcal{X}_i|$. We have

$$1 = \mathbb{E}[\phi_i(X_i)^2] - (\mathbb{E}[\phi_i(X_i)])^2 = \|\mathbf{a}_i\|_2^2 - (\mathbf{a}_i \sqrt{\mathbf{p}_i})^2 = \mathbf{a}_i^T \left(I - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_i}^T \right) \mathbf{a}_i.$$

We next show that the matrix $I - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_i}^T$ is positive semidefinite and the only vectors in its null space are $\mathbf{0}$ and $\sqrt{\mathbf{p}_i}$. This is because

$$\mathbf{x}^T \left(I - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_i}^T \right) \mathbf{x} = \|\mathbf{x}\|_2^2 - (\mathbf{x} \sqrt{\mathbf{p}_i})^2 \geq 0, \quad (2.60)$$

where we use Cauchy-Schwartz and $\|\sqrt{\mathbf{p}_i}\|_2^2 = 1$ to obtain the last inequality (2.60). This inequality becomes an equality if and only if $\mathbf{x} = \mathbf{0}$ or $\mathbf{x} = \sqrt{\mathbf{p}_i}$.

Now consider the objective function of optimization (2.59). We have

$$\begin{aligned} \mathbb{E} [(\phi_i(X_i) - \bar{\phi}_i)(\phi_{i'}(X_{i'}) - \bar{\phi}_{i'})] &= \mathbb{E} [\phi_i(X_i)\phi_{i'}(X_{i'})] - \bar{\phi}_i\bar{\phi}_{i'} \\ &= \mathbf{a}_i^T Q_{i,i'} \mathbf{a}_i - (\mathbf{a}_i^T \sqrt{\mathbf{p}_i})(\mathbf{a}_{i'}^T \sqrt{\mathbf{p}_{i'}}) = \mathbf{a}_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) \mathbf{a}_{i'}. \end{aligned}$$

Therefore, optimization (2.59) (which is equivalent to the NMC optimization (2.3) according to Lemma 2.37) can be written as,

$$\begin{aligned} \max_{\mathbf{a}_i} \quad & \sum_{(i,i') \in E} \mathbf{a}_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) \mathbf{a}_{i'} \quad (2.61) \\ & \mathbf{a}_i^T \left(I - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_i^T} \right) \mathbf{a}_i = 1, \quad 1 \leq i \leq n. \end{aligned}$$

For each i , since $I - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_i^T}$ is positive semidefinite. Thus, we can write $I - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_i^T} = B_i B_i^T$. Recall $\mathbf{b}_i = B_i \mathbf{a}_i$. Thus, constraints of optimization (2.61) can be written as $\mathbf{b}_i^T \mathbf{b}_i = \|\mathbf{b}_i\|_2^2 = 1$. We next write \mathbf{a}_i as a function of \mathbf{b}_i . Note that since B_i is not invertible, there are many choices for \mathbf{a}_i as a function of \mathbf{b}_i characterized as follows:

Let $U_i \Sigma_i U_i^T$ be the singular value decomposition of B_i . The vector $\sqrt{\mathbf{p}_i}$ is the singular vector corresponding to singular value zero ($\sigma_i^{(1)} = 0$).

$$\begin{aligned} \mathbf{a}_i &= \left([U_i^{(2)}, \dots, U_i^{(|\mathcal{X}_i|)}] \text{diag}(1/\sigma_i^{(2)}, \dots, 1/\sigma_i^{(n_i)}) [U_i^{(2)}, \dots, U_i^{(|\mathcal{X}_i|)}]^T \right) \mathbf{b}_i + \alpha_i \sqrt{\mathbf{p}_i} \\ &= A_i \mathbf{b}_i + \alpha_i \sqrt{\mathbf{p}_i}, \end{aligned} \quad (2.62)$$

where α_i can be any scalar.⁵

Below, we show that all choices of \mathbf{a}_i according to (2.62) lead to the same objective

⁵Since B_i is symmetric it has a set of $|\mathcal{X}_i|$ orthonormal eigenvectors and can be written as

$$B_i = \sum_{j=1}^{|\mathcal{X}_i|} \mathbf{v}_j \sigma_i^{(j)} \mathbf{v}_j^T.$$

We have $\mathbf{b}_i = \sum_{j=2}^{|\mathcal{X}_i|} \beta_j \mathbf{v}_j$ and $\mathbf{a}_i = \sum_{j=1}^{|\mathcal{X}_i|} \alpha_j \mathbf{v}_j$. From $\mathbf{b}_i = B_i \mathbf{a}_i$, we obtain that $\alpha_j = \beta_j / \sigma_i^{(j)}$ for $j \geq 2$, where α_1 can be any scalar.

function of optimization (2.61):

$$\begin{aligned}
 \mathbf{a}_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) \mathbf{a}_{i'} &= \mathbf{b}_i^T A_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) A_{i'} \mathbf{b}_{i'} \\
 &\quad + \mathbf{b}_i^T A_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) \alpha_{i'} \sqrt{\mathbf{p}_{i'}} \\
 &\quad + \alpha_i \sqrt{\mathbf{p}_i}^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) A_{i'} \mathbf{b}_{i'} \\
 &\quad + \alpha_i \sqrt{\mathbf{p}_i}^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) \alpha_{i'} \sqrt{\mathbf{p}_{i'}} \\
 &\stackrel{(1)}{=} \mathbf{b}_i^T A_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) A_{i'} \mathbf{b}_{i'} \\
 &\quad + \mathbf{b}_i^T A_i^T Q_{i,i'} \alpha_{i'} \sqrt{\mathbf{p}_{i'}} + \alpha_i \sqrt{\mathbf{p}_i}^T Q_{i,i'} A_{i'} \mathbf{b}_{i'} \\
 &\stackrel{(2)}{=} \mathbf{b}_i^T A_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) A_{i'} \mathbf{b}_{i'} \\
 &\quad + \mathbf{b}_i^T A_i^T \sqrt{\mathbf{p}_i} \alpha_{i'} + \alpha_i \sqrt{\mathbf{p}_{i'}^T} A_{i'} \mathbf{b}_{i'} \\
 &= \mathbf{b}_i^T A_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) A_{i'} \mathbf{b}_{i'},
 \end{aligned}$$

where (1) follows from $\sqrt{\mathbf{p}_{i'}^T} \sqrt{\mathbf{p}_{i'}} = 1$, $\sqrt{\mathbf{p}_i}^T \sqrt{\mathbf{p}_i} = 1$, $\sqrt{\mathbf{p}_i} Q_{i,i'} \sqrt{\mathbf{p}_{i'}} = 1$, $A_i^T \sqrt{\mathbf{p}_i} = 0$, and $\sqrt{\mathbf{p}_{i'}^T} A_{i'} = 0$; and (2) follows from $\sqrt{\mathbf{p}_i}^T Q_{i,i'} = \sqrt{\mathbf{p}_{i'}^T}$ and $Q_{i,i'} \sqrt{\mathbf{p}_{i'}} = \sqrt{\mathbf{p}_i}$.

Therefore, the NMC optimization (2.3) can be written as

$$\begin{aligned}
 \max_{\mathbf{b}_1, \dots, \mathbf{b}_n} \quad & \sum_{(i,i') \in E} \mathbf{b}_i^T A_i^T \left(Q_{i,i'} - \sqrt{\mathbf{p}_i} \sqrt{\mathbf{p}_{i'}^T} \right) A_{i'} \mathbf{b}_{i'} \\
 & \|\mathbf{b}_i\|_2 = 1 \quad 1 \leq i \leq n.
 \end{aligned}$$

□

2.7.6 Proof of Proposition 2.17

Proof. Let $\bar{\beta}$ be a solution of the MEP (2.33). Suppose that there exists an $1 \leq i \leq n$ such that $\lambda_i(\bar{\beta}) < 0$. Let $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_n)$ be defined as: $\hat{\beta}_i = -\bar{\beta}_i$, for any i such that $\lambda_i(\bar{\beta}) < 0$, and $\hat{\beta}_{i'} = \bar{\beta}_{i'}$, otherwise. We show that by flipping the sign of $\bar{\beta}_i$ while

keeping the rest of $\bar{\beta}_{i'}$ for $i' \neq i$ the same, $r(\cdot)$ increases. To show this we have

$$\begin{aligned} r(\hat{\beta}) - r(\bar{\beta}) &\stackrel{(1)}{=} 2\hat{\beta}_i^T \sum_{i'=1}^n C_{i,i'} \hat{\beta}_{i'} - 2\bar{\beta}_i^T \sum_{i'=1}^n C_{i,i'} \bar{\beta}_{i'} \\ &\stackrel{(2)}{=} 2(\hat{\beta}_i^T - \bar{\beta}_i^T) \sum_{i'=1}^n C_{i,i'} \bar{\beta}_{i'} \stackrel{(3)}{=} -2\bar{\beta}_i^T \sum_{i'=1}^n C_{i,i'} \bar{\beta}_{i'} \stackrel{(4)}{=} -2\lambda_i(\bar{\beta}_i) > 0, \end{aligned}$$

where equations (1) and (4) follow from (2.34), equation (2) follows from the fact that $\hat{\beta}_{i'} = \bar{\beta}_{i'}$ for $i' \neq i$, and equation (3) follows from the fact that $\hat{\beta}_i = -\bar{\beta}_i$. This completes the proof. \square

2.7.7 Proof of Theorem 2.21

Proof. For any realization of the partitioning, consider NMC over all sub-graphs G_m ($1 \leq m \leq M$) and denote the corresponding functions by ϕ_i^* for $1 \leq i \leq n$. We have

$$\begin{aligned} \rho_G &= \sum_{(i,i') \in E} \mathbb{E}[\phi_i^*(X_i)\phi_{i'}^*(X_{i'})] = \sum_{(i,i') \in E \setminus E^c} \mathbb{E}[\phi_i^*(X_i)\phi_{i'}^*(X_{i'})] + \sum_{(i,i') \in E^c} \mathbb{E}[\phi_i^*(X_i)\phi_{i'}^*(X_{i'})] \\ &\stackrel{(1)}{=} \sum_{m=1}^M \sum_{(i,i') \in E_m} \mathbb{E}[\phi_i^*(X_i)\phi_{i'}^*(X_{i'})] + \sum_{(i,i') \in E^c} \mathbb{E}[\phi_i^*(X_i)\phi_{i'}^*(X_{i'})] \\ &\stackrel{(2)}{\leq} \sum_{m=1}^M \hat{\rho}_{G_m} + \sum_{(i,i') \in E^c} \mathbb{E}[\phi_i^*(X_i)\phi_{i'}^*(X_{i'})] = \hat{\rho}_G + \sum_{(i,i') \in E^c} \mathbf{1}\{(i,i') \in E^c\} \mathbb{E}[\phi_i^*(X_i)\phi_{i'}^*(X_{i'})], \end{aligned}$$

where equation (1) comes from the graph partitioning Definition 2.18, and inequality (2) comes from the fact that $\hat{\rho}_{G_m}$ is the NMC for the partition $G_m = (V_m, E \cap (V_m \times V_m))$. Therefore, by taking expectation over the partitioning, we obtain $\rho_G \leq \mathbb{E}[\hat{\rho}_G] + \epsilon\rho_G$, which gives us $(1 - \epsilon)\rho_G \leq \mathbb{E}[\hat{\rho}_G]$. \square

2.7.8 Proof of Theorem 2.22

Proof. We first present a lemma in which we bound the difference between Q matrices, according to Definition 2.3, by the difference between probability distributions. We then use a sensitivity analysis of the optimization whose optimal solution is NMC.

Lemma 2.38. Let P and \tilde{P} be the matrix form of two joint probability distribution on $(\mathcal{X}_i, \mathcal{X}_{i'})$, such that $P_{X_i, X_{i'}}(x_i, x_{i'}) = [P]_{x_i, x_{i'}}$ and $\tilde{P}_{X_i, X_{i'}}(x_i, x_{i'}) = [\tilde{P}]_{x_i, x_{i'}}$ and where $D_{X_i}(P) = \text{diag}(P_{X_i}(x_i) : x_i \in \mathcal{X}_i)$. We can bound the difference between $Q_{i, i'}$ and $\tilde{Q}_{i, i'}$ by the difference between P and \tilde{P} , as follows:

$$\begin{aligned} \|Q_{i, i'} - \tilde{Q}_{i, i'}\|_2 &\leq \frac{1}{2\delta^2} \sqrt{K} \|D_{X_i}(P) - D_{X_i}(\tilde{P})\|_\infty + \frac{1}{2\delta^2} \sqrt{K} \|D_{X_{i'}}(P) - D_{X_{i'}}(\tilde{P})\|_\infty \\ &\quad + \frac{1}{\delta} \sqrt{K} \|P - \tilde{P}\|_\infty, \end{aligned} \tag{2.63}$$

where δ is the minimum probability of all elements of \mathcal{X}_i and $\mathcal{X}_{i'}$, under both P and \tilde{P} , and K is the maximum alphabet size. In particular, if $\|P_{X_i, X_{i'}} - \tilde{P}_{X_i, X_{i'}}\|_\infty \leq \gamma$ and $K = \max\{|\mathcal{X}_i|, |\mathcal{X}_{i'}|\}$, then we have

$$\|Q_{i, i'} - \tilde{Q}_{i, i'}\|_2 \leq 2\gamma K^n \frac{1}{\delta^2}. \tag{2.64}$$

Proof. We have

$$\|Q_{i, i'} - \tilde{Q}_{i, i'}\|_2 = \|D_{X_i}(P)^{-\frac{1}{2}} P D_{X_{i'}}(P)^{-\frac{1}{2}} - D_{X_i}(\tilde{P})^{-\frac{1}{2}} \tilde{P} D_{X_{i'}}(\tilde{P})^{-\frac{1}{2}}\|_2.$$

We next insert the terms $\pm D_{X_i}(P)^{-\frac{1}{2}} P D_{X_{i'}}(\tilde{P})^{-\frac{1}{2}}$ and $\pm D_{X_i}(P)^{-\frac{1}{2}} \tilde{P} D_{X_{i'}}(\tilde{P})^{-\frac{1}{2}}$, pair the terms, and use triangle inequality to obtain

$$\begin{aligned} \|Q_{i, i'} - \tilde{Q}_{i, i'}\|_2 &\leq \|D_{X_i}(P)^{-\frac{1}{2}} P\|_2 \|D_{X_{i'}}(P)^{-\frac{1}{2}} - D_{X_{i'}}(\tilde{P})^{-\frac{1}{2}}\|_2 \\ &\quad + \|D_{X_{i'}}(\tilde{P})^{-\frac{1}{2}}\|_2 \|\tilde{P}\|_2 \|D_{X_i}(P)^{-\frac{1}{2}} - D_{X_i}(\tilde{P})^{-\frac{1}{2}}\|_2 \\ &\quad + \|D_{X_{i'}}(\tilde{P})^{-\frac{1}{2}}\|_2 \|D_{X_i}(P)^{-\frac{1}{2}}\|_2 \|P - \tilde{P}\|_2, \end{aligned} \tag{2.65}$$

where each term depends on the difference between P and \tilde{P} . We have

$$\begin{aligned} \|D_{X_i}(P)^{-\frac{1}{2}}P\|_2 &\stackrel{(1)}{\leq} \frac{1}{\sqrt{\delta}}\sqrt{K} & (2.66) \\ \|D_{X_{i'}}(P)^{-\frac{1}{2}} - D_{X_{i'}}(\tilde{P})^{-\frac{1}{2}}\|_2 &\leq \|D_{X_{i'}}(P)^{-\frac{1}{2}} - D_{X_{i'}}(\tilde{P})^{-\frac{1}{2}}\|_\infty \\ &\stackrel{(2)}{\leq} \|D_{X_{i'}}(P) - D_{X_{i'}}(\tilde{P})\|_\infty \frac{1}{2\delta^{3/2}}, \end{aligned}$$

where inequality (1) comes from the Cauchy-Schwartz inequality and the fact that $\|P\|_2 \leq \sqrt{K}\|P\|_1 \leq \sqrt{K}$. Inequality (2) follows from the fact that for $x, y \in \mathbb{R}$, we have

$$\frac{1}{\sqrt{x}} - \frac{1}{\sqrt{y}} = \frac{y-x}{\sqrt{x}\sqrt{y}(\sqrt{x}+\sqrt{y})} \leq |x-y| \frac{1}{2(\min\{x,y\})^{3/2}}. \quad (2.67)$$

Using (2.66) in (2.65) results in (2.63). Moreover, using $\|P_{X_1, \dots, X_n} - \tilde{P}_{X_1, \dots, X_n}\|_\infty \leq \gamma$ and marginalizing over $X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n$ leads to

$$\|D_{X_i}(P) - D_{X_i}(\tilde{P})\|_\infty \leq \gamma K^{n-1}. \quad (2.68)$$

Using (2.68) and the fact that $K > 1$ in (2.63) results in (2.64). This completes the proof of this Lemma. \square

Let P and \tilde{P} be two distributions on $\mathcal{X}_1 \times \dots \times \mathcal{X}_n$. We shall compare the solution of the two following optimization problems.

$$\begin{aligned} \max_{\mathbf{a}_i} \quad &\sum_{(i,i') \in E} \mathbf{a}_i^T Q_{i,i'} \mathbf{a}_i & (2.69) \\ &\|\mathbf{a}_i\|_2 = 1, \quad 1 \leq i \leq n, \\ &\mathbf{a}_i \perp \sqrt{\mathbf{p}_i}, \quad 1 \leq i \leq n, \end{aligned}$$

and

$$\begin{aligned}
 \max_{\mathbf{a}_i} \quad & \sum_{(i,i') \in E} \mathbf{a}_i^T \tilde{Q}_{i,i'} \mathbf{a}_i & (2.70) \\
 & \|\mathbf{a}_i\|_2 = 1, \quad 1 \leq i \leq n, \\
 & \mathbf{a}_i \perp \sqrt{\tilde{\mathbf{p}}_i}, \quad 1 \leq i \leq n.
 \end{aligned}$$

Let ρ_G and $\tilde{\rho}_G$ be the optimal values for (2.69) and (2.70), respectively. Recall that

$$\delta = \min_{1 \leq i \leq n} \left(\min\{\delta_{X_i}(P), \delta_{X_i}(\tilde{P})\} \right).$$

Suppose \mathbf{a}_i^* yields to an optimal solution of optimization (2.69). Based on this solution, we shall construct a feasible solution for optimization (2.70) and then evaluate its objective function. For any i , let $\mathbf{c}_i = \frac{\mathbf{a}_i^* + \nu_i}{\|\mathbf{a}_i^* + \nu_i\|_2}$, where $\nu_i = \sqrt{\tilde{\mathbf{p}}_i} \langle \mathbf{a}_i^*, \sqrt{\mathbf{p}_i} - \sqrt{\tilde{\mathbf{p}}_i} \rangle$. Under the assumption of Theorem 2.22 (i.e., $\gamma \leq \delta^{3/2} \mathcal{K}^{-1}$), we show that $\|\nu_i\| \leq 1/2$.

We have

$$\begin{aligned}
 \|\nu_i\| &= \|\sqrt{\tilde{\mathbf{p}}_i} \langle \mathbf{a}_i^*, \sqrt{\mathbf{p}_i} - \sqrt{\tilde{\mathbf{p}}_i} \rangle\|_2 \stackrel{(3)}{\leq} \|\sqrt{\tilde{\mathbf{p}}_i}\|_2 \|\mathbf{a}_i^*\|_2 \|\sqrt{\mathbf{p}_i} - \sqrt{\tilde{\mathbf{p}}_i}\|_2 & (2.71) \\
 &\stackrel{(4)}{\leq} \|\mathbf{p}_i - \tilde{\mathbf{p}}_i\|_\infty \frac{1}{2\delta^{3/2}} \stackrel{(5)}{\leq} \frac{\gamma K^{n-1}}{2\delta^{3/2}} \stackrel{(6)}{\leq} \frac{1}{2},
 \end{aligned}$$

where Inequality (3) comes from the Cauchy-Schwartz, Inequality (4) comes from the facts that $\|\sqrt{\tilde{\mathbf{p}}_i}\|_2 = 1$, $\|\mathbf{a}_i^*\|_2 = 1$, (2.67), Inequality (5) comes from the fact that $\|\mathbf{p}_i(j) - \tilde{\mathbf{p}}_i(j)\|_\infty \leq \gamma K^{n-1}$, and Inequality (6) comes from the fact that $\frac{\gamma K^n}{\delta^{3/2}} \leq 1$ by the assumption of Theorem 2.22.

We claim \mathbf{c}_i is feasible for optimization (2.70). First note that the norm of each \mathbf{c}_i is one. We next show that each \mathbf{c}_i is orthogonal to $\sqrt{\tilde{\mathbf{p}}_i}$:

$$\langle \mathbf{c}_i, \sqrt{\tilde{\mathbf{p}}_i} \rangle = \frac{\langle \mathbf{a}_i^*, \sqrt{\tilde{\mathbf{p}}_i} \rangle + \langle \mathbf{a}_i^*, \sqrt{\mathbf{p}_i} - \sqrt{\tilde{\mathbf{p}}_i} \rangle \|\sqrt{\tilde{\mathbf{p}}_i}\|_2^2}{\|\mathbf{a}_i^* + \sqrt{\tilde{\mathbf{p}}_i} \langle \mathbf{a}_i^*, \sqrt{\mathbf{p}_i} - \sqrt{\tilde{\mathbf{p}}_i} \rangle\|_2} = 0,$$

where the last equality follows from $\|\sqrt{\tilde{\mathbf{p}}_i}\|_2 = 1$ and $\mathbf{a}_i \perp \sqrt{\tilde{\mathbf{p}}_i}$. We now plug in the

feasible solution \mathbf{c}_i into the objective function of optimization (2.70).

$$\begin{aligned}
 \tilde{\rho}_G &\stackrel{(7)}{\geq} \sum_{(i,i') \in E} \mathbf{c}_i^T \tilde{Q}_{ii'} \mathbf{c}_{i'} & (2.72) \\
 &\stackrel{(8)}{=} \sum_{(i,i') \in E} \left(\mathbf{a}_i^{*T} Q_{i,i'} \mathbf{a}_{i'}^* + \mathbf{c}_i^T (\tilde{Q}_{i,i'} - Q_{i,i'}) \mathbf{c}_{i'} + (\mathbf{c}_i^T - \mathbf{a}_i^{*T}) Q_{i,i'} \mathbf{c}_{i'} + \mathbf{a}_i^{*T} Q_{i,i'} (\mathbf{c}_{i'} - \mathbf{a}_{i'}^*) \right) \\
 &\stackrel{(9)}{\geq} \rho_G - \sum_{(i,i') \in E} \left(\underbrace{|\mathbf{c}_i^T (\tilde{Q}_{i,i'} - Q_{i,i'}) \mathbf{c}_{i'}|}_{\text{Term I}} + \underbrace{|(\mathbf{c}_i^T - \mathbf{a}_i^{*T}) Q_{i,i'} \mathbf{c}_{i'}|}_{\text{Term II}} + \underbrace{|\mathbf{a}_i^{*T} Q_{i,i'} (\mathbf{c}_{i'} - \mathbf{a}_{i'}^*)|}_{\text{Term III}} \right),
 \end{aligned}$$

where Inequality (7) comes from the fact that \mathbf{c}_i for $1 \leq i \leq n$ is a feasible solution of Optimization (2.70), Equality (8) comes from adding and subtracting $\mathbf{a}_i^{*T} Q_{i,i'} \mathbf{a}_{i'}^*$, $\mathbf{c}_i^T Q_{i,i'} \mathbf{c}_{i'}$ and $\mathbf{a}_i^{*T} Q_{i,i'} \mathbf{c}_{i'}$, and Inequality (9) comes from the triangle inequality and the fact that \mathbf{a}_i^* for $1 \leq i \leq n$ is an optimal solution of Optimization (2.69). Next, we bound terms I-III on the righthand side of this relation.

Using Lemma 2.38, for any i, i' , we have

$$|\mathbf{c}_i^T (\tilde{Q}_{i,i'} - Q_{i,i'}) \mathbf{c}_{i'}| \leq \|\tilde{Q}_{i,i'} - Q_{i,i'}\|_2 \leq 2 \frac{\gamma K^n}{\delta^2}. \quad (2.73)$$

We also have

$$|(\mathbf{c}_i^T - \mathbf{a}_i^{*T}) Q_{i,i'} \mathbf{c}_{i'}| \leq \|\mathbf{c}_i^T - \mathbf{a}_i^{*T}\|_2 \|Q_{i,i'}\|_2 \|\mathbf{c}_{i'}\|_2 = \|\mathbf{c}_i^T - \mathbf{a}_i^{*T}\|_2 \leq 2 \frac{\|\nu_i\|_2}{1 - \|\nu_i\|_2},$$

where we use the following inequality

$$\begin{aligned}
 \left\| \frac{\mathbf{a} + \nu_i}{\|\mathbf{a} + \nu_i\|_2} - \mathbf{a} \right\|_2 &= \left\| \frac{\nu_i}{\|\mathbf{a} + \nu_i\|_2} + \mathbf{a} \left(\frac{1}{\|\mathbf{a} + \nu_i\|_2} - 1 \right) \right\|_2 & (2.74) \\
 &\leq \left\| \frac{\nu_i}{\|\mathbf{a} + \nu_i\|_2} \right\|_2 + \|\mathbf{a}\|_2 \left| \frac{1}{\|\mathbf{a} + \nu_i\|_2} - 1 \right| \\
 &\leq \frac{\|\nu_i\|_2}{1 - \|\nu_i\|_2} + \max \left\{ \left| \frac{1}{1 - \|\nu_i\|_2} - 1 \right|, \left| \frac{1}{1 + \|\nu_i\|_2} - 1 \right| \right\} \leq 2 \frac{\|\nu_i\|_2}{1 - \|\nu_i\|_2},
 \end{aligned}$$

where we used $1 - \|\nu_i\|_2 \leq \|\mathbf{a} + \nu_i\|_2 \leq 1 + \|\nu_i\|_2$ and the fact that $\|\nu_i\|_2 \leq 1/2$

according to (2.71). Using (2.71) and (2.74), we obtain

$$|(\mathbf{c}_i^T - \mathbf{a}_i^{*T}) Q_{i,i'} \mathbf{c}_{i'}| \leq 2 \frac{\gamma K^{n-1}}{\delta^{3/2}}, \quad (2.75)$$

Using (2.73) and (2.75) in (2.72) leads to

$$\tilde{\rho}_G \geq \rho_G - \sum_{(i,i') \in E} \left(4 \frac{\gamma K^n}{\delta^2} + 4\gamma K^{n-1} \frac{1}{\delta^{3/2}} \right) \geq \rho_G - \gamma K^n |E| \frac{8}{\delta^2}.$$

Similarly, we have

$$\rho_G \geq \tilde{\rho}_G - \sum_{(i,i') \in E} \left(4 \frac{\gamma K^n}{\delta^2} + 4\gamma K^{n-1} \frac{1}{\delta^{3/2}} \right) \geq \tilde{\rho}_G - \gamma K^n |E| \frac{8}{\delta^2}.$$

Combining the previous two relations, we obtain

$$|\tilde{\rho}_G - \rho_G| \leq \gamma K^n |E| \frac{8}{\delta^2},$$

which completes the proof. \square

2.7.9 Proof of Theorem 2.24

Proof. We use the following theorem in the proof.

Theorem 2.39 ([89, 90]). Let P be a probability distribution on a finite alphabet \mathcal{X} . Also, let $P^{(m)}$ denote the empirical probability distribution of X , obtained from m i.i.d. samples, $\{x_i\}_{i=1}^m$, drawn according to P . We have

$$\mathbb{P} [\|P^{(m)} - P\|_\infty > \gamma] \leq 4e^{-m\frac{\gamma^2}{2}}.$$

Theorem 2.39 establishes that P and $P^{(m)}$ are close to each other for a sufficiently large sample size m . Moreover, one can use Theorem 2.22 to bound the difference between NMC values of distributions P and $P^{(m)}$ by $\|P^{(m)} - P\|_\infty$. However, to prove Theorem 2.24, we also need to bound the minimum sample probability of $P^{(m)}$. We

explain these steps below.

We have

$$\begin{aligned} \mathbb{P}\left[\delta_{X_i}(P^{(m)}) \geq \frac{\delta'}{2}\right] &\geq \mathbb{P}\left[\|P_{X_i}^{(m)} - P_{X_i}\|_\infty \leq \frac{\delta'}{2}\right] \geq \mathbb{P}\left[\|P^{(m)} - P\|_\infty \leq \frac{\delta'}{2K^{n-1}}\right] \\ &\geq 1 - \mathbb{P}\left[\|P^{(m)} - P\|_\infty > \frac{\delta'}{2K^{n-1}}\right] \geq 1 - 4e^{-m\frac{1}{2}\left(\frac{\delta'}{2K^{n-1}}\right)^2}, \end{aligned}$$

where the first inequality follows from $\delta_{X_i}(P^{(m)}) \geq -|\delta_{X_i}(P^{(m)}) - \delta_{X_i}(P)| + \delta_{X_i}(P) \geq \delta - \|P_{X_i}^{(m)} - P_{X_i}\|_\infty$ and the second inequality follows from $\|P_{X_i}^{(m)} - P_{X_i}\|_\infty \leq K^{n-1}\|P^{(m)} - P\|_\infty$. Using the union bound, we obtain

$$\mathbb{P}\left[\bigcap_{i \in V} \left\{\delta_{X_i}(P^{(m)}) \geq \frac{\delta'}{2}\right\}\right] = 1 - \mathbb{P}\left[\bigcup_{i \in V} \left\{\delta_{X_i}(P^{(m)}) \leq \frac{\delta'}{2}\right\}\right] \geq 1 - 4ne^{-m\frac{1}{2}\left(\frac{\delta'}{2K^{n-1}}\right)^2}.$$

By applying the union bound once more, we have

$$\mathbb{P}\left[\left\{\|P^{(m)} - P\|_\infty < \gamma\right\} \cap \left\{\bigcap_{i \in V} \left\{\delta_{X_i}(P^{(m)}) \geq \frac{\delta'}{2}\right\}\right\}\right] \geq 1 - 4ne^{-m\frac{1}{2}\left(\frac{\delta'}{2K^{n-1}}\right)^2} - 4e^{-m\frac{\gamma^2}{2}}.$$

Thus, with a large probability, i.e., $1 - 4ne^{-m\frac{1}{2}\left(\frac{\delta'}{2K^{n-1}}\right)^2} - 4e^{-m\frac{\gamma^2}{2}}$, the minimum probability of $P^{(m)}$ is bounded by $\delta = \delta'/2$ and $\|P^{(m)} - P\|_\infty$ is bounded by γ .

For given $\eta < 1$ and $\epsilon < 1$, we let $\gamma = \frac{\epsilon(\delta')^2}{32|E|}K^{-n}$ and

$$m_0 \geq 2 \left(\frac{32|E|K^n}{\epsilon(\delta')^2}\right)^2 \log\left(\frac{8n}{\eta}\right).$$

Now we use Theorem 2.22 with $\delta = \delta'/2$. First note that $\gamma \leq \delta^{3/2}\mathcal{K}^{-1}$ which is required as an assumption for invoking Theorem 2.22. This leads to

$$|\rho_G^{(m)} - \rho_G| \leq \gamma|E|K^n \frac{8}{\delta^2} \leq \epsilon,$$

with probability (at least)

$$1 - \left(4ne^{-m_0\frac{1}{2}\left(\frac{\delta'}{2K^{n-1}}\right)^2} - 4e^{-m_0\frac{\gamma^2}{2}}\right) \geq 1 - \eta.$$

This completes the proof. \square

2.7.10 Proof of Theorem 2.26

Proof. Let $\mathbf{a} = (\mathbf{a}_1, \dots, \mathbf{a}_n)$. Define

$$U_L(\mathbf{a}) \triangleq \sum_{(i,i') \in E} \sum_{j=1}^L a_{i,j} a_{i',j} (\rho_{i,i'})^j \quad (2.76)$$

$$U(\mathbf{a}) \triangleq \sum_{(i,i') \in E} \sum_{j=1}^{\infty} a_{i,j} a_{i',j} (\rho_{i,i'})^j.$$

Let $\tilde{\mathbf{a}}^{(L)} = (\tilde{\mathbf{a}}_1^{(L)}, \dots, \tilde{\mathbf{a}}_n^{(L)})$ and $\tilde{\mathbf{a}} = (\tilde{\mathbf{a}}_1, \dots, \tilde{\mathbf{a}}_n)$ be optimal solutions of the following optimizations

$$\max_{\mathbf{a}} U_L(\mathbf{a}) \quad (2.77)$$

$$\mathbf{a} \in \mathcal{C}',$$

and

$$\max_{\mathbf{a}} U(\mathbf{a}) \quad (2.78)$$

$$\mathbf{a} \in \mathcal{C}',$$

respectively, where \mathcal{C}' is the feasible set of the optimization (2.46). We first show for any L and any $1 \leq i \leq n$, $\tilde{\mathbf{a}}_i^{(L)} = (0, s_i, \dots, 0)$. We then show for any $1 \leq i \leq n$, $\tilde{\mathbf{a}}_i$ must be equal to $(0, s_i, \dots)$, completing the proof.

Claim 1: We first characterize a global maximizer of optimization (2.77). Note that the constraint set \mathcal{C}' can be truncated w.l.o.g. to consider $1 \leq j \leq L$. Let Λ be the matrix of correlation coefficients where $[\Lambda]_{i,i'} = \rho_{i,i'}$. Diagonal elements of Λ are all zero, as we ignore self-loops. Define

$$\mathbf{x} = (a_{1,1}, a_{2,1}, \dots, a_{n,1}, a_{1,2}, \dots, a_{n,L})^T. \quad (2.79)$$

Moreover, define A_0 as an $nL \times nL$ matrix composed of L^2 blocks of size $n \times n$ whose m -th diagonal block is equal to $2 \Lambda^m$, where $\Lambda^m[i, j] \triangleq (\Lambda[i, j])^m$. Off-diagonal blocks of A_0 are all zeros. Moreover, define A_i for $1 \leq i \leq n$ as an $nL \times nL$ matrix where $A_i[i + (m - 1)n, i + (m - 1)n] = 1$ for $1 \leq m \leq L$, otherwise it is zero. Therefore, optimization (2.77) can be re-written as the following standard quadratic optimization:

$$\begin{aligned} \max_{\mathbf{x}} \quad & \frac{1}{2} \mathbf{x}^T A_0 \mathbf{x} \\ & \frac{1}{2} \mathbf{x}^T A_i \mathbf{x} - \frac{1}{2} \leq 0, \quad 1 \leq i \leq n. \end{aligned} \quad (2.80)$$

Note that equality constraints of optimization (2.77) are replaced by inequality ones in optimization (2.80). This is because, since $(\rho_{i,i'})^2 \geq 0$, optimal solutions of optimization (2.80) occur in the boundary of its feasible set. Optimization (2.80) is a non-convex quadratic optimization with quadratic constraints. Reference [91] (Proposition 3.2)⁶ characterizes necessary and sufficient conditions for global minimizers of a generalized form of optimization (2.80). Let

$$\bar{\mathbf{x}} = [s_1, s_2, \dots, s_n, 0, \dots, 0]^T, \quad (2.81)$$

where $s_i \in \{-1, 1\}$, for $1 \leq i \leq n$. Using Proposition 3.2 and Theorem 3.1 of [91], to have $\bar{\mathbf{x}}$ as a global minimizer of optimization (2.80), we need to have

$$\left(\sum_{i=1}^n \lambda_i A_i - A_0 \right) \bar{\mathbf{x}} = 0 \quad (2.82)$$

and

$$\sum_{i=1}^n \lambda_i A_i - A_0 \succeq 0, \quad (2.83)$$

⁶Optimization (2.80) can be stated as a minimization by replacing $-A_0$ instead of A_0 to use Proposition 3.2 of Reference [91]

where $\lambda_i \geq 0$, and $A \succeq 0$ means that A is a positive semi-definite matrix. Using definitions of A_0 , A_i , and $\bar{\mathbf{x}}$, equation (2.82) is satisfied iff

$$\lambda_i = 2 \sum_{i' \neq i} s_i s_{i'} \rho_{i,i'} \geq 0, \quad 1 \leq i \leq n. \quad (2.84)$$

Using (2.84) and Gerschgorin's circle theorem, if

$$\sum_{i' \neq i} (1 - s_i s_{i'}) \rho_{i,i'} \geq 0, \quad \forall 1 \leq i \leq n, \quad (2.85)$$

$$\sum_{i' \neq i} s_i s_{i'} \rho_{i,i'} \geq \sum_{i' \neq i} \rho_{i,i'}^2, \quad \forall 1 \leq i \leq n, \quad (2.86)$$

conditions (2.83) are satisfied. Thus, $\bar{\mathbf{x}}$ is a global minimizer of optimization (2.80), establishing that for any L and any $1 \leq i \leq n$ we have $\tilde{\mathbf{a}}_i^{(L)} = (0, 1, \dots, 0)$.

Claim 2: Let $\mathbf{a}^* = (\mathbf{a}_1^*, \dots, \mathbf{a}_n^*)$ where $\mathbf{a}_i^* = (0, s_i, \dots)$ for all $1 \leq i \leq n$. We next show that $\tilde{\mathbf{a}} = \mathbf{a}^*$. We proceed by contradiction. Let

$$\Delta \triangleq U(\tilde{\mathbf{a}}) - U(\mathbf{a}^*). \quad (2.87)$$

By the contradiction assumption $\Delta > 0$. Since $(\rho_{i,i'})^L \rightarrow 0$ as $L \rightarrow \infty$ for all $i \neq i'$, $U_L(\cdot)$ converges uniformly to $U(\cdot)$. Thus there exists L_0 such that for $L \geq L_0$ and any \mathbf{a} , we have

$$|U_L(\mathbf{a}) - U(\mathbf{a})| \leq \frac{\Delta}{4}. \quad (2.88)$$

Therefore we have

$$U_L(\tilde{\mathbf{a}}) \geq U(\tilde{\mathbf{a}}) - \frac{\Delta}{4} \quad (2.89)$$

$$\begin{aligned} &= U(\mathbf{a}^*) + \Delta - \frac{\Delta}{4} \\ &\geq U_L(\mathbf{a}^*) - \frac{\Delta}{4} + \Delta - \frac{\Delta}{4} \\ &= U_L(\mathbf{a}^*) + \frac{\Delta}{2} > U_L(\mathbf{a}^*). \end{aligned} \quad (2.90)$$

This is in contradiction with the assumption that \mathbf{a}^* is the maximizer of $U_L(\cdot)$. Putting Claims 1 and 2 together completes the proof. \square

2.7.11 Proof of Proposition 2.27

Proof. Under assumptions of Theorem 2.26 and using the definition of Hermite-Chebyshev polynomials (2.12), we can restrict the feasible set of optimization (2.3) to the set of functions $\phi_i(X_i) = s_i X_i$ where $s_i \in \{-1, 1\}$ for all $1 \leq i \leq n$. Moreover, we have

$$\mathbb{E}[\phi_i(X_i)\phi_{i'}(X_{i'})] = s_i s_{i'} \rho_{i,i'}.$$

Furthermore, $\mathbb{E}[s_i X_i] = \mathbb{E}[X_i] = 0$, and $\mathbb{E}[(s_i X_i)^2] = \mathbb{E}[X_i^2] = 1$. This completes the proof. \square

2.7.12 Proof of Theorem 2.32

Proof. The first part of this Theorem comes directly from the argument of Example 2.6. To prove the second part, we re-write optimization (2.51) as follows:

$$\max \sum_{(i,i')} \mathbb{E}[g_i(f_i(X_i)) g_{i'}(f_{i'}(X_{i'}))], \quad (2.91)$$

$$\mathbb{E}[g_i(f_i(X_i))] = 0, \quad 1 \leq i \leq n,$$

$$\mathbb{E}[g_i(f_i(X_i)^2)] = 1, \quad 1 \leq i \leq n.$$

Define $\phi_i(X_i) = g_i(f_i(X_i))$ for $1 \leq i \leq n$. Since f_i 's are bijective, feasible regions of optimizations (2.91) and (2.3) are equal. Under the assumptions of Corollary 2.29, $\phi_i^*(X_i) = X_i$. Thus, $g_i^*(f_i(X_i)) = X_i$. \square

Chapter 3

Network Integration

The temporal and spatial expression of genes is coordinated by a hierarchy of transcription factors (TFs), whose interactions with each other and with their target genes form directed regulatory networks [4]. In addition to individual interactions, the structure of a regulatory network captures a broad systems-level view of regulatory and functional processes, since genes cluster into modules that perform similar functions [92–94]. Accurate inference of these regulatory networks is important both in the recovery and functional characterization of gene modules, and for comparative genomics of regulatory networks across multiple species [95, 96]. This is especially important because animal genomes, as fly, worm, and mouse are routinely used as models for human disease [13, 97].

Here, we infer regulatory networks of human, and model organisms *D. melanogaster* fly, and *C. elegans* worm, three of the most distant and deeply studied metazoan species. To infer regulatory interactions among transcription factors and target genes in each species, we combine genome-wide transcription factor binding profiles, conserved sequence motif instances [98] and gene expression levels [99, 100] for multiple cell types that have been collected by the ENCODE and modENCODE consortia. The main challenge is to integrate these diverse evidence sources of gene regulation in order to infer robust and accurate regulatory interactions for each species.

Ideally, inference of regulatory networks would involve performing extensive all-against-all experiments of chromatin immune-precipitation (ChIP) assays for every known transcription factor in every cell type of an organism, in order to identify all potential targets of TFs, followed by functional assays to verify that a TF-gene interaction is functional [94, 101]. However, the combinatorial number of pairs of TFs and cell types makes this experiment prohibitively expensive, necessitating the use of methods to reduce dimensionality of this problem. Here, we first infer three types of feature-specific regulatory connections based on functional and physical evidences and then integrate them to infer regulatory interactions in each species (Figure 3-1-a). One feature-specific network is based on using sequence motifs to scan the genome for instances of known binding sites of each TF, and then match predicted binding instances to nearby genes (a motif network). A second approach is to map TFs

to genes nearby their ChIP peaks using a window-based approach (a ChIP binding network). The third feature specific network uses gene expression profiles under different conditions in order to find groups of genes that are correlated in expression and therefore likely to function together (an expression-based network).

Previous work [94] has shown that, while ChIP networks are highly informative of true regulatory interactions, the number of experiments that can be carried out is typically very small, yielding a small number of high confidence interactions. Motif networks tend to be less informative than ChIP networks, but yield more coverage of the regulatory network, while co-expression based networks tend to include many false-positive edges and are the least informative [101, 102]. However, integration of these three networks [3, 93, 103, 104] into a combined network yield better performance than the individual networks in terms of recovering known regulatory interactions, by predicting interactions that tend to be supported by multiple lines of evidence. Here, we use an integration approach that combines interaction ranks across networks [3]. Inferred regulatory interactions show significant overlap with known interactions in human and fly, indicating the accuracy and robustness of the used inference pipeline. In the following, we explain our network inference framework with more details ¹.

3.1 Inference of Feature Specific Regulatory Networks

For each species, we form feature-specific regulatory networks using functional (gene expression profiles) and physical (motif sequences and ChIP peaks) evidences as follows:

Functional association networks. Expression-based networks represent interactions among TFs and target genes which are supported by correlation in gene expression levels across multiple samples [4, 106–108]. There are several methods to infer regulatory networks using gene expression profiles [3]. The input for these algorithms is a gene by condition matrix of expression values. The output of these methods are expression-based regulatory networks. We use the side information of

¹Some parts of this chapter have also been presented in [105].

TF lists to remove outgoing edges from target genes (in fact, TF lists are used as inputs to network inference algorithms to enhance their performance by limiting search space of the methods.).

To reduce bias and obtain a single expression-based network for each species, we combine results of two different expression-based network inference methods (Figure 3-1-a): one method is CLR [99] (context likelihood of relatedness) which constructs expression networks using mutual information among gene expression profiles along with a correction step to eliminate background correlations. The second method used is GENIE3 [100] (Gene Network Inference with Ensemble of Trees) which is a tree-based ensemble method that decomposes the network inference problem to several feature selection subproblems. In each subproblem, it identifies potential regulators by performing a regression analysis using random forest. GENIE3 has been recognized as the top-performing expression based inference method in the DREAM5 challenge [3].

Table 3.1 summarizes the number of genes and TFs in expression-based regulatory networks. These numbers refer to genes and TFs that are mapped to Entrez Gene IDs [109], the standard IDs that we use throughout our analysis. As it is illustrated in this table, expression based networks cover most of potential regulatory edges from TFs to targets. Despite a high coverage, however, the quality of inferred expression networks are lower than the one for physical networks [101]. This can be partially owing to indirect effects and transitive interactions in expression-based regulatory networks [3].

Physical association networks. We form two physical regulatory networks for each of the considered species using two types of physical evidences as our inference features: In the first approach, we use conserved occurrences of known sequence motifs [98], while in the second approach, we use experimentally defined TF binding occupancy profiles from ChIP assays of ENCODE and modENCODE [94, 101]. Table 3.2 shows the number of TFs associated to motifs as well as the number of TFs with genome-wide ChIP profiles in human, fly and worm. TSS coordinates are based on the genome annotations from ENCODE and modENCODE for human and worm, respectively, and the FlyBase genome annotations (FB5.48) for fly.

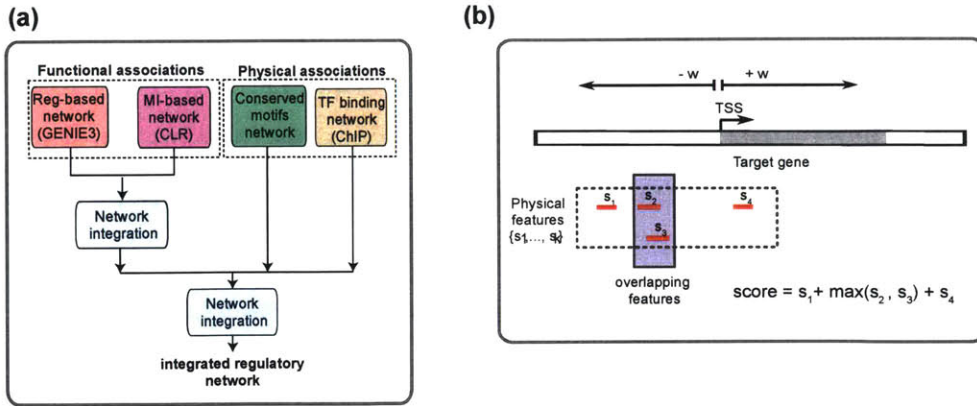


Figure 3-1: (a) The proposed framework to infer integrative regulatory networks. (b) The proposed framework to infer physical feature-specific regulatory networks.

Each physical feature is assigned to a score: motif sequence features are assigned to conservation scores according to a phylogenetic framework [98], while sequence read density of TFs determines scores of ChIP peaks. Further, two physical features are called overlapping if their corresponding sequences have a minimum overlap of 25% in relation to their lengths (Jaccard Index > 0.25).

Our inference algorithm is based on occurrence of these features (motif sequences or ChIP peaks) within a fixed window size around the transcription start site (TSS) of target genes (Figure 3-1-b). We use a fixed window of 5kb around the transcription start site (TSS) in human and 1kb in fly and worm. Then, we apply a max-sum algorithm to assign weights to TF-target interactions in each case: we take the maximum score of overlapping features and sum the scores of non-overlapping ones. In ChIP networks, because read densities are not comparable across different TFs, we normalize TF-target weights for each TF by computing z-scores.

3.2 Inference of Integrated Regulatory Networks

Feature specific networks have certain biases and shortcomings. While Physical networks (motif and ChIP networks) show high quality considering overlap of their interactions with known interactions [101], their coverage of the entire network is pretty low mostly owing to the cost of the experiments. On the other hand, while

expression based networks have a larger coverage of regulatory networks compared to physical ones, they include many false-positive edges partially owing to indirect information flows [102]. To overcome these limitations, we therefore integrate these feature-specific regulatory interactions into a single integrated network [3,93,103,104] (Figure 3-1-a).

Suppose there are K input feature-specific regulatory networks, each with n genes and m TFs (only TF nodes can have out-going edges in the network). Let $w_{i,j}^l$ and $w_{i,j}$ represent interaction weights between TF i and target gene j in the input network l and in the integrative network, respectively. We use a rank-based (borda) integration technique to infer integrated networks in considered species. In this approach, integrative weights are computed as follows:

$$w_{i,j} = 1/K \sum_{l=1}^K r_{i,j}^l, \quad (3.1)$$

where $r_{i,j}^l$ represents the rank of interactions $i \rightarrow j$ in the input network l . An edge with the maximum weight is mapped to the rank nm . We also assume non-existent interactions are mapped to rank 0 (if $w_{i,j}^l = 0$, then $r_{i,j}^l = 0$) [3]. Moreover, ties are broken randomly among edges with same weights.

We find that top-ranking integrative interactions in human and fly networks are primarily supported by ChIP and motif evidences, while worm interactions are primarily supported by co-expression edges, consistent with the lower coverage of worm ChIP and motif interactions (Figure 3-2).

To validate inferred integrated networks, we use known interactions in TRANSFAC [110], REDfly [111] and EdgeDB [112] as human, fly and worm benchmarks, respectively. We assess the quality of various networks by using (a) the area under the receiver operating characteristic curve (AUROC); and (b) the area under the precision recall curve (AUPR), for each benchmark network (Figures 3-3). Let $TP(k)$ and $FP(k)$ represent the number of true positives and false positives in top k predictions, respectively. Suppose the total number of positives and negatives in the gold standard are represented by P and N , respectively. Then, an ROC curve plots

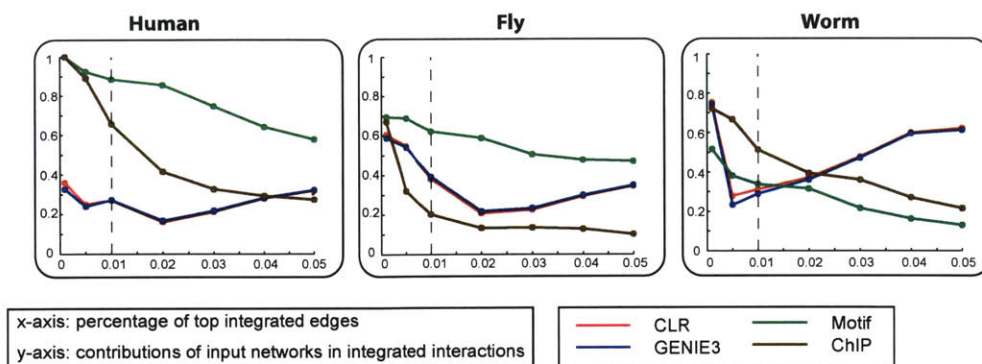


Figure 3-2: Contributions of input feature-specific networks in integrated interactions.

	Human	Fly	Worm
Genes	19,088	12,897	19,277
TFs	2,757	675	905

Table 3.1: Number of genes and TFs covered by gene expression data.

true positive rate vs. false positive rate ($TP(k)/P$ vs. $FP(k)/N$), while a PR curve plots precision ($TP(k)/k$) vs. recall ($TP(k)/P$). A high AUPR value indicates that, top predictions significantly overlap with known interactions, while a high AUROC value indicates the advantage of inferred predictions in discriminating true and false positives compared to random predictions (AUROC of a random predictor is 0.5).

Figure 3-3 illustrates AUROC and AUPR scores for feature-specific and integrative networks, in different cut-offs, and in all three considered species. Considering the top 5% of interactions in each weighted network as predicted edges, according to AUROC metric, integrative networks outperform feature-specific networks in all three species. In fact, AUROC values of integrative networks are 0.58 in human, 0.62 in fly, and 0.52 in worm, respectively. AUPR values of integrative networks are 0.019 in human, 0.047 in fly, and 0.037 in worm, respectively. Notably, all methods have low scores over the EdgeDB (worm) benchmark, which can be partially owing to sparse physical networks and/or systematic bias of EdgeDB interactions.

As the cut-off (network density) increases, AUROC values of integrative networks tend to increase while their AUPR scores are decreasing in general. This is because of the fact that, the rate of true positives is lower among medium ranked interactions

	Human	Fly	Worm
Motif network	485	221	30
ChIP network	165	51	88

Table 3.2: Number of TFs covered by evolutionary conserved motifs and TF binding datasets.

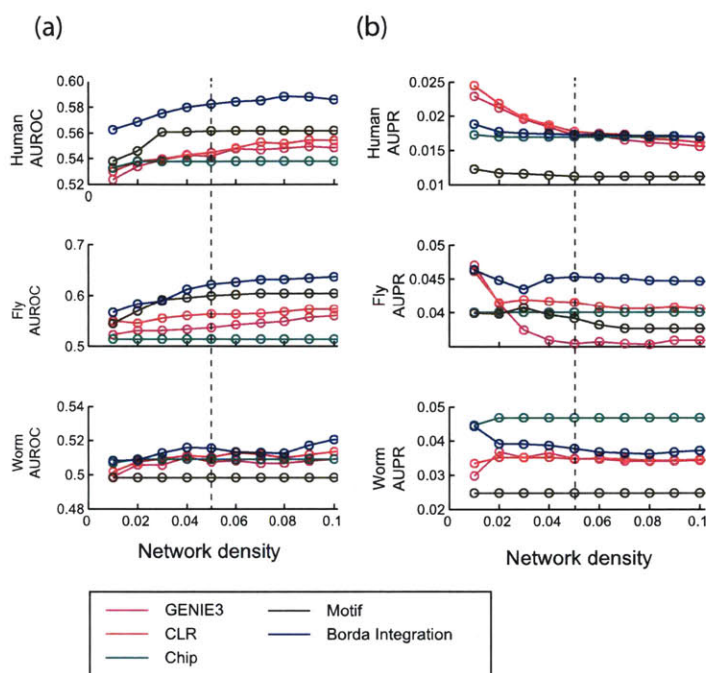


Figure 3-3: AUROC and AUPR scores of feature-specific and integrated regulatory networks in human, fly and worm species.

compared to top ones. Considering both AUROC and AUPR curves for all species, we binarize networks using their top 5% interactions which leads to balanced values of AUROC and AUPR in all inferred networks. This results in $2.6M$ interactions in human, $469k$ in fly and $876k$ in worm. In integrative networks, the median number of targets for each TF is 253 in human, 290 in fly and 640 in worm, with a median of 132 regulators per gene in human, 29 in fly, and 43 in worm.

Chapter 4

Network Deconvolution

Network science has been widely adopted in recent years in diverse settings, including molecular and cell biology [4], social sciences [7], information science [8], document mining [113], and other data mining applications. Networks provide an efficient representation for variable interdependencies, represented as weighted edges between pairs of nodes, with the edge weight typically corresponding to the confidence or the strength of a given relationship. Given a set of observations relating the values that elements of the network take in different conditions, a network structure is typically inferred by computing the pairwise correlation, mutual information, or other similarity metrics between each pair of nodes.

The resulting edges include numerous indirect dependencies due to transitive effects of correlations. For example, if there is a strong dependency between nodes 1 and 2, and between nodes 2 and 3 in the true (direct) network, high correlations will also be visible between nodes 1 and 3 in the observed (direct and indirect) network, thus inferring an edge from node 1 to node 3, even though there is no direct information flow between them (Figure 4-1). Moreover, even if a true relationship exists between a pair of nodes, its strength may be over-estimated due to additional indirect relationships, and distinguishing the convolved direct and indirect contributions is a daunting task. As the size of networks increases, a very large number of indirect edges may be due to second-order, third-order, and higher-order interactions, resulting in diffusion of the information contained in the direct network, and leading to inaccurate network structures and network weights in many applications [4,106,114–119].

Several approaches have been proposed to infer direct dependencies among variables in a network. For example, partial correlations have been used to characterize conditional relationships among small sets of variables [120–122], and probabilistic approaches such as maximum entropy models have been used to identify informative network edges [118,123,124]. Other works use graphical models and message passing algorithms to characterize direct information flows in a network [125] [59], or variations of Granger causality [126] to capture the dynamic relationships among variables [127–129]. Alternative approaches formulated the problem of separating direct from indirect dependencies as a general feature selection problem [100,130,131],

using Bayesian networks [132–134], or using an information theoretic approach to eliminate indirect information flow in the network [135]. These methods are limited to relatively low-order interaction terms, or are computationally very expensive, or are designed for specific applications [136, 137], thus limiting their applicability.

Here our goal is to develop a method for inferring the direct dependencies in a network, corresponding to true interactions, and removing the effects of transitive relationships that result from indirect effects. Our method is based on the following model: when the matrix of direct dependencies is known, all transitive relationships can be computed by summing this direct matrix and all its powers, corresponding to the transitive closure of a weighted adjacency matrix, which convolves all direct and indirect paths at all lengths (Figure 4-1). Given an observed matrix of correlations that contains both direct and indirect effects, our task is to recover the original direct matrix that gave rise to the observed matrix. For a weighted network where edge weights represent the confidence, mutual information, or correlation strength relating two elements in the network, the inverse problem seeks to recognize the fraction of the weight of each edge attributable to direct vs. indirect contributions, rather than to keep or remove unit-weight edges. This inverse problem is dramatically harder than the forward problem of transitive closure, as the original matrix is not known.

We introduce a heuristic algorithm for Network Deconvolution (ND) that can solve the inverse problem of transitive closure of a weighted adjacency matrix, by use of decomposition principles of eigenvectors and eigenvalues, and by exploiting the closed form solution of infinite Taylor series. We demonstrate the effectiveness of this approach and our algorithm in several large-scale networks from different domains and with different properties. First, we seek to distinguish likely direct targets in gene regulatory networks as a post-processing step for diverse gene network inference methods, and show that ND improves both global and local network quality. Second, we show effectiveness of network deconvolution in distinguishing directly-interacting amino-acid residues based on pairwise mutual information data in multi-species protein alignments. Third, we apply ND to a social network setting using a co-authorship network that contains solely connectivity information, and show that

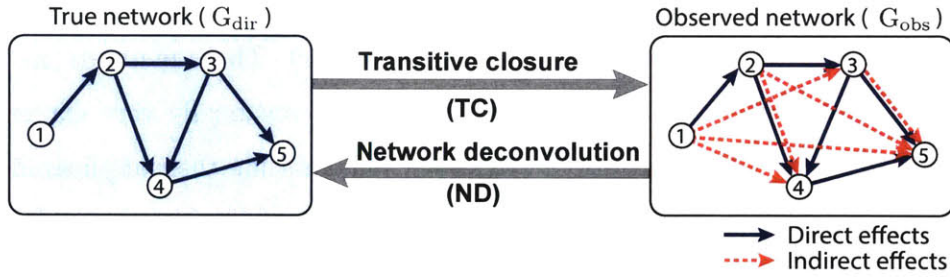


Figure 4-1: Direct edges in a network (solid blue arrows) can lead to indirect relationships (dashed red arrows) due to transitive information flow. These indirect contributions can be of length two (e.g. $1 \rightarrow 2 \rightarrow 3$), three (e.g. $1 \rightarrow 2 \rightarrow 3 \rightarrow 5$) or higher, and can combine both direct and indirect effects (e.g. $2 \rightarrow 4$ and $2 \rightarrow 3 \rightarrow 4$), and multiple indirect effects along varying paths (e.g. $2 \rightarrow 3 \rightarrow 5$, $2 \rightarrow 4 \rightarrow 5$, $2 \rightarrow 3 \rightarrow 4 \rightarrow 5$). Network deconvolution seeks to reverse the effect of transitive information flow across all indirect paths, in order to recover the true direct network (blue edges) based on the observed network (combined blue and red edges).

the resulting edge weights are able to distinguish strong and weak ties independently inferred based on the number of joint papers and additional co-authors.

4.1 Network Deconvolution

4.1.1 Algorithm

Network deconvolution is a systematic approach of computing direct dependencies in a network by use of local edge weights. Let G_{obs} represents an observed dependency matrix a properly scaled similarity matrix among variables. The linear scaling depends on the largest absolute eigenvalue of the un-scaled similarity matrix and is discussed in more details in Section 4.1.3. Components of G_{obs} can be derived by use of different pairwise similarity metrics such as correlation or mutual information. In particular, $g_{i,j}^{obs}$, the (i,j) -th component of G_{obs} , represents the similarity value between the *observed* patterns of variables i and j in the network.

A perennial challenge to inferring networks is that, both direct and indirect dependencies arise together. A direct information flow modeled by an edge in G_{dir} can give rise to two or higher level indirect information flows captured in G_{indir} :

$$G_{indir} = G_{dir}^2 + G_{dir}^3 + \dots \quad (4.1)$$

The power associated with each term in G_{indir} corresponds to the level of indirection contributed by that term. We assume that the observed dependency matrix, G_{obs} , comprises both direct and indirect dependency effects (i.e., $G_{obs} = G_{dir} + G_{indir}$). Further, we assume that G_{dir} is an $n \times n$ decomposable matrix.

The following network deconvolution algorithm finds a solution for direct dependency weights by using the observed dependencies:

Algorithm 4.1. Network deconvolution has three main steps:

- **Linear Scaling Step:** The observed dependency matrix is scaled linearly so that all eigenvalues of the direct dependency matrix are between -1 and 1 .
- **Decomposition Step:** The observed dependency matrix G_{obs} is decomposed to its eigenvalues and eigenvectors such that $G_{obs} = U\Sigma_{obs}U^{-1}$.
- **Deconvolution step:** A diagonal eigenvalue matrix Σ_{dir} is formed whose i -th component is $\lambda_i^{dir} = \frac{\lambda_i^{obs}}{\lambda_i^{obs} + 1}$. Then, the output direct dependency matrix is $G_{dir} = U\Sigma_{dir}U^{-1}$.

Below, we show how the network deconvolution algorithm proposed in Algorithm 4.1 finds a solution for direct dependency weights by using the observed ones.

Suppose U and Σ_{dir} represent eigenvectors and a diagonal matrix of eigenvalues of G_{dir} , where λ_i^{dir} is the i -th diagonal component of the matrix Σ_{dir} . Also, suppose the largest absolute eigenvalue of G_{dir} is strictly smaller than one. This assumption holds by using a linear scaling function over the un-scaled observed dependency network and is discussed in Section 4.1.3. By using the eigen decomposition principle, we have $G_{dir} = U\Sigma_{dir}U^{-1}$. Therefore,

$$\begin{aligned}
 G_{dir} + G_{indir} &\stackrel{(a)}{=} G_{dir} + G_{dir}^2 + \dots & (4.2) \\
 &\stackrel{(b)}{=} (U\Sigma_{dir}U^{-1}) + (U\Sigma_{dir}^2U^{-1}) + \dots \\
 &= U(\Sigma_{dir} + \Sigma_{dir}^2 + \dots)U^{-1} \\
 &= U \begin{pmatrix} \sum_i (\lambda_1^{dir})^i & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sum_i (\lambda_n^{dir})^i \end{pmatrix} U^{-1} \\
 &\stackrel{(c)}{=} U \begin{pmatrix} \frac{\lambda_1^{dir}}{1-\lambda_1^{dir}} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{\lambda_n^{dir}}{1-\lambda_n^{dir}} \end{pmatrix} U^{-1}.
 \end{aligned}$$

Equality (a) follows from the definition of diffusion model of equation (4.1). Equality (b) follows from the eigen decomposition of matrix G_{dir} . Equality (c) uses geometric series to compute the infinite summation in a closed-form since $|\lambda_i^{dir}| < 1$ for all i .

By using the eigen decomposition of the observed network, G_{obs} , we have $G_{obs} = U\Sigma_{obs}U^{-1}$, where

$$\Sigma_{obs} = \begin{pmatrix} \lambda_1^{obs} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_n^{obs} \end{pmatrix}. \quad (4.3)$$

Therefore, from equations (4.2) and (4.3), if

$$\frac{\lambda_i^{dir}}{1-\lambda_i^{dir}} = \lambda_i^{obs} \quad \forall 1 \leq i \leq n, \quad (4.4)$$

the error term $G_{obs} - (G_{dir} + G_{indir}) = 0$. Rewriting equation (4.4) leads to a non-linear filter over eigenvalues in the network deconvolution algorithm:

$$\lambda_i^{dir} = \frac{\lambda_i^{obs}}{1 + \lambda_i^{obs}} \quad \forall 1 \leq i \leq n. \quad (4.5)$$

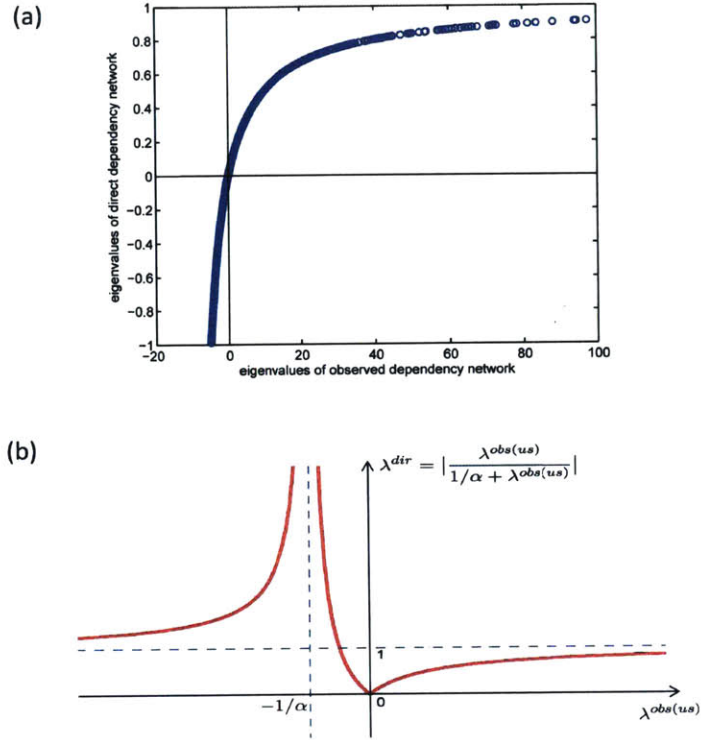


Figure 4-2: Panel (a) shows an example of how eigenvalues of the observed matrix change by application of network deconvolution. In general, ND decreases magnitudes of high positive eigenvalues that were inflated due to transitivity effects over the network. Panel (b) shows the network deconvolution nonlinear mapping function that is applied to eigenvalues of the unscaled observed dependency matrix to compute the eigenvalues of the direct dependency matrix.

This nonlinear filter over eigenvalues allows network deconvolution to eliminate indirect dependency effects modeled as in equation (4.1) and compute a globally optimal solution for direct dependencies over the network with zero error.

4.1.2 Modeling assumptions and intuitions

Network deconvolution can be viewed as a nonlinear filter that is applied on eigenvalues of the observed dependency matrix. Figure 4-2-a shows an example of how eigenvalues of the observed matrix change by application of network deconvolution. In general, ND decreases magnitudes of large positive eigenvalues of the observed matrix

since transitivity effects modeled as in equation (4.1) inflate these positive eigenvalues. Network deconvolution reverses this effect by using the described nonlinear filter.

There are some modeling assumptions and limitations of network deconvolution that we explain below:

- Network deconvolution assumes that networks are *linear time-invariant flow-preserving* operators which excludes nonlinear dynamic networks as well as hidden sinks or sources. Under this condition, indirect flows over different paths are modeled by using power series of the direct matrix, while observed dependencies are modeled as the sum of direct and indirect effects.
- We assume that the maximum absolute eigenvalue of the direct dependency matrix is strictly smaller than one. This assumption can be satisfied by linearly scaling the observed dependency matrix as we explained in Section 4.1.3.

4.1.3 Scaling effects

As discussed in Section 4.1.1, to have convergence of the right-hand side of equation (4.1), the largest absolute eigenvalue of the direct dependency matrix G_{dir} (say β) should be strictly smaller than one (i.e., $\beta < 1$). However, this matrix is unknown in advance and in fact, the whole purpose of having the described model in equation (4.1) is to compute the direct dependency matrix. In this section, we show that by linearly scaling the un-scaled observed dependency matrix G_{obs}^{us} , it can be guaranteed that the largest absolute eigenvalue of the direct dependency matrix is strictly less than one. In the following, we describe how the linear scaling factor α is chosen where $G_{obs} = \alpha G_{obs}^{us}$.

Suppose $\lambda_+^{obs(us)}$ and $\lambda_-^{obs(us)}$ are the largest positive and smallest negative eigenvalues of G_{obs}^{us} . Then, by having

$$\alpha \leq \min \left(\frac{\beta}{(1 - \beta)\lambda_+^{obs(us)}}, \frac{-\beta}{(1 + \beta)\lambda_-^{obs(us)}} \right), \quad (4.6)$$

the largest absolute eigenvalue of G_{dir} will be less or equal to $\beta < 1$.

In the following, we show how inequality (4.6) is derived. Say λ^{dir} and $\lambda^{obs(us)}$ are eigenvalues of direct and unscaled observation matrices. By a similar argument as the one of equation (4.2), we have

$$\lambda^{dir} = \frac{\lambda^{obs(us)}}{\frac{1}{\alpha} + \lambda^{obs(us)}}. \quad (4.7)$$

This function is depicted in Figure 4-2-b. We consider three cases: $\lambda^{obs(us)} \geq 0$, $-1/\alpha < \lambda^{obs} < 0$ and $\lambda^{obs(us)} < -1/\alpha$ (in the case of $\lambda^{obs(us)} = -1/\alpha$, the function is undefined. Hence, we choose α to avoid this case).

- **Case 1:** when $\lambda^{obs(us)} \geq 0$: In this case, we have

$$\frac{\lambda^{obs(us)}}{\frac{1}{\alpha} + \lambda^{obs(us)}} \leq \beta < 1 \Rightarrow \alpha \leq \frac{\beta}{(1 - \beta)\lambda^{obs(us)}}$$

- **Case 2:** when $-1/\alpha < \lambda^{obs} < 0$: In this case, we have

$$\frac{-\lambda^{obs(us)}}{\frac{1}{\alpha} + \lambda^{obs(us)}} \leq \beta < 1 \Rightarrow \alpha \leq \frac{-\beta}{(1 + \beta)\lambda^{obs(us)}}$$

- **Case 3:** when $\lambda^{obs(us)} < -1/\alpha$:

In this case, we have

$$\frac{\lambda^{obs(us)}}{\frac{1}{\alpha} + \lambda^{obs(us)}} \leq \beta < 1 \Rightarrow \frac{1}{\alpha} \leq \frac{1 - \beta}{\beta} \lambda^{obs(us)} < 0,$$

which is not possible since $\alpha > 0$. Therefore, α should be chosen so that eigenvalues of the unscaled observation matrix do not fall in this regime. In other words, for negative $\lambda^{obs(us)}$, $\alpha < \frac{-1}{\lambda^{obs(us)}}$. This condition trivially holds as a consequence of Case 2.

Putting all three cases together, inequality (4.6) is derived which guarantees the largest absolute eigenvalue of the direct dependency matrix is less than or equal to $0 < \beta < 1$.

In equation (4.1), diffusion effects decay with respect to the indirect path length. For example, second order indirect effects decay proportional to β^2 where β is the largest absolute eigenvalue of G_{dir} . Therefore, smaller β means faster convergence of the right-hand side of equation (4.1) and in turn faster decay of diffusion effects. In other words, if β is very small, higher order interactions play insignificant roles in observed dependencies since they decay proportionally to β^k where k is the order of indirect interaction.

Figure 4-3 demonstrates the effect of β in the performance of three different applications: inferring gene regulatory network, inferring protein evolutionary constraints and inferring weak and strong social ties. As it is illustrated in this figure, choosing β close to one (i.e., considering higher order indirect interactions) leads to high performances in these applications. Note that, β has higher effects in the gene regulatory networks and protein evolutionary constraint networks applications compared to the co-authorship network application, which means low-order indirect interactions are sufficient to characterize indirect flows in the co-authorship network compared to two other applications. This is mainly due to the special structure that we observe in the co-authorship network: The co-authorship network can be decomposed locally into densely connected subgraphs and thus most transitive edges happen locally within each module, leading to small indirect effects for longer paths, and thus small β values seem sufficient to characterize their indirect flows. Later in this chapter, we will explain different approaches to choose this parameter for each application.

4.1.4 Asymmetric decomposable matrices

The network deconvolution framework described in Algorithm 4.1 is based on eigen decomposition principles and therefore only can be applied on so called *diagonalizable* matrices. A matrix G is diagonalizable if it can be decomposed to its eigenvectors U and a diagonal matrix of eigenvalues Σ such that $G = U\Sigma U^{-1}$. Many symmetric

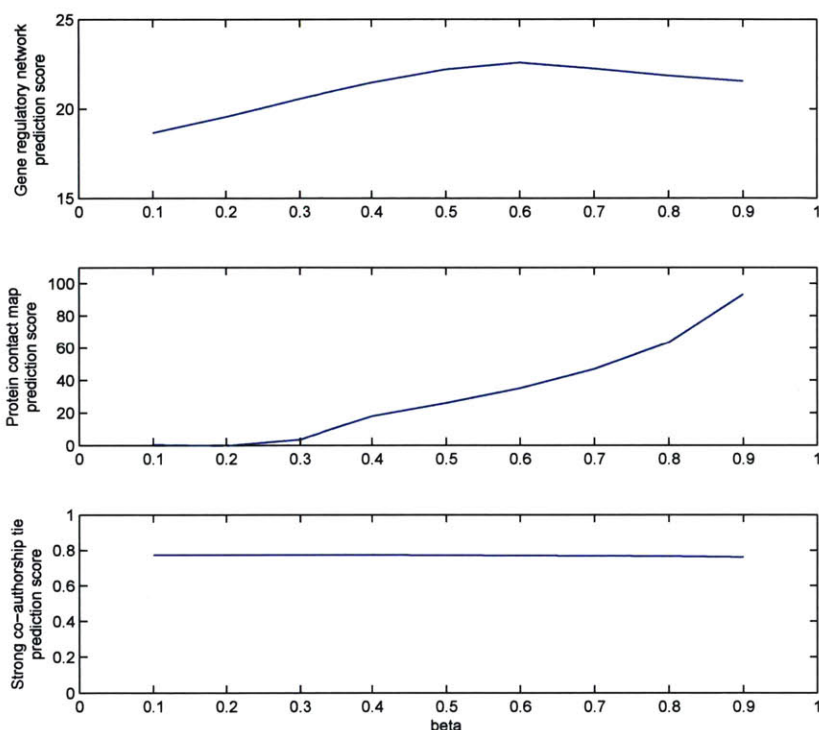


Figure 4-3: This figure illustrates the effect of β in the performance of network deconvolution in various considered applications: Panel (a) demonstrates the overall score of network deconvolution as a gene regulatory network inference method, over three considered networks *E. coli*, *in silico*, and *S.cerevisiae*. This overall score is computed by combining AUPR and AUROC scores as explained in Section 4.2. Panel (b) demonstrates the effect of β in the performance of network deconvolution in inferring protein evolutionary constraints. Y-axis shows the overall contact map prediction rate gains of top 100 predictions over all 15 tested proteins as explained in Section 4.3. Panel (c) illustrates the effect of β in the performance of network deconvolution in inferring weak and string social ties. Y-axis shows the strong tie prediction ratio over the considered co-authorship network as explained in Section 4.4.

matrices and some asymmetric matrices are diagonalizable. Here, we present an example of an asymmetric matrix whose structure is used in the application of network deconvolution in inferring gene regulatory networks where the observed dependency matrix is asymmetric.

Suppose G is a $n \times n$ diagonalizable matrix with the following structure:

$$G = \left[\begin{array}{c|c} G_1 & G_2 \\ \hline 0 & 0 \end{array} \right], \quad (4.8)$$

where G_1 is a $m \times m$ symmetric matrix and G_2 is a $m \times (n - m)$ matrix which can be asymmetric. We show that, eigenvalues of G are real.

Suppose λ is an eigenvalue of G which corresponds to the eigenvector v . By definition, $Gv = \lambda v$. Say

$$v = [v_1 \mid v_2]',$$

where v_1 contains the first m elements of v . Therefore,

$$Gv = \lambda v \Rightarrow \left[\begin{array}{c} G_1 v_1 + G_2 v_2 \\ 0 \end{array} \right] = \lambda \left[\begin{array}{c} v_1 \\ v_2 \end{array} \right]. \quad (4.9)$$

From equation (4.9), we have $v_2 = 0$ and $G_1 v_1 = \lambda v_1$. In other words, λ is also an eigenvalue for the matrix G_1 . Since the matrix G_1 is symmetric, λ is real.

This example illustrates a matrix structure which is common in gene regulatory network inference problems. In that application, the first m rows correspond to transcription factor genes and hence they can have outgoing and incoming edges. The last $n - m$ rows correspond to genes that are not transcription factors and therefore only have in-coming edges. This example illustrates that the proposed network deconvolution algorithm based on eigen decomposition principles can be applied in gene regulatory network inference problems.

4.1.5 Computational complexity analysis of network deconvolution

The proposed network deconvolution algorithm based on eigen decomposition requires computations of eigen decomposition of the observed dependency matrix. This step has order of $\mathcal{O}(n^3)$ computational complexity where n is the number of nodes of the network. Further, the gradient descent based algorithm is iterative and can be very costly to use. Therefore, using baseline network deconvolution algorithms for very large networks can be unattractive due to this required computational power. Here, we propose two extensions to network deconvolution for very large networks. First, many of these large networks are usually low rank and have many zero eigenvalues and therefore the eigen decomposition step can be performed efficiently by exploiting this sparsity. Second, network deconvolution can be parallelized by being applied on local subgraphs of the network. In this section, we explain these two algorithms in more detail.

To show the effectiveness of these extensions, we generate simulated scale-free networks with various levels of eigen sparsity and number of nodes. We add some transitive edges to the network probabilistically by only considering two-hop and three-hop indirect paths. Then, we apply ND to discover transitive edges.

Using sparsity in eigen decomposition

In many applications, networks are low rank (many of network eigenvalues are zero or negligible.). In those cases, eigen decomposition can be performed efficiently by only computing non-zero eigenvalues [138]. Here, we use MATLAB's sparse eigen decomposition algorithm [139,140]. Figure 4-4-a shows CPU time usage of performing network deconvolution on networks with 1000 to 50000 nodes and various eigen sparsity levels. As it is illustrated in this figure, network deconvolution can be performed efficiently even on very large networks if the network is low rank. Figure 4-4-b shows the CPU time usage of ND for various network sizes averaged over networks whose eigen sparsity is less than or equal to 100. This figure demonstrates that unlike the

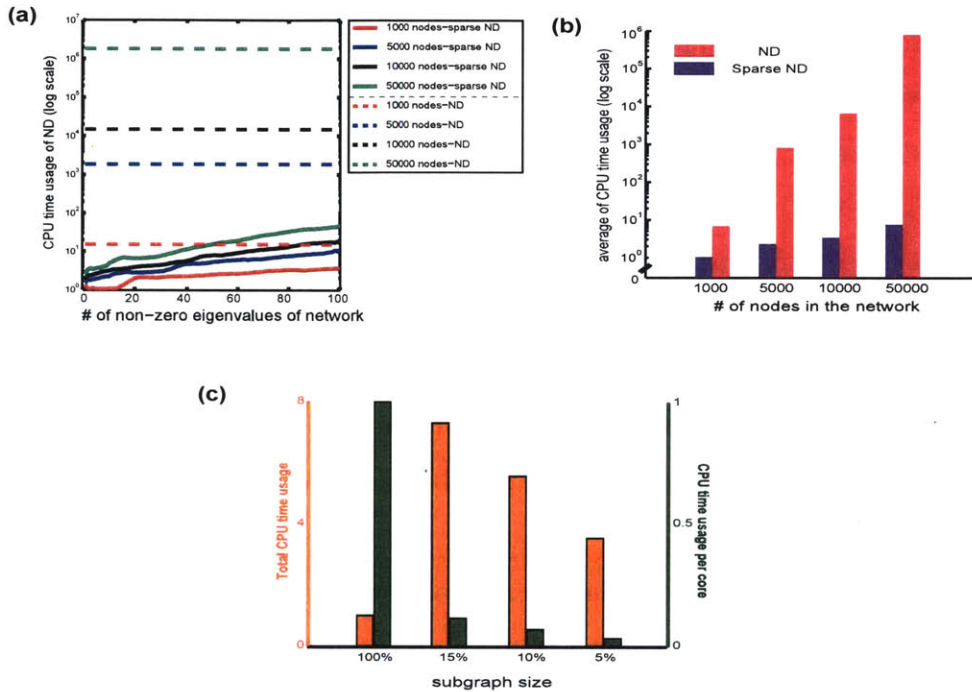


Figure 4-4: Panel (a) shows CPU time usage of performing network deconvolution on networks with 1000 to 50000 nodes and various eigen sparsity levels. As it is illustrated in this figure, network deconvolution can be performed efficiently even on very large networks if the network is low rank. Panel (b) shows the CPU time usage of ND for various network sizes averaged over networks whose eigen sparsity is less than or equal to 100. This figure demonstrates that unlike the baseline network deconvolution algorithm which has a complexity of order $\mathcal{O}(n^3)$, complexity of performing network deconvolution on low rank networks is approximately linear with respect to the number of nodes in the network. (c) Parallelization of network deconvolution is important to apply it effectively on very large networks. This panel demonstrates the total CPU time usage and also CPU time usage per core (per subgraph) for different subgraph sizes of a network with 5000 nodes.

baseline network deconvolution algorithm which has order $\mathcal{O}(n^3)$ complexity, the complexity of performing network deconvolution on low rank networks is approximately linear with respect to the number of nodes in the network.

Parallelization of network deconvolution

Parallelization of algorithms that require large computational power is important in practice. Here, we explain how network deconvolution can be performed in parallel for very large networks:

First, network is divided into subgraphs (perhaps overlapping) with certain sizes. For each subgraph, first and second neighbors of nodes that are not in that subgraph are called marginal nodes of that subgraph. To make sure that edges that are close to subgraph borders do not suffer from network partitioning, we add marginal nodes to each subgraph and then perform network deconvolution. Then, we keep direct dependencies only for subgraph edges (and not for marginal ones). To merge results of different subgraphs, if edges are overlapping, we take the mean of their direct dependencies obtained from different subgraphs. Figure 4-4-c demonstrates the total CPU time usage and also CPU time usage per core (per subgraph) for different subgraph sizes of a network with 5000 nodes.

We next apply our network deconvolution approach to three settings of (1) inferring gene regulatory networks, (2) inferring protein structural constraints, and (3) inferring weak and strong ties in social networks.

4.2 Inferring Gene Regulatory Networks by Network Deconvolution

We first apply our network deconvolution algorithm to gene regulatory networks, which are pervasively used in molecular biology to describe regulatory relationships between transcription factors (regulators) and their target genes [4]. Regulatory network inference from high-throughput gene expression data [4, 106, 107, 114], or by

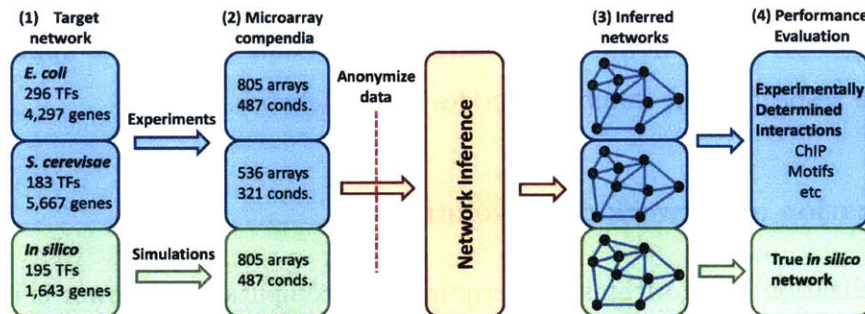


Figure 4-5: DREAM5 participants were solicited to infer gene regulatory interactions from three different expression compendia. For each network, a list of directed, unsigned edges are predicted, ordered according to confidence scores. To score each of these predicted networks, organism specific gold standards containing the known transcription factor to target gene interactions (true positives) are used as compiled in the DREAM project [3]. These gold standards are not available in the inference part. For the evaluation, all transcription factor-target gene pairs that are not part of the gold standards are considered as negatives.

integrating complementary types of datasets [93, 104, 141], is a well-studied problem in computational molecular biology [142, 143], enabling us to benefit from available datasets and community efforts for direct method comparisons [4, 114]. Perhaps the largest such comparison is the recently published network inference challenge part of the Dialogue on Reverse Engineering Assessment and Methods (DREAM) project [3].

4.2.1 DREAM framework

We use datasets and the framework from the DREAM5 project [3] to evaluate the performance of the proposed network deconvolution algorithm as a regulatory inference method. Inferring genome-scale transcriptional regulatory networks from gene expression microarray compendia for a prokaryotic model organism (*E. coli*), a eukaryotic model organism (*S. cerevisiae*) and an *in silico* benchmark is considered in our evaluations. Each of these compendia is represented as an expression matrix of genes and their chip measurements. Further, a list of candidate transcription factors for each compendium is provided. These datasets are fully anonymized (Figure 4-5).

For each network, a list of directed, unsigned edges are predicted, ordered ac-

ording to confidence scores. The confidence scores are only to order predicted edges and are not used directly in the scoring. To score each of these predicted networks, organism specific gold standards containing the known transcription factor to target gene interactions (true positives) are used as compiled in the DREAM project [3]. These gold standards are not available in the inference part. For the evaluation, all transcription factor-target gene pairs that are not part of the gold standards are considered as negatives, although, as the gold standards are based on incomplete knowledge, they might contain yet unknown true interactions.

Two considered biological networks (*E. coli* and *S. cerevisiae*) have 4297 and 5667 genes where 296 and 183 of them are transcription factors, respectively. The *in silico* network has 1643 genes where 195 of them are transcription factors. For each of these networks, gene expressions are reported under several conditions such as time courses, perturbations, knock-outs, etc.

It has been shown that, preprocessing the data by assigning different weights to various expression data types can enhance the quality of predicted networks [3]. However, we avoid these preprocessing to demonstrate the gain obtained by network deconvolution algorithm alone.

A key assumption in inferring gene regulatory networks from gene expression data is that, mRNA levels of regulators and their targets show some mutual dependencies. Therefore, the dependency content within interactions of a network can be used as a coarse estimate of inference difficulty. *E. coli* and *in silico* networks demonstrate higher information contents over interacting pairs compared to non-interacting ones [3]. This is true regardless whether Pearson's correlation or mutual information is used to measure dependencies. However, in *S. cerevisiae* network, the dependency distribution of interacting and non-interacting pairs are almost identical. This shows it is more difficult to infer the regulatory network from gene expression data. This fact is reflected in low scores of various inference methods for this network.

We use the same performance evaluation metrics as in the DREAM5 challenge [3]. Network predictions are evaluated as binary classification tasks where edges are predicted to be present or absent. Then, standard performance metrics from machine

learning are used: precision-recall (PR) and receiver operating characteristic (ROC) curves. Then, predictions are evaluated statistically and an overall score summarizing the performance over all three networks is computed. Note that, similar to DREAM5 framework, for each network only top 100,000 edges are considered in the performance evaluation.

For each predicted network, the performance was assessed by the area under the ROC curve (AUROC, true positive rate vs. false positive rate), and the area under the precision vs. recall curve (AUPR). Expressions for true positive rate (TPR), false positive rate (FPR), precision and recall as a function of the cutoff (k) in the edge list are as follows [3]:

$$\begin{aligned} \text{recall}(k) &= \frac{TP(k)}{P}, \\ \text{precision}(k) &= \frac{TP(k)}{TP(k) + FP(k)} = \frac{TP(k)}{k}, \end{aligned}$$

where $TP(k)$ and $FP(k)$ are the number of true positives and false positives in the top k predictions, and P is the number of positives in the gold standard. The *true positive rate*, $TPR(k)$, is similar to *recall*. The *false positive rate* is the fraction of negatives incorrectly predicted in the top k predictions:

$$FPR(k) = \frac{FP(k)}{N}, \tag{4.10}$$

where N is the number of negatives in the gold standard.

AUROC and AUPR values were separately transformed into p -values by simulating a null distribution for a large number (25,000) of random networks [3]. Overall ROC and PR scores are derived by calculating the geometric mean of the network specific p -values:

$$\begin{aligned}\text{ROC score} &= \frac{1}{3} \sum_{i=1}^3 -\log_{10} p_{ROC_i} \\ \text{PR score} &= \frac{1}{3} \sum_{i=1}^3 -\log_{10} p_{PR_i}\end{aligned}$$

where p_{ROC_i} and p_{PR_i} represent p -values of ROC and PR curves for network i . Then, an overall score summarizing the performance over all three networks is calculated as the mean of the ROC and PR scores.

4.2.2 Evaluation Results

The network deconvolution algorithm has two parameters: β which is the eigenvalue scaling parameter as explained in Section 4.1.3, and α which is the density parameter and determines the fraction of edges in the input network. The density parameter α is a preprocessing parameter, set in advance for each application, which represents the expected density of the input network to be deconvolved, namely the fraction of non-zero edges relative to a fully connected network. For the gene regulatory network inference application, we used $\alpha = 0.1$, based on the expectation that regulators target, on average, approximately 10% of all genes, consistent with the observed sparsity of ChIP-seq experiments in diverse species (setting this value to 1 would correspond to each regulator targeting 100% of all genes). From a ranking perspective, this corresponds to relative ranks in the top 10% of regulatory connections being meaningful, but the remaining 90% of ranks being uninformative.

The β parameter represents the rate at which indirect effects are expected to propagate through the network. For small values of beta, indirect effects decay rapidly and higher-order transitive interactions play a smaller role. For large values of beta, indirect effects persist for longer periods and higher-order interactions play a more important role. For this application, we used $\beta = 0.5$, based on the intuition that transitive regulatory effects require rounds of transcription and translation that should delay indirect effects and thus cause them to decay more rapidly. We demonstrate

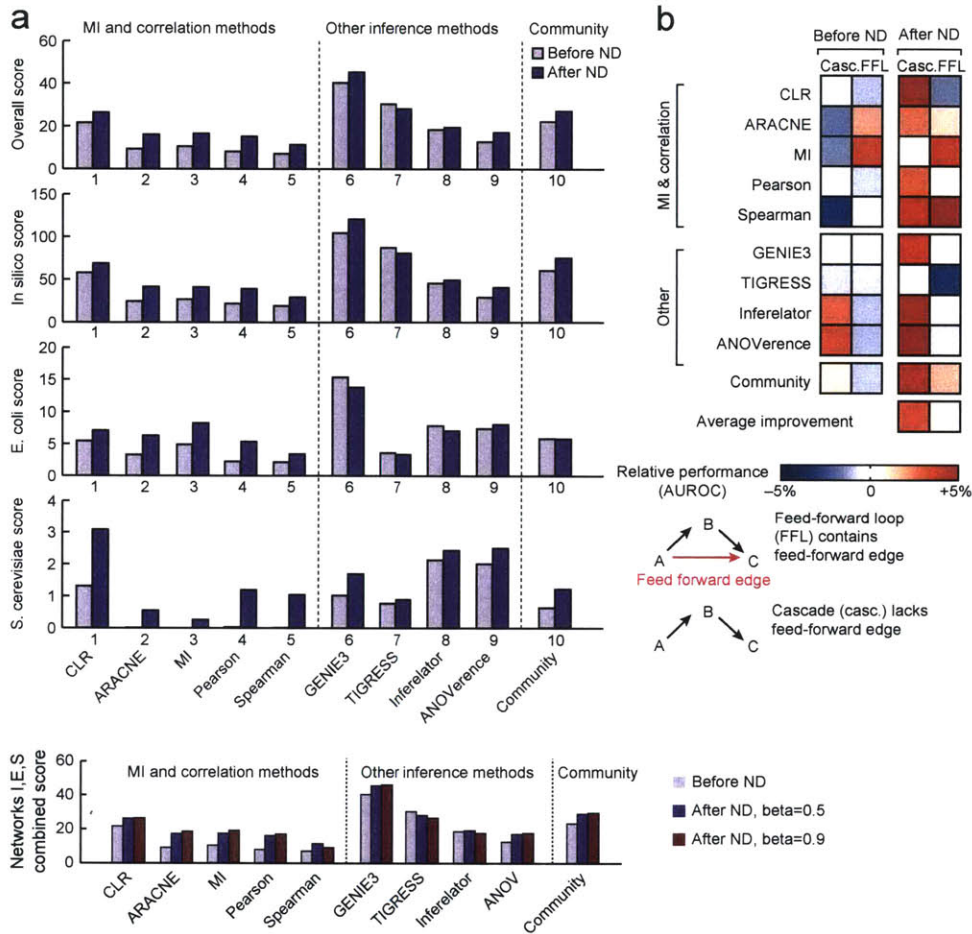


Figure 4-6: (a) ND applied to the inferred networks of top-scoring methods from DREAM5 leads to consistent improvements for mutual information (MI) and correlation based methods. ND also improves other top-scoring methods, including the best performing method of the DREAM5 challenge (GENIE3). Moreover, the community network obtained by integrating network predictions from individual methods (1-9) before ND is outperformed by the community network based on deconvolved networks. (b) Network motif analysis showing the relative performance of inference methods for cascades (casc.) and feed-forward loops (FFL) before and after ND. Red/blue corresponds to increased/decreased prediction accuracy of the two motif types relative to the overall performance of the method before ND (measured by the area under the ROC curve, AUROC). The original methods (before ND, left side) have different relative performances for cascades and FFLs, e.g., mutual information-based network inference (MI) tends to include feed-forward edges (red arrow), resulting in higher accuracy for FFLs but lower accuracy for cascades, while the opposite is true for the Inferelator and ANOVereance. The deconvolved networks (after ND, right side) show significantly higher accuracy for true cascade network motifs for all methods, and moderately improved accuracy for FFLs on average.

results for other values of these parameters along with a cross validation analysis in Section 4.2.3.

In the practical implementation of network deconvolution (ND), matrices are made symmetric, although this was not needed in our three examples as the input matrices were already symmetric (except some cases in regulatory inference application). Network deconvolution provides direct weights for both observed and non-observed interactions, but their use depends on the application. In this chapter, we focused on weights for the subset of observed edges, as our goal was to deconvolve indirect effects on the existing interactions. For gene networks we used the top 100,000 edges. For asymmetric networks that are not diagonalizable, a small random matrix (components smaller than epsilon) is added to the input matrix, to eliminate perfectly correlated subspaces.

The DREAM5 challenge provides an ideal benchmark for evaluating ND, given the uniform benchmarks for network reconstruction used, and the participation of many of the research teams at the forefront of network inference research, with a total of 35 different prediction methods applied across a wide array of methodologies.

Given that ND is designed as a way to eliminate indirect edge weights in mutual information based and correlation based networks, we first applied it to the networks predicted by the top-scoring such methods, including CLR [99], ARACNE [135] and basic mutual information (Relevance networks) [144]. In all cases, we found that network deconvolution substantially improved the performance of each method according to all metrics used and for all networks tested in DREAM5 (Figure 4-6-a). It is notable that ARACNE, which seeks to remove transitive edges by studying feed-forward loops directly, showed a significant improvement by network deconvolution, indicating that these indirect effects are not always detectable at the local level but instead require a global network deconvolution approach. As information theoretic methods are among the most widely-used network inference approaches [3], their use in combination with ND can be of great general use.

We next applied ND to other top performing inference methods that are not based on mutual information or correlation. These include ANOVarence [145] that

uses a non-parametric non-linear similarity metric between transcription factors and target genes, GENIE3 [146] that uses regression and a tree-based ensemble method, TIGRESS [147] that uses a sparse regression formulation and feature selection, and Inferelator [107] that uses regression and variable selection based on expression data. We found that network deconvolution was effective even when applied to these methods. The performance was increased for three of the four methods, including for the top performing method (GENIE3). As GENIE3 was the overall top-performing method, this suggests that the combination of GENIE3 and ND provides the new top-performing method [3]. We also applied ND in combination with the community prediction method from DREAM5. We found that community prediction after ND showed greater performance than community prediction on the original networks, suggesting that network deconvolution maintains the complementary aspects of these networks important in community prediction approaches. We note that the community prediction approach is not the best predictor here, with or without ND, likely due to the insufficiently diverse nature of the original networks. Overall, these results suggest that despite the ability of even the best-performing methods to recover high-quality networks, strong indirect effects remain, which can be reduced by use of ND.

We next studied how ND affects the prediction of local network connectivity patterns. We specifically focused on the ability to correctly predict feed-forward loops, that truly contain both an indirect path and a feed-forward edge, and regulatory cascades, for which A and C are only connected through B. We performed a network motif analysis [3, 106] to compare the ability of inference methods to correctly predict feed-forward loops and regulatory cascades before and after ND (Figure 4-6-b). We used the refined approach introduced by Petri et al. [148], which measures performance for different motif types based on the area under the ROC curve (AUROC) [3]. Briefly, this involves: (1) identifying the subset of edges pertaining to a given motif type in the true network; (2) evaluating the performance for this subset of edges using the AUROC. Note that edges are counted only once for a given motif type, even if they are part of several overlapping motif instances of the same type. Using this approach,

we computed for each network inference method the motif-specific AUROC of regulatory cascades and feed-forward loops. The prediction bias is given by the difference between the overall AUROC of a given inference method (computed over the complete set of edges as described in Methods) and the motif-specific AUROC. Figure 4-6-b shows the average prediction biases over the *in silico* and *E. coli* networks (following the DREAM project [3], we excluded *S. cerevisiae* because the inference methods did not recover sufficient true edges for a meaningful network motif analysis).

Consistent with previous studies [3], we found that network inference methods tend to perform better on one or the other network motif, based on their approach for dealing with indirect information (Figure 4-6-b). For example, mutual information-based network inference (MI) is biased towards including feed-forward edges, leading to increased accuracy for feed-forward loops but many spurious transitive edges for cascades, while the Inferelator and ANOverence are biased towards excluding feed-forward edges, leading to increased accuracy for cascades but many missing feed-forward edges in feed-forward loops. Interestingly, the ARACNE algorithm that seeks to directly remove transitive edges shows a decreased performance for FFL relative to MI, highlighting the difficulty of distinguishing transitive edges from true feed-forward edges. If ND can accurately identify spurious indirect edges but preserve true feed-forward edges, we should expect significantly increased accuracy for cascades, and no decrease in accuracy for feed-forward edges. Indeed, we found that deconvolved networks lead to improved prediction accuracy for true cascades for each method, thus correctly eliminating spurious edges. Importantly, the improved performance on cascades did not lead to an increased error rate on feed-forward loops, where prediction accuracy remained similar or improved in most deconvolved networks, with the exception of TIGRESS, which was also the only method where ND did not lead to an improved overall performance. Taken together, these results show that ND effectively distinguishes direct from indirect edges, improving the predictions of a wide range of gene regulatory network inference approaches.

$\beta = 0.5$				$\beta = 0.9$			
	Overall	MI/Corr.	Other		Overall	MI/Corr.	Other
$\alpha = 0.1$	142%	166%	111%	$\alpha = 0.1$	140%	166%	108%
$\alpha = 0.5$	132%	153%	107%	$\alpha = 0.5$	132%	157%	100%
$\alpha = 1.0$	127%	144%	106%	$\alpha = 1.0$	146%	182%	101%

Table 4.1: Effect of varying both density (α) and eigenvalue scaling (β) parameters on the DREAM5 networks.

4.2.3 Effect of Varying ND Parameters

To understand the effect of varying both density and scaling, we evaluated the average performance of ND relative to other DREAM5 methods for $\alpha = 0.1, 0.5, 1$, and $\beta = 0.5, 0.9$ eigenvalue scaling in Table 4.1. We found that for all parameter combinations, ND led to improved average performance relative to the other DREAM5 methods (127%-146% relative performance on average). The improvement was consistently higher for Mutual Information (MI) and Correlation (Corr) based methods, which were the primary motivation behind ND, indicative of their abundant indirect effects. The highest relative performance overall (146% relative performance) was obtained for the combination of parameters closest to those used in the protein folding and social network applications ($\alpha = 1, \beta = 0.9$), but all parameter combinations tested led to similar performance improvements.

In practice, one can choose these parameters using cross-validation, by separating the data in a training and a testing set, using the training set to select an appropriate set of parameters, and evaluating performance in the test set using those selected parameters [149]. Figure 4-7 illustrates a cross validation analysis to choose the eigenvalue scaling parameter β in the application of ND on DREAM5 networks (for $\alpha = 0.1$). We train this parameter on 1/3 of input networks and evaluate the performance on the remaining 2/3. Experiments have been repeated 100 times.

4.3 Inferring Protein Structural Contacts

We next applied ND to infer structural constraints between pairs of amino-acids for protein structure prediction [150–152]. Prior work used evolutionary information to

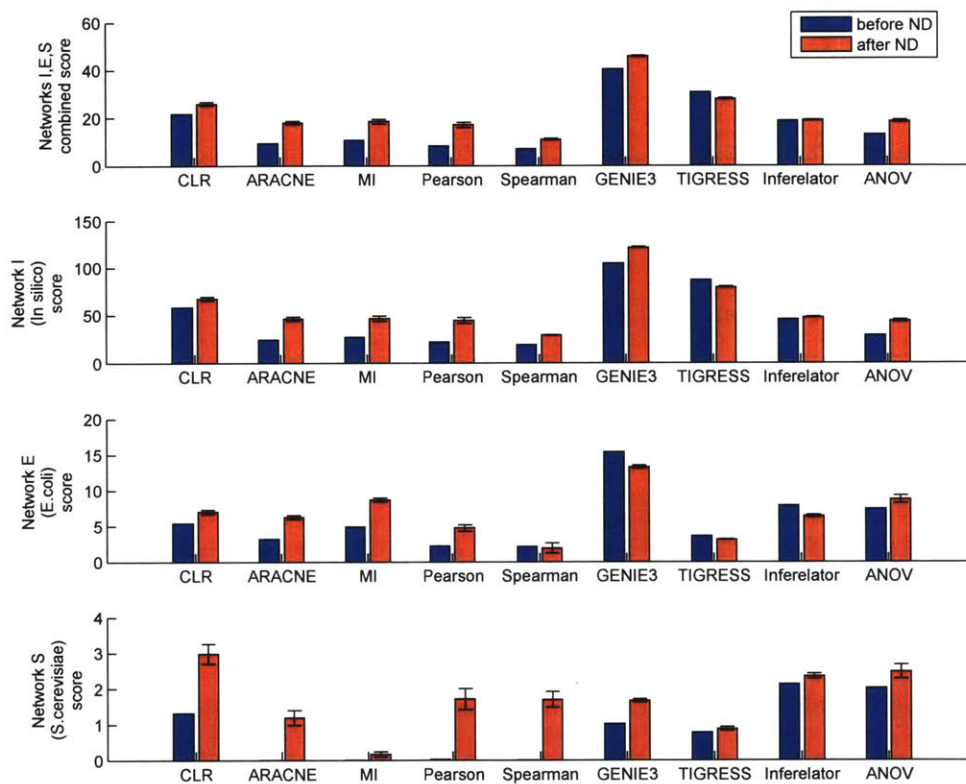


Figure 4-7: A cross validation analysis to choose the eigenvalue scaling parameter β in the application of ND on DREAM5 networks. We train this parameter on 1/3 of input networks and evaluate the performance on the remaining 2/3 of networks. Experiments have been repeated 100 times.

Protein PDB ID	Protein name	Protein PDB ID	Protein name	Protein PDB ID	Protein name
1hzx	OPSD BOVIN	1bkr	SPTB2 HUMAN	1odd	OMPR ECOLI
3tgi	TRY2 RAT	2it6	A8MVQ9 HUMAN	1g2e	ELAV4 HUMAN
5p2l	RASH HUMAN	1rqm	THIO ALIAC	1wvn	PCBP1 HUMAN
1f2l	RNH ECOLI	2o72	CADH1 HUMAN	5pti	BPT1 BOVIN
1e6k	CHEY ECOLI	1r9h	O45418 CAEEL	2hda	YES HUMAN

Table 4.2: Tested protein names and PDB IDs

reveal pairs of amino acid residues that are proximal in the three-dimensional protein structure. However, the pairwise evolutionary correlation matrix may contain many transitive relationships between pairs of residues [115–118, 153–156]. For example, if two amino-acid residues both interact with an intermediate residue, but are not directly interacting with each other, they will show high mutual information due to indirect effects. One approach to remove transitive noise is to use a probabilistic maximum entropy solution [118] that is specifically designed for inferring directly interacting residues [123, 124]. Our aim here is to demonstrate effectiveness of using network deconvolution as a general method to infer directly interacting residues over protein contact networks.

As strong clusters of high mutual information have been shown to hinder identification of directly-interacting residues, we reasoned that network deconvolution may be able to break up these clusters and reveal directly-interacting residues, by distinguishing those correlations that can be explained by transitive relationships. Here, we build on an approach which uses comparative genomics information of residue co-variation across evolutionarily-diverged species.

We applied network deconvolution to predict contact maps on fifteen proteins in different folding classes with sizes ranging from 50 to 260 residues [123]. The list of these proteins are presented in Table 4.2. In our input network, the nodes represent amino acid residues, and each edge between a pair of residues represents their co-variation across multiple sequence alignments spanning 2,000-72,000 sequences, quantified by their mutual information. Applying ND to a mutual information network leads to a systematic and substantial increase in the discovery rate of interacting amino-acids, based on non-adjacent amino-acid contact maps for known structures (Figures 4-8, 4-9, 4-10). We also applied ND to a weighted interaction network based

on direct information [123]. Although using ND over direct information led to a small but consistent improvement over the top predictions especially for non-redundant interacting pairs, a robust performance assessment requires comparison of predicted proteins 3D structures which is beyond the scope of this manuscript. Below we explain details of these experiments.

4.3.1 Evaluation Results

Mutual information (MI) and direct information [123] (DI), a method recently developed for the protein-structure prediction problem that seeks to explicitly maximally distribute the interaction probabilities across different residue pairs were used to compute co-variation similarities among different residues. We used the full MI and DI matrices computed from the alignments at <http://evofold.org/> using the code provided on the website. Then, network deconvolution is applied on MI and DI matrices.

We say two amino acid pairs *contact* if their CA atom distance is less than 7 Angstrom. Trend of our results is robust to the choice of threshold (Figure 4-11). However, choosing a small threshold leads to sparse contact maps over which significant statistical results cannot be inferred, while considering large thresholds may introduce false contacts. We also report the results for both smaller (5 Angstrom, Figure 4-11-a) and larger (8 Angstrom, Figure 4-11-b) thresholds. Distant contact maps are the contacts whose amino acid indices differ at least by 4 with two modifications. We exclude contacts in alpha-helices between positions i and $i + 8$ when there are already contacts between positions i and $i + 4$, and between $i + 4$ and $i + 8$. We also exclude contacts between positions i and $i + 4$ to remove trivial contacts due to alpha-helices. More elaborate contact refinement schemes can also be used as well [123].

For evaluation, we consider fraction of discovered true positive contacts among top predictions. We also say a discovered contact is *non-redundant* when it cannot be explained by previously discovered ones (i.e., there is no indirect path connecting them over the network). Non-redundant contacts provide important constraints in

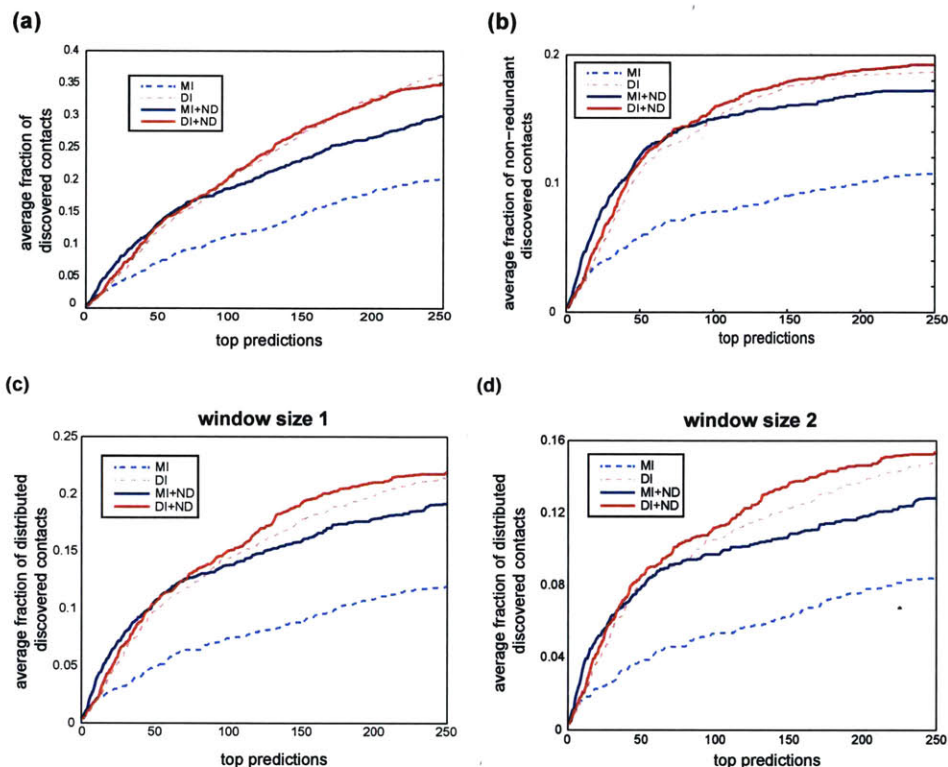


Figure 4-8: Overall total and non-redundant contact discovery rates. (a) Recovery of truly interacting amino acid residues with increasing numbers of top-scoring pairs for mutual information (MI), direct information (DI), and network deconvolution (ND). Recovery after ND is consistently higher than MI for each of tested proteins. ND leads to a small but consistent improvement across the top predictions of direct information (DI). In the strongest-scoring predictions, ranked between 0 and 50 for each method, ND outperforms DI, both independently and as a post-processing step. (b) Discovery of non-redundant contacts, that cannot be explained by previously discovered ones. ND shows improvement over both MI and DI specially in top predictions. (c),(d) We show here an alternate definition of redundant contacts, using a window-based approach (predictions that are within a distance window of previously-discovered contacts are considered redundant). In this metric, ND shows consistent performance increase at all cut-offs.

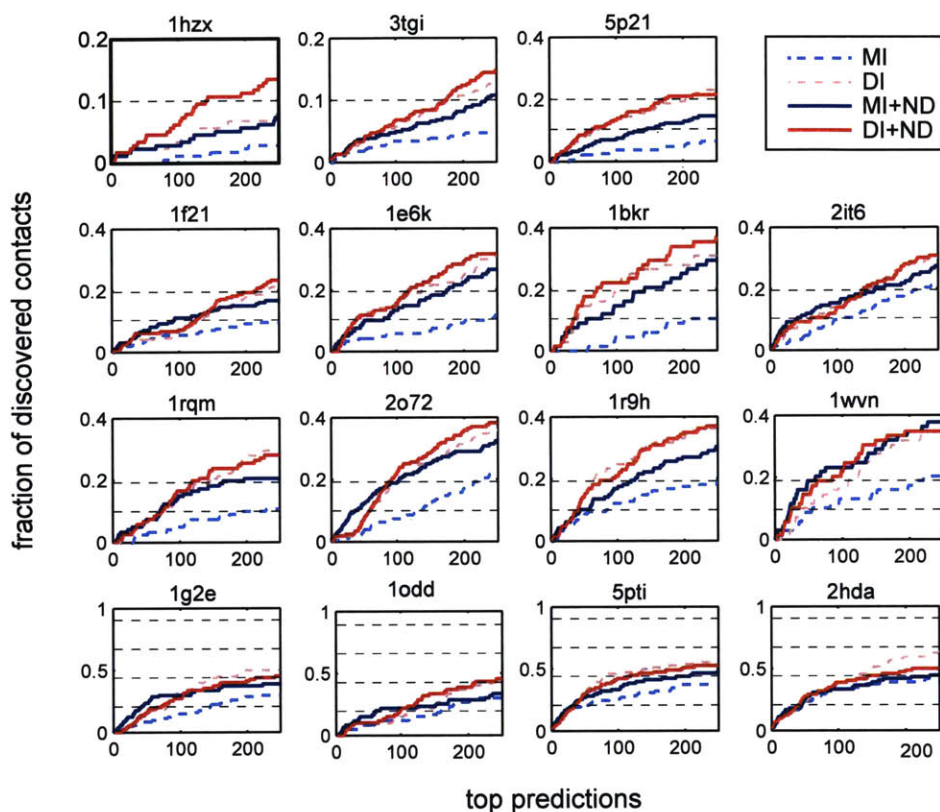


Figure 4-9: Contact recovery rate after ND is consistently higher than MI for each of 15 tested proteins. ND leads to a small but consistent improvement across the top predictions of direct information (DI). Remarkably, in the strongest-scoring predictions, ranked between 0 and 50 for each method, ND outperforms DI, both independently and as a post-processing step. We excluded protein 1hzx (highlighted in the figure) from overall performance comparisons as all methods had very low performance on it possibly due to its poor alignment.

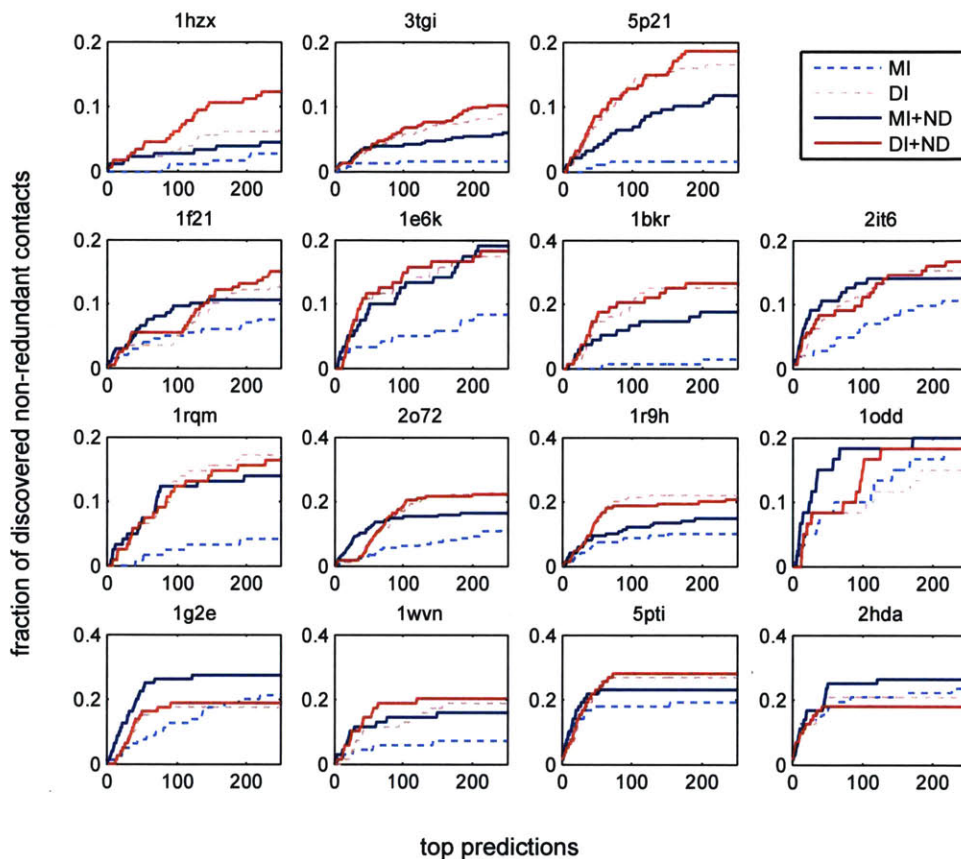


Figure 4-10: Focusing on non-redundant contacts, that cannot be explained by previously discovered ones, ND shows a significant improvement over the performance of MI. The improvement of ND over DI is even more pronounced when considering non-redundant interactions.

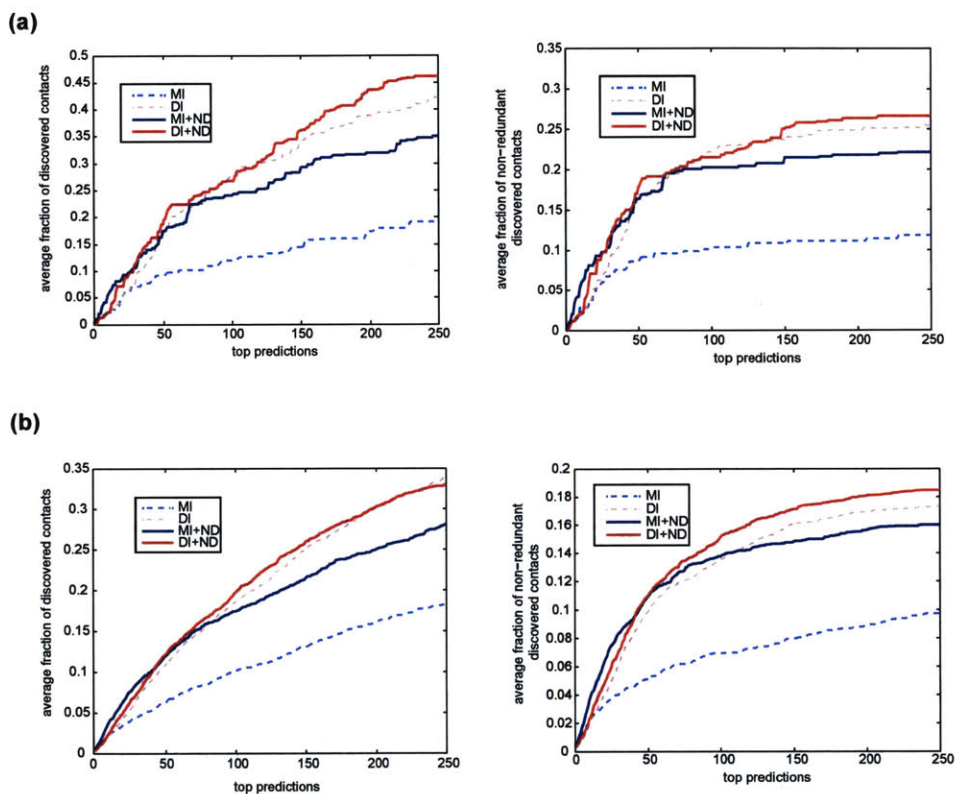


Figure 4-11: Performance trend of ND is robust to the choice of proximity threshold to define contacts. However, choosing a small threshold leads to sparse contact maps over which significant statistical results cannot be inferred, while considering large thresholds may introduce false contacts. We choose threshold 7 Angstrom and also report the results for both smaller (5 Angstrom, part a) and larger (8 Angstrom, part b) thresholds.

predicting 3D structure of proteins.

We compare the performance of network deconvolution (ND) against mutual information (MI) and direct information (DI). For the ND algorithm we use $\alpha = 1$ (i.e., full density networks) and $\beta = 0.99$. Effects of various β parameters on the ND performance have been illustrated in Figure 4-3-b. Figure 4-8 shows average total and non-redundant contact discovery rates for 14 tested proteins. We excluded protein 1hzx as all methods had very low performance on it possibly due to its poor alignment.

Recovery rate of contacts after ND is consistently higher than MI for each of 15 tested proteins, leading to an average increase of 90%, 65% and 50% in the top 50, 100 and 250 predictions. ND leads to a small but consistent improvement across the top predictions of direct information (DI). Remarkably, in the strongest-scoring predictions, ranked between 0 and 50 for each method, ND outperforms DI, both independently and as a post-processing step (Figure 4-8-a and 4-9). Focusing on discovery of non-redundant contacts, that cannot be explained by previously discovered ones, ND shows an even greater improvement over the performance of MI, by 116%, 95%, and 60% in the top 50, 100, and 250 predictions, respectively (4-8-b). Considering non-redundant interactions, ND shows improvement of 9%, 8%, and 6% over DI for the top 50, 100, and 250 predictions, respectively (Figure 4-8-b and Figure 4-10). We also developed an alternate definition of redundant contacts, using a window-based approach (i.e., predictions that are within a distance window of previously-discovered contacts are considered redundant). In this metric, ND shows consistent performance increase at all cut-offs over MI and DI (Figure 4-8-c,d).

4.4 Application to Co-authorship Collaboration Relationships

We next applied our network deconvolution approach to a social network of co-authorship information [157] to distinguish strong and weak collaborations, that can

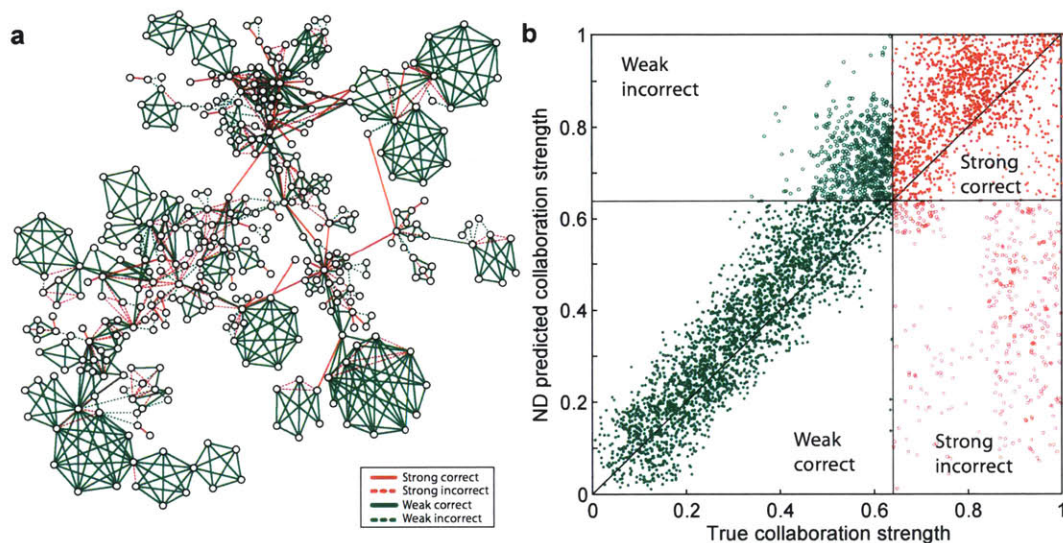


Figure 4-12: (a) Use of network deconvolution to distinguishing strong ties (red) and weak ties (green) in the largest connected component of a co-authorship network containing 379 authors. True collaboration strengths were computed by summing the number of co-authored papers and down-weighting each paper by the number of additional co-authors. ND only had access to unweighted co-authorship edges, but exploiting transitive relationships to weigh down weak ties resulting in accurate predictions (solid lines), demonstrating that this information lies within the network edges, and that ND is well-suited for discovering it. (b) Beyond the binary classification of strong and weak ties, we found a strong correlation across all 2742 edges connecting 1589 authors, between the weights assigned by ND (y-axis) and the true collaboration strengths (x-axis) obtained using additional publication details.

play different key roles in social networks [16, 158, 159]. Given the recent surge of social networks like Facebook or ResearchGate, recognizing weak and strong ties is increasingly important for recommending friends or colleagues, recognizing conflicts of interest, or evaluating an author's contribution to a team. Previous approaches have defined strong ties using shared indirect contacts, edges that increase network distance upon removal, or edges connecting nodes within the same module [158]. In co-authorship networks, strong ties have been defined by using additional information beyond network connectivity, including the number of co-authored papers, and the number of other co-authors of these papers [157].

We used an unweighted input network of 1589 scientists working in the field of network science [157], in which two authors are connected by an edge if they have co-authored at least one paper. We then applied our network deconvolution approach directly on the edges provided by the co-authorship network, to recognize whether network connectivity information alone is sufficient to capture additional information about strong and weak ties previously computed on the same network. Our assumption is that edges resulting from indirect paths likely correspond to weak collaborations, diluted over many other co-authors, while edges with low indirect contributions are more likely to correspond to meaningful collaborations. Application of ND to this unit-weight network led to a weighted network whose transitive closure most closely captures the input network information, and whose weights represent the inferred strength of likely direct interactions. We then ranked all co-authorship edges according to the weight assigned to each by the ND approach.

For the ND algorithm we used $\alpha = 1$ (i.e., full density networks) and $\beta = 0.95$. Effects of various β parameters on the ND performance have been illustrated in Figure 4-3-c. We found that the resulting edge weights indeed capture co-authorship tie strengths previously computed by summing the number of co-authored papers and down-weighting each paper by the number of additional co-authors [157]. We defined true strong ties based on Newman's weight ≥ 0.5 that incorporates additional publication information, and our predictions based on the network deconvolution weight corresponding to the same fraction of edges. We found that network deconvolu-

tion correctly recovered more than 77% of strong co-authorship ties solely by use of the network topology, demonstrating that additional information about collaboration strength lies within network connectivity information, and that ND is very well-suited for discovering it (Figure 4-12-a). Beyond the binary classification of edges into strong and weak, we found a strong overall agreement between the rank obtained by the true collaboration strength and the rank provided by the ND weight. Note that in Figure 4-12-b, tie strengths are mapped to be between 0 and 1 by using their edge ranks. We have added some noise to points for visualization purposes. The exception was a population of edges that had strong collaboration scores but weak ND weights, likely due to the number of co-authored publications that factors in the collaboration score but is not available in the ND network input.

Chapter 5

Network Alignment

The term *network alignment* encompasses several distinct but related problem variants [95]. In general, network alignment aims to find a bijective mapping across two (or more) networks so that if two nodes are connected in one network, their images are also connected in the other network(s). If such an errorless alignment scheme exists, network alignment can be simplified to the problem of graph isomorphism [160]. However, in general, an errorless alignment scheme may not be feasible across two networks. In that case, network alignment aims to find a mapping with the minimum error and/or the maximum overlap.

Network alignment has a broad range of applications in systems biology, social sciences, computer vision, and linguistics. For instance, network alignment has been used frequently as a comparative analysis tool in studying protein-protein interaction networks across different species [161–166]. In computer vision, network alignment has been used in image recognition by matching similar images [167, 168]. It has also been applied in ontology alignment to find relationships among different representations of a database [169, 170], and in user de-anonymization to infer user/sample identifications using similarity between datasets [171].

A network alignment optimization seeks an assignment of nodes and edges across multiple networks to maximize (or alternatively minimize) a cost function with quadratic terms. This problem is closely related to the quadratic assignment problem (QAP) [172] which is computationally challenging to solve exactly. Reference [173] shows that approximating a solution of maximum quadratic assignment problem within a factor better than $2^{\log^{1-\epsilon} n}$ is not feasible in polynomial time in general. However, owing to numerous applications of quadratic assignment problems in different areas, several algorithms have been designed to solve it approximately: some methods use exact search approaches based on branch-and-bound [174] and cutting plane [175]. These methods can only be applied to very small problem instances owing to their high computational complexity. Some methods attempt to solve the underlying QAP by linearizing the quadratic term and transforming the optimization into a mixed integer linear program (MILP) [176–179]. In practice, the very large number of introduced variables and constraints in linearization of the QAP objective function poses

an obstacle for solving the resulting MILP efficiently. Some methods use convex relaxations of the QAP to compute a bound on its optimal value [180–184]. The provided solution by these methods may not be a feasible solution for the original quadratic assignment problem. Other methods to solve the network alignment optimization include semidefinite [184, 185] or non-convex [186] relaxations, Bayesian inference [187], message passing [188] or other heuristics [161, 162, 164, 189–193]. We will review these methods in Section 5.1.2. For more details on these methods, we refer readers to references [172, 194].

Spectral inference methods have received significant attention in various network science problems such as network clustering and low dimensional embedding problems [2, 195–198]. However, the use of spectral techniques in the network alignment problem has been limited [161, 162, 189, 190] partially owing to the lack of principled connections between existing spectral network alignment methods and relaxations of the underlying QAP. In fact, the performance of existing spectral network alignment methods have been assessed mostly through limited simulations and/or validations with real data where an analytical performance characterization is lacked even in simple cases.

In this chapter, we propose a network alignment framework that uses an orthogonal relaxation of the underlying QAP in a maximum weight bipartite matching optimization. Our method simplifies, in a principled way, the network alignment optimization to simultaneous alignment of eigenvectors of (transformations of) adjacency graphs scaled by corresponding eigenvalues. We show that our framework not only can be employed to provide a theoretical justification for existing heuristic spectral network alignment methods, but it also leads to a new scalable network alignment algorithm which outperforms existing ones over various synthetic and real networks. We prove that our solution is asymptotically exact with high probability for Erdős-Rényi graphs [199], under some general conditions. Proofs are based on a characterization of eigenvectors of Erdős-Rényi graphs, along with some spectral perturbation analysis. For an analytical performance characterization of our proposed method, we consider asymptotically large Erdős-Rényi graphs owing to their

tractable spectral characterization. Note that finding an isomorphic mapping across asymptotically large Erdős-Rényi graphs is a well studied problem and can be solved efficiently through canonical labeling [200]¹. However, note that in the network alignment problem graphs can be non-isomorphic as well. Also note that in our optimality analysis we only consider finite and asymptotically large graphs. For arguments on infinite graphs, see Section 5.1.3.

Moreover, we generalize the objective function of the network alignment problem to consider both matched and mismatched interactions in a standard QAP formulation. In general, existing scalable network alignment methods only consider maximizing the number of overlapping edges (matches) across two networks and ignore the number of resulting mismatches (interactions that exist only in one of the networks). This limitation can be critical particularly in applications where networks have different sizes and low similarity. Through analytical performance characterization, simulations on synthetic networks, and real-data analysis, we show that the combination of both algorithmic and qualitative aspects of our network alignment framework leads to improved performance of our proposed algorithm compared to existing network alignment methods. On simple examples, we isolate multiple aspects of the proposed algorithm and show that each aspect is critical in the performance of the method. The proposed algorithmic and qualitative improvements can also be adapted to existing network alignment packages.

Some network alignment formulations aim to align paths [165] or subgraphs [166, 202] across two (or multiple) networks. The objective of these methods is different from our network alignment optimization where a bijective mapping across nodes of two networks is desired according to a quadratic assignment problem. Solutions of these different methods may be related. For instance, a bijective mapping across nodes of two networks can provide information about conserved pathways and/or subgraphs across networks, and vice versa.

Many networks have modular structures in which groups of nodes tend to interact

¹Laszlo Babai has recently outlined his proof that the computational complexity of the general graph isomorphism problem is Quasipolynomial time [201].

with each other more strongly compared to the rest of the network [92, 203]. One can take advantage of such modular structure to split a large network alignment optimization into small subproblems and use semidefinite programming (SDP) over each subproblem, which yields a more accurate approximation to the underlying QAP at the expense of more computation [180–185]. Here, we consider a special class of modular networks and propose an extension of the spectral network alignment algorithm where spectral information is used to split a large QAP into smaller ones over which a SDP-based optimization [185] is performed. This hybrid method has high performance similar to SDP, while enjoying significantly less computational complexity.

We illustrate the effectiveness of our spectral network alignment algorithm against some existing network alignment methods over various synthetic network models including Erdős-Rényi, power law, regular, and stochastic block structures, under different noise models. Having illustrated the efficiency of our algorithm both theoretically and through simulations, we apply it to two real-data applications: We use network alignment to compare gene regulatory networks across human, fly and worm species which we infer by integrating genome-wide functional and physical genomics datasets from ENCODE and modENCODE consortia (see Chapter 3 for more details). We show that our network alignment method infers conserved regulatory interactions across these species despite long evolutionary distances separating these organisms. Moreover, we find strong conservation of centrally-connected genes and biological pathways, especially for human-fly comparisons. In a second application, we show that our spectral network alignment algorithm improves our ability to de-anonymize small subgraphs of the Twitter follower network sampled in years 2008 and 2009, where we de-anonymize user IDs in 2009 using user IDs in 2008. This application illustrates the extent of personal information that can be retrieved from network structures, and raises additional considerations that need to be addressed in different privacy-related applications.

The rest of this chapter is organized as follows. In Section 5.1, we present the network alignment problem and review existent network alignment techniques. In Section 5.2, we introduce our proposed algorithm and discuss its relationship with

the underlying quadratic assignment problem. Moreover, we present the optimality of our method over random graphs, under some general conditions. In Section 5.3, we consider the trace formulation of the network alignment optimization and introduce a network alignment algorithm which uses higher-order eigenvectors of adjacency graphs to align network structures. In Section 5.4, we consider the network alignment problem of modular networks and introduce an algorithm which solves it efficiently. In Section 8.3, we compare performance of our method with existent network alignment methods over different synthetic network structures. In Section 5.6, we illustrate applications of our method in comparative analysis of regulatory networks across species. In Section 5.7, we explain an application of network alignment in user de-anonymization over twitter follower subgraphs. In Section 5.8, we present proofs for the main results of the chapter.

5.1 Network Alignment Problem Setup

5.1.1 Problem Formulation and Notation

In this section, we introduce the network alignment problem formulation. Let $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ be two graphs (networks) where V_a and E_a represent set of nodes and edges of graph $a = 1, 2$, respectively. By abuse of notation, let G_1 and G_2 be their matrix representations as well where $G_a(i, j) = 1$ iff $(i, j) \in E_a$, for $a = 1, 2$. Suppose network a has n_a nodes, i.e., $|V_a| = n_a$. We assume that networks are un-weighted (binary), and possibly directed. The proposed framework can be extended to the case of weighted graphs as well.

Let X be an $n_1 \times n_2$ binary matrix where $x(i, j') = 1$ means that node i in network 1 is mapped to node j' in network 2. The pair (i, j') is called a mapping edge across two networks. In the network alignment setup, each node in one network can be mapped to at most one node in the other network, i.e., $\sum_i x(i, j') \leq 1$ for all j' , and similarly, $\sum_{j'} x(i, j') \leq 1$ for all i .

Let \mathbf{y} be a vectorized version of X . That is, \mathbf{y} is a vector of length $n_1 n_2$ where,

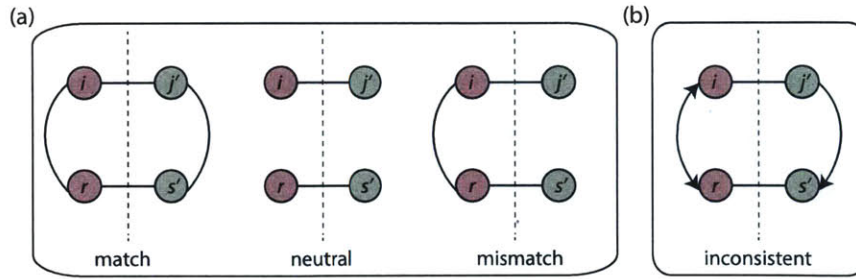


Figure 5-1: (a) An illustration of matched, mismatched, and neutral mappings, for undirected graphs. (b) An illustration of inconsistent mappings, for directed graphs, where they are matches in one direction, and mismatches in the other direction.

$y(i + (j' - 1)n_1) = x(i, j')$. To simplify notation, define $y_{i,j'} \triangleq x(i, j')$. Two mappings (i, j') and (r, s') can be matches which cause overlaps, can be mismatches which cause errors, or can be neutrals (Figure 5-1-a).

Definition 5.1. Suppose $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ are undirected graphs. Let $\{i, r\} \subseteq V_1$ and $\{j', s'\} \subseteq V_2$ where $x(i, j') = 1$ and $x(r, s') = 1$. Then,

- (i, j') and (r, s') are *matches* if $(i, r) \in E_1$ and $(j', s') \in E_2$.
- (i, j') and (r, s') are *mismatches* if only one of the edges (i, r) and (j', s') exists.
- (i, j') and (r, s') are *neutrals* if none of the edges (i, r) and (j', s') exists.

Definition 5.1 can be extended to the case where G_1 and G_2 are directed graphs. In this case, mappings (i, j') and (r, s') are matches/mismatches if they are matches/mismatches in one of the possible directions. However, it is possible to have these mappings be matches in one direction, while they are mismatches in the other direction (Figure 5-1-b). These mappings are denoted as *inconsistent mappings*, defined as follows:

Definition 5.2. Let $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ be two directed graphs and $\{i, r\} \subseteq V_1$ and $\{j', s'\} \subseteq V_2$ where $x(i, j') = 1$ and $x(r, s') = 1$. If edges $i \rightarrow r$, $r \rightarrow i$, and $j' \rightarrow s'$ exist, however, $s' \rightarrow j'$ does not exist, then mappings (i, j') and (r, s') are *inconsistent*.

Existing network alignment formulations aim to find a mapping matrix X which maximizes the number of matches between networks. However, these formulations

can lead to mappings which cause numerous mismatches, especially if networks have different sizes and low similarity. In this chapter, we propose a more general formulation for the network alignment problem which considers both matches and mismatches simultaneously.

For a given alignment matrix X across networks G_1 and G_2 , we assign an *alignment score* by considering the number of matches, mismatches and neutrals caused by X :

$$\text{Alignment Score}(X) = s_1(\# \text{ of matches}) + s_2(\# \text{ of neutrals}) + s_3(\# \text{ of mismatches}), \quad (5.1)$$

where s_1 , s_2 , and s_3 are scores assigned to matches, neutrals, and mismatches, respectively. Note that existing alignment methods ignore effects of mismatches and neutrals by assuming $s_2 = s_3 = 0$ which is restrictive particularly in aligning graphs with different number of nodes and low similarity. In the following, we re-write (5.1) as a quadratic assignment formulation.

Consider two undirected graphs $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$. We form an *alignment network* represented by adjacency matrix A in which nodes are mapping edges across the original networks, and the edges capture whether the pair of mapping edges induce matches, mismatches or neutrals (Figure 5-2).

Definition 5.3. Let $\{i, r\} \subseteq V_1$ and $\{j', s'\} \subseteq V_2$, where $x(i, j') = 1$ and $x(r, s') = 1$.

$$A[(i, j'), (r, s')] = \begin{cases} s_1, & \text{if } (i, j') \text{ and } (r, s') \text{ are matches,} \\ s_2, & \text{if } (i, j') \text{ and } (r, s') \text{ are neutrals,} \\ s_3, & \text{if } (i, j') \text{ and } (r, s') \text{ are mismatches,} \end{cases} \quad (5.2)$$

where s_1 , s_2 , and s_3 are scores assigned to matches, neutrals, and mismatches, respectively.

We can re-write (5.2) as follows:

$$A[(i, j'), (r, s')] = (s_1 + s_2 - 2s_3)G_1(i, r)G_2(j', s') + (s_3 - s_2)(G_1(i, r) + G_2(j', s')) + s_2. \quad (5.3)$$

We can summarize (5.2) and (5.3) as follows:

$$A = (s_1 + s_2 - 2s_3)(G_1 \otimes G_2) + (s_3 - s_2)(G_1 \otimes \mathbf{1}_{n_2}) + (s_3 - s_2)(\mathbf{1}_{n_1} \otimes G_2) + s_2(\mathbf{1}_{n_1} \otimes \mathbf{1}_{n_2}), \quad (5.4)$$

where \otimes represents matrix Kronecker product, and $\mathbf{1}_n$ is an $n \times n$ matrix whose elements are all ones.

A similar scoring scheme can be used for directed graphs. When graphs are directed, some mappings can be inconsistent according to Definition 5.2, i.e., they are matches in one direction and mismatches in another. Scores of inconsistent mappings can be assigned randomly to matched/mismatched scores, or to an average score of matches and mismatches (i.e., $(s_1 + s_3)/2$). For random graphs, inconsistent mappings are rare events. For example, suppose network edges are distributed according to a Bernoulli distribution with parameter p . Then, the probability of having an inconsistent mapping for a particular pair of paired nodes across networks is equal to $4p^3(1 - p)$. Therefore, their effect in network alignment is negligible particularly for large sparse networks. Throughout the chapter, for directed graphs, we assume inconsistent mappings have negligible effect unless it is mentioned explicitly.

Alignment scores s_1 , s_2 and s_3 of (5.2) can be arbitrary in general. However, in this chapter we consider the case where $s_1 > s_2 > s_3 > 0$ with the following rationale: Suppose a mapping matrix X has a total of κ non-zero edges. For example, if networks have $n_1 = n_2 = n$ nodes and there is no unmapped nodes across two networks, $\kappa = n$. The total number of matches, mismatches and neutrals caused by this mapping is equal to $\binom{\kappa}{2}$. Thus, for mapping matrices with the same number of mapping edges, without loss of generality, one can assume that, alignment scores are strictly positive $s_1, s_2, s_3 > 0$ (otherwise, a constant can be added to the right-hand side of (5.2)). In general, mappings with high alignment scores might have slightly different number of

mapping edges owing to unmapped nodes across the networks which has a negligible effect in practice. Moreover, in the alignment scheme, we wish to encourage matches and penalize mismatches. Thus, throughout this chapter, we assume $s_1 > s_2 > s_3 > 0$.

In practice, some mappings across two networks may not be possible owing to additional side information. The set of possible mappings across two networks is denoted by $\mathcal{R} = \{(i, j') : i \in V_1, j' \in V_2\}$. If $\mathcal{R} = V_1 \times V_2$, the problem of network alignment is called *unrestricted*. However, if some mappings across two networks are prevented (i.e., $y_{i,j'} = 0$, for $(i, j') \notin \mathcal{R}$), then the problem of network alignment is called *restricted*.

In the following, we present the network alignment optimization which we consider in this chapter:

Definition 5.4 (Network Alignment Problem Setup). Let $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ be two binary networks. Network alignment aims to solve the following optimization:

$$\begin{aligned}
 \max_{\mathbf{y}} \quad & \mathbf{y}^T A \mathbf{y}, & (5.5) \\
 & \sum_i y_{i,j'} \leq 1, \quad \forall i \in V_1, \\
 & \sum_{j'} y_{i,j'} \leq 1, \quad \forall j' \in V_2, \\
 & y_{i,j'} \in \{0, 1\}, \quad \forall (i, j') \in V_1 \times V_2, \\
 & y_{i,j'} = 0, \quad \forall (i, j') \notin \mathcal{R},
 \end{aligned}$$

where A is defined according to (5.3), and $\mathcal{R} \subseteq V_1 \times V_2$ is the set of possible mappings across two networks.

In the following, we re-write (5.5) using the trace formulation of a standard QAP. Here, we consider undirected networks $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ with n_1 and n_2 nodes, respectively. Without loss of generality, we assume $n_1 \leq n_2$. Moreover, we assume that all nodes of the network G_1 are mapped to nodes of the network G_2 (i.e., there are no un-aligned nodes in G_1).

Define $\xi_1 \triangleq s_1 + s_2 - 2s_3$, $\xi_2 \triangleq s_3 - s_2$ and $\xi_3 \triangleq s_2$. We can rewrite the objective function of Optimization (5.5) using (5.4) as follows:

$$\begin{aligned}
 \mathbf{y}^T A \mathbf{y} &= \xi_1 \mathbf{y}^T (G_1 \otimes G_2) \mathbf{y} + \xi_2 \mathbf{y}^T (G_1 \otimes \mathbf{1}_{n_2}) \mathbf{y} + \xi_2 \mathbf{y}^T (\mathbf{1}_{n_1} \otimes G_2) \mathbf{y} + \xi_3 \mathbf{y}^T (\mathbf{1}_{n_1} \otimes \mathbf{1}_{n_2}) \mathbf{y} \\
 &= \xi_1 \text{Tr}(G_1 X G_2 X^T) + \xi_2 \text{Tr}(G_1 X \mathbf{1}_{n_2} X^T) + \xi_2 \text{Tr}(\mathbf{1}_{n_1} X G_2 X^T) + \xi_3 \text{Tr}(\mathbf{1}_{n_1} X \mathbf{1}_{n_2} X^T) \\
 &= \text{Tr}\left(\left(G_1 + \frac{\xi_2}{\xi_1}\right) X (\xi_1 G_2 + \xi_2) X^T\right) + \left(\xi_3 - \frac{\xi_2^2}{\xi_1}\right) \text{Tr}(\mathbf{1}_{n_1} X \mathbf{1}_{n_2} X^T) \\
 &= \text{Tr}(G'_1 X G'_2 X^T) + \left(\xi_3 - \frac{\xi_2^2}{\xi_1}\right) (\min(n_1, n_2))^2 \\
 &= \text{Tr}(G'_1 X G'_2 X^T) + \text{constant},
 \end{aligned}$$

where

$$\begin{aligned}
 G'_1 &\triangleq G_1 + \frac{\xi_2}{\xi_1} = G_1 + \frac{s_3 - s_2}{s_1 + s_2 - 2s_3}, \\
 G'_2 &\triangleq \xi_1 G_2 + \xi_2 = (s_1 + s_2 - 2s_3) G_2 + (s_3 - s_2).
 \end{aligned}$$

Thus, the network alignment optimization (5.5) can be reformulated as follows:

$$\begin{aligned}
 \max_X \quad & \text{Tr}(G'_1 X G'_2 X^T), \tag{5.6} \\
 & \sum_i x_{i,j'} \leq 1, \quad \forall i \in V_1, \\
 & \sum_{j'} x_{i,j'} \leq 1, \quad \forall j' \in V_2, \\
 & x_{i,j'} \in \{0, 1\}, \quad \forall (i, j') \in V_1 \times V_2, \\
 & x_{i,j'} = 0, \quad \forall (i, j') \notin \mathcal{R}.
 \end{aligned}$$

5.1.2 Prior Work

Network alignment problem (5.5) is an example of a quadratic assignment problem (QAP) [172]. Reference [173] shows that approximating a solution of maximum quadratic assignment problem within a factor better than $2^{\log^{1-\epsilon} n}$ is not feasible in polynomial time in general. However, owing to various applications of QAP in dif-

ferent areas, many methods exist for approximate solutions. In the following, we briefly summarize previous works by categorizing them into four groups and explain advantages and shortcomings of each. For more details on these methods, we refer readers to references [172, 194].

- **Exact search methods:** These methods provide a global optimal solution for the quadratic assignment problem. However, owing to their high computational complexity, they can only be applied to very small problem instances. Examples of exact algorithms include methods based on branch-and-bound [174] and cutting plane [175].
- **Linearizations:** These methods attempt to solve QAP by eliminating the quadratic term in the objective function of Optimization (5.5), transforming it into a mixed integer linear program (MILP). An existing MILP solver is applied to find a solution for the relaxed problem. Examples of these methods are Lawler’s linearization [176], Kaufmann and Broeckx linearization [177], Frieze and Yadegar linearization [178], and Adams and Johnson linearization [179]. These linearizations can provide bounds on the optimal value of the underlying QAP [173]. In general, linearization of the QAP objective function is achieved by introducing many new variables and new linear constraints. In practice, the very large number of introduced variables and constraints poses an obstacle for solving the resulting MILP efficiently.
- **Semidefinite/convex relaxations and bounds:** These methods aim to compute a bound on the optimal value of the network alignment optimization, by considering the alignment matrix in the intersection of the sets of orthogonal and stochastic matrices. The provided solution by these methods may not be a feasible solution of the original quadratic assignment problem. Examples of these methods include orthogonal relaxations [180], projected eigenvalue bounds [181], convex relaxations [182–184], and matrix splittings [185]. In the computer vision literature, [189, 190] use spectral techniques to approximately solve QAP by inferring a cluster of assignments over the feature network. Then,

they use a greedy approach to reject assignments with low associations.

In particular, [185] introduces a convex relaxation of the underlying network alignment optimization based on matrix splitting which provides bounds on the optimal value of the underlying QAP. The proposed SDP method provides a bound on the optimal value and additional steps are required to derive a feasible solution. Moreover, owing to its computational complexity, it can only be used to align small networks, limiting its applicability to alignment of large real networks ². In Section 5.4, we address these issues and introduce a hybrid method based on our proposed scheme in Section 5.2, and the semidefinite relaxation of [185] to align large modular network structures with low computational complexity.

- **Other methods:** There are several other techniques to approximately solve network alignment optimization. Some methods use a Lagrangian relaxation [204], Bayesian framework [187], or message passing [188], or some other heuristics [161, 162, 164]. In Section 8.3, we assess the performance of some of these network alignment techniques through simulations.

5.1.3 Network Alignment and Graph Isomorphism

Network alignment optimization (5.5) is closely related to the problem of *graph isomorphism*, defined as follows:

Definition 5.5 (Graph Isomorphism). Let $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ be two binary networks. G_1 and G_2 are isomorphic if there exists a permutation matrix P such that $G_1 = PG_2P^T$.

The computational problem of determining whether two finite graphs are isomorphic is called the *graph isomorphism problem*. Moreover, given two isomorphic networks G_1 and G_2 , the problem of graph isomorphism aims to find the permutation

²In our simulations, networks should have approximately less than 70 nodes to be able to run it on an ordinary laptop.

matrix P such that $G_1 = PG_2P^T$. The computational complexity of this problem is unknown [205].

Network alignment and graph isomorphism problems are related to each other. Loosely speaking, network alignment aims to minimize the distance between permuted versions of two networks (or, alternatively to maximize their overlap). Therefore, if the underlying networks are isomorphic, an optimal solution of the network alignment optimization should be the same (or close) to the underlying permutation matrix P , where $G_1 = PG_2P^T$. In the following lemma, we formalize such a connection between the network alignment optimization and the classical graph isomorphism problem:

Lemma 5.6. Let G_1 and G_2 be two isomorphic Erdős-Rényi graphs [199] such that $Pr[G_1(i, j) = 1] = p$ and $G_2 = PG_1P^T$, where P is a permutation matrix. Let $p \neq 0, 1$. Then, for any selection of scores $s_1 > s_2 > s_3 > 0$, P maximizes the expected network alignment objective function of Optimization (5.5).

Proof. The proof is presented in Section 5.8.1. □

The result of Lemma 5.6 can be extended to the case where edges of graphs are flipped through a random noise matrix:

Lemma 5.7. Let G_1 be an Erdős-Rényi graphs such that $Pr[G_1(i, j) = 1] = p$. Let \tilde{G}_1 be a graph resulting from flipping edges of G_1 independently and randomly with probability q . Suppose $G_2 = P\tilde{G}_1P^T$ where P is a permutation matrix. Let $0 < p < 1/2$ and $0 \leq q < 1/2$. Then, for any selection of scores $s_1 > s_2 > s_3 > 0$, P maximizes the expected network alignment objective function of Optimization (5.5).

Proof. The proof is presented in Section 5.8.1. □

Finding an isomorphic mapping across sufficiently large Erdős-Rényi graphs can be done efficiently with high probability (w.h.p.) through canonical labeling [200]. Canonical labeling of a network consists of assigning a unique label to each vertex such that labels are invariant under isomorphism. The graph isomorphism problem can then be solved efficiently by mappings nodes with the same canonical labels to each other [206]. One example of canonical labeling is the degree neighborhood of a

vertex defined as a sorted list of neighborhood degrees of vertices [200]. Note that, network alignment formulation is more general than the one of graph isomorphism; network alignment aims to find an optimal mappings across two networks which are not necessarily isomorphic.

Remark 5.8. In [207], Babai, Erdős, and Selkow derive an interesting and perhaps a counter-intuitive result on random graph isomorphism. Based on their result, *any two infinite random graphs are isomorphic with probability one*. For instance, consider two infinite Erdős-Rényi graphs G_1 and G_2 where $Pr[G_1(i, j) = 1] = p_1$ and $Pr[G_2(i, j) = 1] = p_2$ and $p_1, p_2 \neq 0, 1$. The result of [207] indicates that, G_1 and G_2 are isomorphic even if $p_1 \neq p_2$. This may seem counter-intuitive as two networks may seem to have different densities. However, this result is only true for *infinite graphs*, not asymptotically large ones. The difference may seem subtle but significant. In infinite graphs, notions of graph size, graph density, etc. are different than the ones for finite graphs. Throughout this chapter, we only consider finite or asymptotically large graphs, not infinite ones. In the following, we give some intuition on the results of [207] about infinite graphs.

The proof of the result in [207] is based on a notion for infinite graphs called *the extension property* which roughly states that, an infinite graph G has the extension property if for any two disjoint finite subsets of nodes V_1 and V_2 , there exists a vertex $v \in V - V_1 - V_2$ such that v is connected to all vertices in V_1 and to no vertices in V_2 . If two infinite graphs have extension properties, they are isomorphic with probability one. One way to prove this is to construct an infinite sequences of mappings by considering disjoint subsets of nodes V_1 and V_2 , and adding a node v according to the extension property.

Finally, it is straightforward to show that an infinite Erdős-Rényi graph with parameter $p \neq 0, 1$ has the extension property with probability one. This is analogous to the monkey-text book problem: if a monkey randomly types infinite letters, with probability one, he will type any given text book. This is because the number of letters in a text book is finite and therefore, the probability that the letter sequence of a given text book does not appear in a random infinite letter sequence is zero.

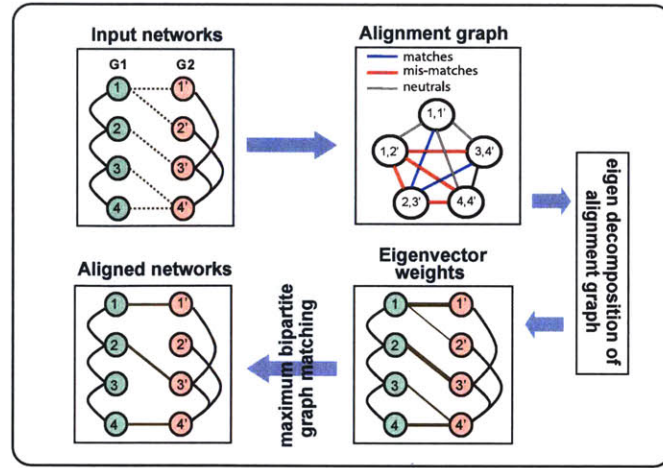


Figure 5-2: Framework of EigenAlign algorithm 5.9.

Now consider two finite disjoint subsets V_1 and V_2 over an infinite graph. Note that being finite is key here. Similarly, the probability that a vertex v exists such that it is connected to all vertices in V_1 and to no vertices in V_2 is one because its complement set has zero probability (or it has zero Lebesgue measure if we map binary sequences to numbers in $[0, 1]$).

5.2 EigenAlign Algorithm

We now introduce *EigenAlign*, an algorithm which solves a relaxation of the network alignment optimization (5.5) leveraging spectral properties of networks. Unlike other alignment methods, EigenAlign considers both matches and mismatches in the alignment scheme. Moreover, we prove its optimality (in an asymptotic sense) in aligning Erdős-Rényi graphs under some technical conditions. In the following, we describe the EigenAlign algorithm:

Algorithm 5.9 (EigenAlign Algorithm). Let $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ be two binary networks whose corresponding alignment network is denoted by A according to (5.3). EigenAlign algorithm solves the network alignment optimization (5.5) in two steps:

- **Eigenvector Computation Step:** In this step, we compute \mathbf{v} , an eigenvector of the alignment network A with the maximum eigenvalue.
- **Linear Assignment Step:** In this step, we solve the following maximum weight bipartite matching optimization:

$$\begin{aligned}
\max_{\mathbf{y}} \quad & \mathbf{v}^T \mathbf{y}, \\
& \sum_{j'} y_{i,j'} \leq 1, \quad \forall i \in V_1, \\
& \sum_i y_{i,j'} \leq 1, \quad \forall j' \in V_2, \\
& y_{i,j'} \in \{0, 1\}, \quad \forall (i, j') \in V_1 \times V_2, \\
& y_{i,j'} = 0, \quad \forall (i, j') \notin \mathcal{R}.
\end{aligned} \tag{5.7}$$

This framework is depicted in Figure 5-2. In the rest of this section, we provide intuition on different steps of the EigenAlign algorithm through both quadratic assignment relaxation argument as well as a fixed point analysis. In Section 5.2.3, we discuss optimality of EigenAlign over random graphs. In Section 5.3, we will introduce an extension of the EigenAlign algorithm which uses higher-order eigenvectors of adjacency graphs to align network structures.

5.2.1 EigenAlign as a Relaxation of Quadratic Assignment

In this section, we explain EigenAlign as a relaxation of the underlying quadratic assignment optimization (5.5). For simplicity, we assume all mappings across networks are possible (i.e., $\mathcal{R} = \{(i, j') : \forall i \in V_1, \forall j' \in V_2\}$). In the restricted network alignment setup, without loss of generality, one can eliminate rows and columns of the alignment matrix corresponding to mappings that are not allowed.

In the eigen decomposition step of EigenAlign, we ignore bijective constraints (i.e., constraints $\sum_i y_{i,j'} \leq 1$ and $\sum_{j'} y_{i,j'} \leq 1$) because they will be satisfied in the second step of the algorithm through a linear optimization. By these assumptions,

Optimization (5.5) can be simplified to the following optimization:

$$\begin{aligned} \max_{\mathbf{y}} \quad & \mathbf{y}^T A \mathbf{y}, \\ & y_{i,j'} \in \{0, 1\}, \quad \forall (i, j') \in V_1 \times V_2. \end{aligned} \tag{5.8}$$

To approximate a solution of this optimization, we relax integer constraints to constraints over a hyper-sphere restricted by hyper-planes (i.e., $\|\mathbf{y}\|_2 \leq 1$ and $\mathbf{y} \geq 0$). Using this relaxation, Optimization (5.8) is simplified to the following optimization:

$$\begin{aligned} \max_{\mathbf{y}} \quad & \mathbf{y}^T A \mathbf{y}, \\ & \|\mathbf{y}\|_2 \leq 1, \\ & \mathbf{y} \geq 0. \end{aligned} \tag{5.9}$$

In the following, we show that \mathbf{v} , the leading eigenvector of the alignment matrix A is an optimal solution of Optimization (5.9). Suppose \mathbf{y}_1 is an optimal solution of Optimization (5.9). Let \mathbf{y}_2 be a solution of the following optimization which ignores non-negativity constraints:

$$\begin{aligned} \max_{\mathbf{y}} \quad & \mathbf{y}^T A \mathbf{y}, \\ & \|\mathbf{y}\|_2 \leq 1. \end{aligned} \tag{5.10}$$

Following the Rayleigh – Ritz formula, the leading eigenvector of the alignment matrix is an optimal solution of Optimization (5.10) (i.e., $\mathbf{y}_2 = \mathbf{v}$). Now we use the following theorem to show that in fact $\mathbf{y}_1 = \mathbf{v}_1$:

Theorem 5.10 (Perron–Frobenius Theorem [160]). Suppose A is a matrix whose elements are strictly positive. Let \mathbf{v} be an eigenvector of A corresponding to the largest eigenvalue. Then, $\forall i, v_i > 0$. Moreover, all other eigenvectors must have at least one negative, or non-real component.

Since \mathbf{y}_2 is a solution of a relaxed version of Optimization (5.10), we have $\mathbf{y}_2^T A \mathbf{y}_2 \geq$

$\mathbf{y}_1^T A \mathbf{y}_1$. Using this inequality along with Perron-Frobenius Theorem lead to $\mathbf{y}_1 = \mathbf{v}$, as the unique solution of optimization (5.9).

The solution of the eigen decomposition step assigns weights to all possible mappings across networks ignoring bijective constraints. However, in the network alignment setup, each node in one network can be mapped to at most one node in the other network. To satisfy these constraints, we use eigenvector weights in a linear optimization framework of maximum weight bipartite matching setup of Optimization (5.7) [208].

Remark 5.11. Many existing network alignment techniques are based on iterative algorithms [161, 162, 188]. The EigenAlign relaxation of the underlying quadratic assignment problem can be viewed as the following fixed point solution using a linear mapping function which can be solved iteratively:

$$y(t_1) = \frac{1}{\lambda} \sum_{t_2} A(t_1, t_2) y(t_2), \quad (5.11)$$

where λ is a constant and t_i represents a node in the alignment network A . Note that, this mapping function is not necessarily a contraction. Instead of using an iterative approach, EigenAlign solves this relaxation in a closed form using the leading eigenvector of the alignment network. In particular, (5.11) can be re-written as,

$$A\mathbf{y} = \lambda\mathbf{y}, \quad (5.12)$$

where \mathbf{y} is the vector of mapping weights whose elements represent a measure of centrality of nodes in the alignment network. If A is diagonalizable, eigenvectors of matrix A provide solutions for this equation, where λ is the eigenvalue corresponding to eigenvector \mathbf{y} . However, eigenvectors of matrix A can have negative components which are not allowed in the mapping function since $y_{i,j'} \geq 0$. Since components of matrix A are all positive, similarly to discussions of (5.9), according to Perron-Frobenius Theorem 5.10, all components of the eigenvector associated with the maximum positive eigenvalue are positive. This result holds for non-negative matrices as

well if they are strongly connected. Therefore, the leading eigenvector of matrix A satisfies desired conditions for the mapping function and provides alignment scores for the nodes in the alignment network.

Remark 5.12. Suppose G_1 and G_2 are two undirected graphs with n_1 and n_2 nodes, respectively. IsoRank [161] is a network alignment method which aligns networks using neighborhood similarities of nodes. Let I be the identity matrix of size $n_1 \times n_1$. Suppose D_1 and D_2 are diagonal matrices of sizes $n_1 \times n_1$ and $n_2 \times n_2$ whose diagonal elements are node degrees of networks G_1 and G_2 , respectively. The first step of the IsoRank algorithm solves the following fixed point equation:

$$\mathbf{y} = (I \otimes D_1^{-1})A(I \otimes D_2^{-1})\mathbf{y}, \quad (5.13)$$

where \otimes represents matrix Kronecker product, and $s_2 = s_3 = 0$ in (5.3). IsoRank solves (5.13) using a power iteration method. The IsoRank algorithm has several differences with EigenAlign in both algorithmic and qualitative aspects. From a qualitative aspect, IsoRank ignores effects of mismatches across networks by assuming $s_2 = s_3 = 0$. IsoRank uses a stationary centrality measure in the alignment network, having a degree scaling step which is based on the assumption of uniformly random distribution of alignment scores over potential mappings. Moreover, it uses a greedy method to select bijective mappings across networks, instead of the linear optimization framework used in the EigenAlign algorithm. In Sections 8.3 and 5.6, over both synthetic and real networks, we show that, the EigenAlign solution has significantly better performance compared to the one of IsoRank. Moreover, in Section 8.3, we isolate multiple aspects of the EigenAlign algorithm on some toy examples, and evaluate their individual contributions in its performance. We note that, existent network alignment methods such as IsoRank may be improved by borrowing algorithmic and/or qualitative ideas of the proposed EigenAlign method. However, these extensions are beyond the scope of this work.

5.2.2 Computational Complexity of EigenAlign

In this part, we analyze computational complexity of the EigenAlign Algorithm 5.9. Suppose the number of nodes of networks G_1 and G_2 are $\mathcal{O}(n)$. Let $k = \mathcal{O}(|\mathcal{R}|)$ be the number of possible mappings across two networks. In an unrestricted network alignment setup, we have $k = \mathcal{O}(n^2)$. However, in a restricted network alignment, k may be significantly smaller than n^2 . EigenAlign has three steps:

- (i) The alignment network A should be formed which has a computational complexity of $\mathcal{O}(k^2)$ because all pairs of possible mappings should be considered.
- (ii) In the eigen decomposition step, we need to compute the leading eigenvector of the alignment network. This operation can be performed in almost linear time in k using QR algorithms and/or power methods [209]. Therefore, the worst case computational complexity of this part is $\mathcal{O}(k)$.
- (iii) Finally, we use eigenvector weights in a maximum weight bipartite matching algorithm which can be solved efficiently using linear programming or Hungarian algorithm [208]. The worst case computational complexity of this step is $\mathcal{O}(n^3)$. If the set \mathcal{R} has a specific structure (e.g., small subsets of nodes in one network are allowed to be mapped to small subsets of nodes in the other network), this cost can be reduced significantly. In Section 5.6, we see this structure in aligning regulatory networks across species as genes are allowed to be aligned to homologous genes within their gene families.

Proposition 5.13. The worst case computational complexity of the EigenAlign Algorithm is $\mathcal{O}(k^2 + n^3)$.

Remark 5.14. For large networks, to reduce the overall computational complexity, the linear assignment optimization may be replaced by a greedy bipartite matching algorithm. In the greedy matching approach, at each step, the heaviest possible mapping is added to the current matching until no further mappings can be added. It is straightforward to show that this greedy algorithm finds a bipartite matching

whose weight is at least half the optimum. The computational complexity of this greedy algorithm is $\mathcal{O}(k \log(k) + nk)$.

5.2.3 Performance Characterization of EigenAlign Over Erdős-Rényi Graphs

In this section, we analyze optimality of the EigenAlign algorithm over Erdős-Rényi graphs, for both isomorphic and non-isomorphic cases, and under two different noise models. While real networks often have different structures than Erdős-Rényi graphs, the analytical tractability and the spectral characterization of these graphs makes it possible to actually prove that the proposed relaxations are asymptotically tight, at least in special cases, providing a rigorous foundation for the proposed EigenAlign algorithm. In this section, we only consider finite and asymptotically large graphs. For arguments on infinite graphs, see Section 5.1.3.

Suppose $G_1 = (V_1, E_1)$ is an undirected Erdős-Rényi graph with n nodes where $Pr[G_1(i, j) = 1] = p$. Self-loops are allowed as well. Suppose \tilde{G} is a noisy version of the graph G_1 . We consider two different noise models in this section:

- Noise Model I: In this model, we have,

$$\tilde{G}_1 = G_1 \odot (1 - Q) + (1 - G_1) \odot Q, \quad (5.14)$$

where \odot represents the Hadamard product, and Q is a binary random matrix whose edges are drawn i.i.d. from a Bernoulli distribution with $Pr[Q(i, j) = 1] = p_e$. In words, the operation $G_1 \odot (1 - Q) + (1 - G_1) \odot Q$ flips edges of G_1 independently randomly with probability p_e .

- Noise Model II: In this model, we have,

$$\tilde{G}_1 = G_1 \odot (1 - Q) + (1 - G_1) \odot Q', \quad (5.15)$$

where Q and Q' are binary random matrices whose edges are drawn i.i.d. from

a Bernoulli distribution with $Pr[Q(i, j) = 1] = p_e$ and $Pr[Q'(i, j) = 1] = p_{e_2}$. Under this model, edges of G_1 flip independently randomly with probability p_e , while non-connecting tuples in G_1 will be connected in \tilde{G}_1 with probability p_{e_2} . Because G_1 is an Erdős-Rényi graph with parameter p , choosing

$$p_{e_2} = \frac{pp_e}{1-p}, \quad (5.16)$$

leads to the expected density of networks G_1 and G_2 be p .

We define G_2 as follows:

$$G_2 = P\tilde{G}_1P^T, \quad (5.17)$$

where P is a permutation matrix. Throughout this section, we assume that we are in the restricted network alignment regime: we desire to choose n mappings $i \leftrightarrow j'$ across two networks among $|\mathcal{R}| = kn$ possible mappings where $i \in V_1$, $j' \in V_2$, and $k > 1$. n true mappings ($i \leftrightarrow i'$ if $P = I$) are included in \mathcal{R} , while the remaining $(k-1)n$ mappings are selected independently randomly. Moreover, we choose scores assigned to matches, neutrals and mismatches as $s_1 = \alpha + \epsilon$, $s_2 = 1 + \epsilon$ and $s_3 = \epsilon$, respectively, where $\alpha > 1$ and $0 < \epsilon\lambda 1$. These selections satisfy score conditions $s_1 > s_2 > s_3 > 0$.

Theorem 5.15 (EigenAlign over Erdős-Rényi graphs). Let \tilde{P} be the solution of the EigenAlign Algorithm 5.9. Then, under both noise models (5.14) and (5.15), if $0 < p < 1/2$, and $0 \leq p_e < 1/2$, then as $n \rightarrow \infty$, the error probability goes to zero:

$$Pr\left[\frac{1}{n}\|\tilde{P} - P\| \rightarrow 0\right].$$

Theorem 5.15 states that, the EigenAlign algorithm is able to recover the underlying permutation matrix which relates networks G_1 and G_2 to each other according to (5.17). On the other hand, according to Lemma 5.7, this permutation matrix is in fact optimizes the expected network alignment score.

Proposition 5.16. Under conditions of Theorem 5.15, the permutation matrix inferred by EigenAlign maximizes the expected network alignment objective function defined according to Optimization (5.5).

In noise models (5.14) and (5.15), if we put $p_e = 0$, then G_2 is isomorphic with G_1 because there exists a permutation matrix P such that $G_2 = PG_1P^T$. For this case, we have the following Corollary:

Corollary 5.17 (EigenAlign on Isomorphic Erdős-Rényi graphs). Let G_1 and G_2 be two isomorphic Erdős-Rényi graphs with n nodes such that $G_1 = PG_2P^T$, where P is a permutation matrix. Under conditions of Theorem 5.15, as $n \rightarrow \infty$, the error probability of EigenAlign solution goes to zero.

We present proofs of Theorem 5.15 and Corollary 5.17 in Sections 5.8.2 and 5.8.3. In the following, we sketch main ideas of their proofs: Since input networks G_1 and G_2 are random graphs, the alignment network formed according to (5.3) will be a random graph as well. The first part of the proof is to characterize the leading eigenvector of this random alignment network. To do this, we first characterize the leading eigenvector of the expected alignment network which in fact is a deterministic graph. In particular, we prove that eigenvector scores assigned to true mappings is strictly larger than the ones assigned to false mappings. To prove this, we characterize top eigenvalues and eigenvectors of the expected alignment network algebraically. The restricted alignment condition (i.e., $|\mathcal{R}| = kn$) is necessary to have this bound. Then, we use Wedin Sin Theorem [210] from perturbation theory, Gershgorian circle Theorem [211] from spectral matrix theory, and Chernoff bound to characterize the leading eigenvector of the random alignment network for sufficiently large n . Finally, we use Chebyshev's inequality to show that the error probability of the EigenAlign algorithm is asymptotically zero w.h.p.

Note that finding an isomorphic mapping across asymptotically large Erdős-Rényi graphs (Corollary 5.17) is a well studied problem and can be solved efficiently through canonical labeling [200]. However, those techniques do not address a more general network alignment problem similar to the setup considered in Theorem 5.15. For

more details, see Section 5.1.3.

Theorem 5.15 and Corollary 5.17 consider a restricted network alignment case where $|\mathcal{R}| = kn$. As explained briefly in the proof sketch and with more details in Lemma 5.29, this technical condition is necessary to show that, expected eigenvector scores of true mappings are strictly larger than the ones of false mappings as $n \rightarrow \infty$. In Section 8.3 and through simulations, we show that error of the EigenAlign algorithm is empirically small even in an unrestricted network alignment setup.

5.3 Low Rank Alignment of Networks

The EigenAlign algorithm introduced in Section 5.2 uses the leading eigenvector of the alignment graph to align graph structures. EigenAlign can be viewed as the rank-one approximation of the linearization of the underlying quadratic assignment problem. In this section, we introduce an extension of EigenAlign which uses higher-order eigenvector of adjacency graphs to align network structures. We refer to this extension as *LowRank Align*. LowRank Align can be useful specially in cases where leading eigenvectors of graphs are not informative. This case occurs for instance in the alignment of regular graph structures. Moreover, LowRank Align does not require an explicit formation of the alignment graph which can be costly for large networks if all mappings across networks are possible.

Suppose G_1 and G_2 are input networks with n_1 and n_2 nodes. For simplicity, in this section we assume that they are symmetric. Recall that X represents the mapping matrix whose size is $n_1 \times n_2$. In this section, we use the trace formulation of the quadratic assignment problem.

The maximum overlap network alignment optimization can be written as follows:

$$\begin{aligned}
 \max \quad & Tr(G_1 X G_2 X^T), \\
 \sum_{i=1}^{n_1} x_{i,j} & \leq 1, \quad \forall 1 \leq j \leq n_2, \\
 \sum_{j=1}^{n_2} x_{i,j} & \leq 1, \quad \forall 1 \leq i \leq n_1, \\
 x_{i,j} & \in \{0, 1\}, \quad \forall 1 \leq i \leq n_1, \forall 1 \leq j \leq n_2.
 \end{aligned} \tag{5.18}$$

Optimization (5.18) finds a mapping matrix X that maximizes the number of overlapping edges (matches) across two graphs. However, such a mapping can cause numerous mismatches. We discuss a balanced network alignment optimization later in this section.

For simplicity, we assume $n_1 = n_2 = n$. All discussions can be extended to the case where $n_1 \neq n_2$. Let Π be the set of all permutation matrices of size $n \times n$. Thus, optimization (5.18) can be written as follows:

$$\begin{aligned}
 \max \quad & Tr(G_1 X G_2 X^T), \\
 X & \in \Pi.
 \end{aligned} \tag{5.19}$$

Let X^* be an optimal solution of optimization (5.19). Finding an optimal solution of this optimization is known to be NP-hard [173].

If $X \in \Pi$, we have,

$$Tr(G_1 X G_2 X^T) = Tr((G_1 + \delta_1 I) X (G_2 + \delta_2 I) X^T) + \text{constant}. \tag{5.20}$$

In other words, we can add and subtract multiples of identity to make the resulting symmetric matrices positive definite, without changing the structure of the problem. Thus, without loss of generality, we assume that matrices G_1 and G_2 are positive semi-definite.

Several algorithms have considered relaxations/approximations of this optimiza-

tion as follows:

$$\begin{aligned} \max \quad & Tr(G_1 X G_2 X^T), \\ & X \in \Gamma. \end{aligned} \tag{5.21}$$

Suppose X_0 is a solution of optimization (5.21) such that $\|X^* - X_0\|_2 \leq \epsilon$ (i.e., its distance from the optimal solution is bounded by ϵ). X_0 may not be a valid permutation matrix. One way to find a permutation matrix using X_0 is to project X_0 over the space of permutation matrices Π :

$$\begin{aligned} \max \quad & Tr(X X_0^T), \\ & X \in \Pi. \end{aligned} \tag{5.22}$$

However, it has been shown that an optimal solution of optimization (5.22) has a poor perform in practice [212]. In the following, we propose an alternative algorithm to efficiently infer a permutation matrix using X_0 with a certain performance guarantee:

Consider the following optimization:

$$\begin{aligned} \max \quad & Tr(G_1 X_0 G_2 X^T), \\ & X \in \Pi. \end{aligned} \tag{5.23}$$

This is a maximum weight bipartite matching optimization which can be solved exactly using linear programming. Let X_{lin}^* be an optimal solution of optimization (5.23). Define

$$\begin{aligned} f(X) &\triangleq Tr(G_1 X G_2 X^T), \\ \tilde{f}(X) &\triangleq Tr(G_1 X_0 G_2 X_0^T) + 2Tr(G_1 X_0 G_2 (X - X_0)^T). \end{aligned} \tag{5.24}$$

Theorem 5.18. Let X^* and X_{lin}^* be optimal solutions of optimizations (5.18) and

(5.23), respectively. We have,

$$|f(X^*) - \tilde{f}(X_{lin}^*)| \leq \epsilon^2 \sum_{i=1}^n \sigma_i(G_1)\sigma_i(G_2), \quad (5.25)$$

where ϵ is a bound on the relaxation gap ($\|X^* - X_0\|_2 \leq \epsilon$), and $\sigma_i(G_a)$ represents the i -th largest singular value of matrix G_a , for $a = 1, 2$.

Proof. See Section 5.8.4. □

If Γ is assumed to be the set of orthogonal matrices (i.e., $\Gamma = \mathcal{O}$), an optimal solution of optimization (5.21) can be found using eigen decomposition of matrices G_1 and G_2 as follows:

Theorem 5.19. Suppose v_i and u_i are eigenvectors of symmetric matrices G_1 and G_2 corresponding to the i -th largest eigenvalues $\sigma_i(G_1)$ and $\sigma_i(G_2)$, respectively. Let V and U be eigenvector matrices whose i -th columns are v_i and u_i , respectively. Then,

$$X_0 = VU^T = \sum_{i=1}^n v_i u_i^T, \quad (5.26)$$

is an optimal solution of optimization (5.21) over orthogonal matrices (i.e., $\Gamma = \mathcal{O}$).

Proof. See reference [180]. □

Remark 5.20. Theorem 5.19 characterizes an optimal solution of the orthogonal relaxation of the network alignment optimization. However, for any v_i , $-v_i$ is also an eigenvector of the matrix G_1 . A similar argument can be mentioned for eigenvectors of the matrix G_2 . Let

$$\mathcal{X}_0 \triangleq \left\{ X_0 : X_0 = \sum_{i=1}^n s_i v_i u_i^T, \mathbf{s} \in \{-1, 1\}^n \right\}, \quad (5.27)$$

where s_i is the i -th component of the vector \mathbf{s} . The set \mathcal{X}_0 represents multiple optimal solutions of optimization (5.21) when $\Gamma = \mathcal{O}$. \mathcal{X}_0 can have at most 2^n distinct members.

Suppose $X_0 \in \mathcal{X}_0$. Then, optimization (5.23) can be simplified to the following optimization which finds a valid permutation matrix using the orthogonal relaxation of the network alignment optimization:

$$\begin{aligned} \max \quad & Tr\left(\left(\sum_{i=1}^n \sigma_i(G_1)\sigma_i(G_2)v_i u_i^T\right)X^T\right), \\ & X \in \Pi. \end{aligned} \tag{5.28}$$

This optimization can be solved efficiently using linear programming.

Remark 5.21. A projection of $X_0 \in \mathcal{X}_0$ over permutation matrices according to optimization (5.22) leads to the following optimization:

$$\begin{aligned} \max \quad & Tr\left(\left(\sum_{i=1}^n v_i u_i^T\right)X^T\right), \\ & X \in \Pi. \end{aligned} \tag{5.29}$$

Unlike optimization (5.29), the objective function of (5.28) takes into account eigenvalues of matrices G_1 and G_2 as well. In other words, eigenvectors with largest eigenvalues contribute more to the objective function compared to the ones with small eigenvalues. To this end, the objective function of optimization (5.28) simplifies the network alignment problem to the simultaneous alignment of eigenvectors whose contributions in the overall alignment score are weighed by their corresponding eigenvalues.

One efficient way to approximate a solution for the original network alignment optimization (5.19) would be to obtain a solution X_0 for its orthogonal relaxation using eigen decomposition of input matrices (Theorem 5.19), and then to use optimization (5.28) to obtain a feasible solution with bounded error (Theorem 5.18). However, as we discussed above, there are possibly exponentially many optimal solutions for optimization (5.28) and obtaining their resulting permutation matrices would be computationally infeasible. However, because contributions of eigenvectors with small eigenvalues to the objective function of optimization (5.28) is small, one

can instead solve the following optimization based on the low rank approximation of the objective function:

$$\begin{aligned} \max \quad & Tr\left(\left(\sum_{i=1}^k \sigma_i(G_1)\sigma_i(G_2)v_iu_i^T\right)X^T\right), \\ & X \in \Pi, \end{aligned} \tag{5.30}$$

where objective function considers only top k eigenvectors of input matrices (i.e., its affinity matrix is low rank). The following optimization summarizes the *LowRank Align algorithm*:

$$\begin{aligned} \max \quad & Tr\left(\left(\sum_{i=1}^k s_i\sigma_i(G_1)\sigma_i(G_2)v_iu_i^T\right)X^T\right), \\ & X \in \Pi, \\ & s_i \in \{-1, 1\}, \quad \forall 1 \leq i \leq k. \end{aligned} \tag{5.31}$$

where k is a constant that determines the rank of the affinity matrix.

In the restricted network alignment setup, some mapping pairs may not be allowed. In that case, one can set the affinity weights of such pairs in optimization (5.31) to be $-\infty$.

Remark 5.22. Unlike EigenAlign, LowRank Align does not require an explicit formation of the alignment graph which can be costly for large networks if all mappings across them are possible. Moreover, in such cases, the linear assignment step can be performed using a greedy method described in Remark 5.14 to reduce the overall computational complexity of the alignment method.

Optimization (5.18) finds a mapping across two graphs G_1 and G_2 to maximize the number of overlapping edges (matches) across them. However, the resulting mapping matrix can cause numerous mismatches, particularly for graphs with different sizes, because the objective function of optimization (5.18) does not penalize the number of mismatches. The number of mismatches caused by a mapping matrix X across two

graphs G_1 and G_2 is equal to:

$$\begin{aligned} \text{number of mismatches} &= \text{Tr}(G_1 X (1 - G_2) X^T + G_1 X (1 - G_2) X^T) \\ &= -2\text{Tr}(G_1 X G_2 X^T) + \text{Tr}(G_1 X \mathbf{1} X^T) + \text{Tr}(\mathbf{1} X G_2 X^T), \end{aligned} \quad (5.32)$$

where $\mathbf{1}$ represents a matrix of all ones. The following optimization considers both matches and mismatches caused by a mapping matrix X :

$$\begin{aligned} \max \quad & (\beta_1 + 2\beta_2)\text{Tr}(G_1 X G_2 X^T) - \beta_2\text{Tr}(G_1 X \mathbf{1} X^T + \mathbf{1} X G_2 X^T), \\ & X \in \Pi, \end{aligned} \quad (5.33)$$

where Π is the set of all permutation matrices of size $n \times n$. β_1 and β_2 are match and mismatch contribution weights in the objective function of optimization (5.33). This optimization is similar to an alignment optimization where the number of mismatches is bounded. Moreover, the objective function of Optimization (5.33) is equivalent to the one introduced in (5.1) for some choices of s_1 , s_2 and s_3 . One can simply this optimization to the following optimization:

$$\begin{aligned} \max \quad & \text{Tr}((G_1 - \alpha)X(G_2 - \alpha)X^T), \\ & X \in \Pi, \end{aligned} \quad (5.34)$$

where

$$\alpha = \frac{\beta_2}{\beta_1 + \beta_2} = \frac{1}{2 + \beta_1/\beta_2} = \frac{1}{2 + R}. \quad (5.35)$$

To obtain optimization (5.34), we ignore the term $\text{Tr}(\mathbf{1} X \mathbf{1} X^T)$ since it is a constant term over permutation matrices $X \in \Pi$. This is also true even when the size of two graphs are not equal (i.e., $n_1 \neq n_2$) if all nodes of the smaller graph are mapped to some nodes in the larger graph. To solve optimization (5.34), one needs to choose α (or equivalently $R = \beta_1/\beta_2$). This parameter tunes contributions of matches and mismatches in the network alignment objective function. Choosing $\alpha \rightarrow 0$ ($R \rightarrow \infty$)

results in merely maximizing the number of matches across two graphs, while $\alpha = 1/2$ only considers minimizing the number of mismatches across two networks. To have balance between maximizing matches and minimizing mismatches across two networks, one can choose,

$$R = \frac{\beta_1}{\beta_2} = \frac{\text{number of mismatches}}{\text{number of matches}}. \quad (5.36)$$

5.4 Alignment of Modular Networks

Many networks have modular structures in which groups of nodes tend to interact with each other more strongly compared to the rest of the network [92, 203]. One can take advantage of such modular structure to split the large network alignment optimization into small subproblems and use a computationally expensive, but tight, semidefinite programming (SDP) [185] over each subproblem.

As we discuss in Section 5.2.1, EigenAlign provides a simple relaxation of the QAP in (5.5), based on the eigenvector centrality of the alignment graph. This can be thought of as a linearization of the quadratic assignment cost along the direction of the eigenvector centrality. EigenAlign relaxes integer constraints to constraints over a hyper-sphere restricted by hyper-planes. This relaxation leads to an efficient method to align large networks. There are several methods based on convex relaxations of the underlying QAP that seek solutions in the intersection of orthogonal and stochastic matrices [180–185]. In general, these relaxations are tighter (hence better) than the one used in the EigenAlign algorithm and lead to a better approximation of the optimal solution. However, these methods have high computational complexity, which limits their applicability in aligning large networks. For instance, as reported in [185] and also according to our experiments, to be able to run the SDP-based method proposed in [185] on an ordinary laptop, networks should not have more than around 70 nodes. Although this threshold can be improved by a better implementation of the method, this SDP-based method is not scalable for large networks. Moreover, distributed first-order methods to solve large SDPs [213] based on low-rank factor-

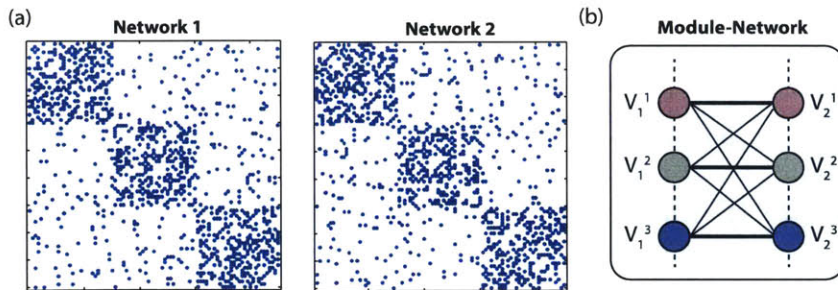


Figure 5-3: (a) Example modular network structures of Definition 5.23. (b) An illustration of the module-network bipartite graph whose nodes correspond to modules in the original networks.

ization [214], block coordinate descent methods [215], spectral bundle methods [216], and eigenvalue saddle point transformations [217] can be enhanced by splitting a large optimization into smaller ones, particularly for modular network structures.

Here, we propose a method, which we term *EigenAlign+SDP*, that uses the EigenAlign algorithm along with an SDP-based relaxation of the underlying QAP over modular network structures. The proposed algorithm uses the EigenAlign solution to identify small subgraphs (modules, or blocks) across two networks which are likely to be aligned to each other. Then, it uses a more computationally expensive SDP-based method to solve multiple small-size network alignment problems, in parallel. The proposed method is based on this key insight that, the EigenAlign solution provides a robust mapping of modules across networks, which enables the use of more expensive and accurate optimization methods over small sub-problems.

In this section, we consider stochastic block network structures:

Definition 5.23. Let $G_1 = (V_1, E_1)$ be an undirected graph. Suppose there is a partition of nodes $\{V_1^1, V_1^2, \dots, V_1^m\}$, where nodes in the same partition are connected to each other independently randomly with probability p , while nodes across partitions are connected to each other independently randomly with probability q .

Let \tilde{G}_1 be a noisy version of G_1 . Here, we consider the noise model II (5.15) as it is more general, while all arguments can be extended to the noise model I (5.14). We assume G_2 is a permuted version of \tilde{G}_1 according to (5.17), i.e., $G_2 = P\tilde{G}_1P^T$.

Figure 5-3-a demonstrates example networks according to this model.

In the following, we present the *EigenAlign+SDP* method which aims to infer an optimal alignment across G_1 and G_2 .

Algorithm 5.24 (EigenAlign+SDP Algorithm). Let $G_1 = (V_1, E_1)$ be a stochastic block matrix of Definition 5.23, and $G_2 = (V_2, E_2)$ is defined according to (5.17). The EigenAlign+SDP algorithm solves the network alignment optimization (5.5) in the following steps:

- **EigenAlign Step:** In this step, we compute the EigenAlign solution across G_1 and G_2 .
- **Spectral Partitioning Step:** In this step, we use a spectral partitioning method [2] to partition each network to m blocks.
- **Module Alignment Step:** In this step, we use the EigenAlign solution in a maximum weight bipartite matching optimization to compute an optimal alignment of modules across G_1 and G_2 .
- **SDP-based Alignment Step:** In this step, we use an SDP-based relaxation of the underlying QAP [185] followed by a maximum weight bipartite matching step, to compute node alignments in each module pairs.

In the following, we explain the above steps in more detail. In the EigenAlign step, we use EigenAlign algorithm 5.9 to find a mapping across networks G_1 and G_2 , denoted by X^* . In the spectral partitioning step, we use the algorithm proposed in [2] to find m blocks (modules) in each network. These modules are denoted by $\{V_1^1, V_1^2, \dots, V_1^m\}$ and $\{V_2^1, V_2^2, \dots, V_2^m\}$, in networks G_1 and G_2 , respectively. Note that other network clustering methods can be used in this step alternatively. Moreover, here we assume that, the number of clusters (blocks) m is known. Otherwise, it can be learned using Bayesian [218] or Monte Carlo [219] techniques. In the module alignment step, we use the EigenAlign solution X^* to infer a mapping across inferred modules. To do this, we form a bipartite module-network whose nodes in the first and second layers correspond to $\{V_1^1, V_1^2, \dots, V_1^m\}$, and $\{V_2^1, V_2^2, \dots, V_2^m\}$,

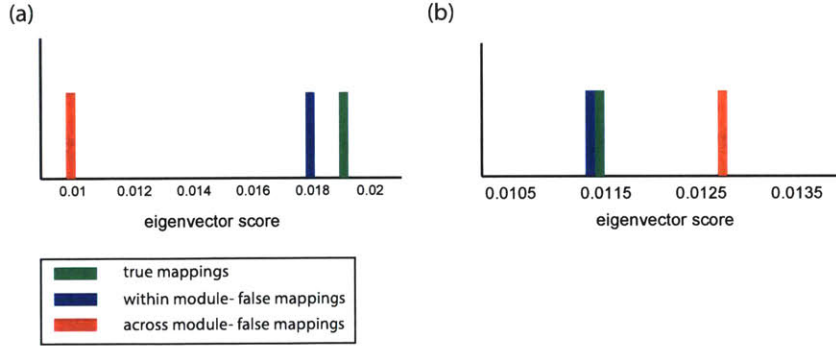


Figure 5-4: Eigenvector scores of within and across module mappings of the expected alignment network with parameters $n = 20$, $m = 4$, $p = 0.3$, $q = 0.01$, and (a) $\alpha = 100$, (b) $\alpha = 10$. Parameters of panel (a) satisfy conditions of Theorem 5.25.

respectively (Figure 5-3-b). The edge weight connecting node i in the first layer of this graph to node j' in the second layer (i.e., $w(i, j')$) is computed as follows:

$$w(i, j') = \sum_{\substack{a \in V_1^i \\ b \in V_2^{j'}}} X^*(a, b). \quad (5.37)$$

We use these weights in a maximum bipartite matching optimization to infer a one-to-one mapping across modules of two networks. Other methods based on random walks over aggregation of Markov chains [220] can be used in this step to find module mappings. In the last step, we use an SDP-based relaxation of the underlying QAP [185] along with a linear programming step of maximum weight bipartite matching to compute node alignments in each module-pair. In this step, other methods based on convex/semidefinite relaxations of the underlying QAP can be used as well.

In the following, we prove that the EigenAlign solution provides a robust module-level mapping across networks even in the high-noise regime.

Define,

$$\begin{aligned} \mathcal{M}_{in} &\triangleq \{(i, j') : i \in V_1^a, j' \in V_2^a, 1 \leq a \leq m\} \\ \mathcal{M}_{across} &\triangleq \{(i, j') : i \in V_1^a, j' \in V_2^b, 1 \leq a, b \leq m, a \neq b\}, \end{aligned} \quad (5.38)$$

where \mathcal{M}_{in} and \mathcal{M}_{across} represent mappings within and across modules in two networks, respectively. We wish to show that eigenvector centrality scores assigned to \mathcal{M}_{in} mappings in the EigenAlign algorithm are strictly larger than the ones assigned to \mathcal{M}_{across} mappings. Suppose G_1 and G_2 have m blocks, each with size n . Thus, the total number of nodes in each network is mn . Let A represent the alignment network across G_1 and G_2 , according to (5.3). This network has m^2n^2 number of nodes, out of which mn^2 are within module mappings (i.e., $|\mathcal{M}_{in}| = mn^2$), and the rest are across module mappings (i.e., $|\mathcal{M}_{across}| = (m^2 - m)n^2$). Let \bar{A} be the expected alignment matrix, where $\bar{A}(t_1, t_2) = \mathbb{E}[A(t_1, t_2)]$. To highlight key optimality conditions, in this section, we focus on the eigenvector analysis of the matrix \bar{A} . Similarly to Section 5.2.3, these properties can be extended to the alignment network A by using perturbation analysis of random matrices.

We choose scores assigned to matches, neutrals and mismatches as $s_1 = \alpha + \epsilon$, $s_2 = 1 + \epsilon$ and $s_3 = \epsilon$, respectively, where $\alpha > 1$ and $0 < \epsilon \ll 1$. To simplify notation, we assume ϵ is sufficiently small so that its effect is negligible. Let v correspond to the leading eigenvector of the expected alignment matrix \bar{A} .

Theorem 5.25. For the noiseless case ($p_e = 0$), if $p > q$, $\alpha > 1/q - 1$, and $m > 2$, then,

$$v(i) > v(j), \quad \forall i \in \mathcal{M}_{in}, \forall j \in \mathcal{M}_{across}. \quad (5.39)$$

Proof. A proof is presented in Section 5.8.5. □

The condition $p > q$ implies that the connectivity density of nodes in modules are larger than the one across modules. This condition is also essential in the spectral partitioning step of the EigenAlign+SDP algorithm to infer modules reliably. The condition $m > 2$ is a technical condition required in the proof of the strict inequality of (5.39). In Section 5.5.3, through simulations, we show that this condition is not essential for the algorithm. The condition $\alpha > 1/q - 1$ guarantees that the expected alignment score of two mapping pairs where one belongs to \mathcal{M}_{in} and the other one

belongs to \mathcal{M}_{across} is strictly larger than the one where both mappings belong to \mathcal{M}_{across} . More details on these optimality conditions can be found in Section 5.8.5.

Now, we consider the case when G_2 is related to a noisy version of G_1 :

Theorem 5.26. Suppose $0 \leq p_e^2 \lambda_1$ and $0 \leq p_{e_2}^2 \lambda_1$. Let $p > q$, $m > 2$ and $\alpha \Gamma 1/q$.

- If $q \leq \frac{1}{4}$ and $p < \frac{1}{1+p_e}$,
- or, if $q > \frac{1}{4}$ and $p < \min\left(\frac{1}{1+p_e}, \frac{6q - \sqrt{4q(4q-1)(1-q)}}{2(1+2q)}\right)$,

then,

$$v(i) > v(j), \quad \forall i \in \mathcal{M}_{in}, \forall j \in \mathcal{M}_{across}. \quad (5.40)$$

Proof. A proof is presented in Section 5.8.6. □

If $p_e = 0$, the condition $p < \frac{1}{1+p_e}$ is simplified to $p < 1$. For higher noise level, this condition guarantees that correct within module mappings are assigned to higher scores than incorrect within-module ones. The additional constraint $p < \frac{6q - \sqrt{4q(4q-1)(1-q)}}{2(1+2q)}$ is required to guarantee that, the expected alignment score of two mapping pairs, one in \mathcal{M}_{in} and the other one in \mathcal{M}_{across} , is strictly larger than the one where both mappings belong to \mathcal{M}_{across} . We illustrate Theorems 5.25 and 5.26 in Figure 5-4. More details on these sufficient optimality conditions can be found in Section 5.8.6.

Theorems 5.25 and 5.26 indicate that eigenvector scores of within module mappings are strictly larger than the ones of across module mappings, in the expected alignment network. However, definitions of \mathcal{M}_{in} and \mathcal{M}_{across} are based on the permutation matrix P which relates networks G_1 and G_2 according to (5.17). In general, P may not be necessarily the optimal solution of the network alignment optimization (5.5). In the following, we provide conditions where the permutation matrix P in fact maximizes the expected objective function of the network alignment optimization (5.5).

Theorem 5.27. Under the conditions of Theorem 5.25, if

$$p \leq \frac{1 + \sqrt{1 + (\alpha^2 - 1)q}}{1 + \alpha}, \quad (5.41)$$

then, the permutation matrix P maximizes the expected objective function of the network alignment optimization (5.5).

Proof. A proof is presented in Section 5.8.7. \square

For instance, setting $\alpha = 1/q$ and $q\lambda 1$, the condition of Theorem 5.27 results to the condition $p^2 \leq q$.

Theorem 5.28. Under the conditions of Theorem 5.26, if

$$p^2 \leq \frac{q(1 - p_e)}{1 + p_e^2}, \quad (5.42)$$

then, the permutation matrix P maximizes the expected objective function of the network alignment optimization (5.5).

Proof. A proof is presented in Section 5.8.8. \square

Conditions of Theorems 5.27 and 5.28 are based on an upper bound on the parameter p . To explain this condition intuitively, suppose the permutation matrix P maps modules V_1^a in network G_1 to modules V_2^a in network G_2 , where $1 \leq a \leq m$. In the case of large number of modules, the expected alignment score of the permutation matrix P is dominated by mapping-pairs (V_1^a, V_2^a) and (V_1^b, V_2^b) , where $a \neq b$. These scores depend on the connectivity density across modules, namely q . On the other hand, consider an incorrect permutation matrix \tilde{P} where $\frac{1}{nm} \|P - \tilde{P}\| > 0$. The alignment score of \tilde{P} is upper bounded by scores of incorrect within module mappings which depend on p . To have a sufficient optimality condition, we wish the expected alignment score of P to be larger than the upper bound of the expected alignment score of \tilde{P} . More details on these optimality conditions are presented in Sections 5.8.7 and 5.8.8.

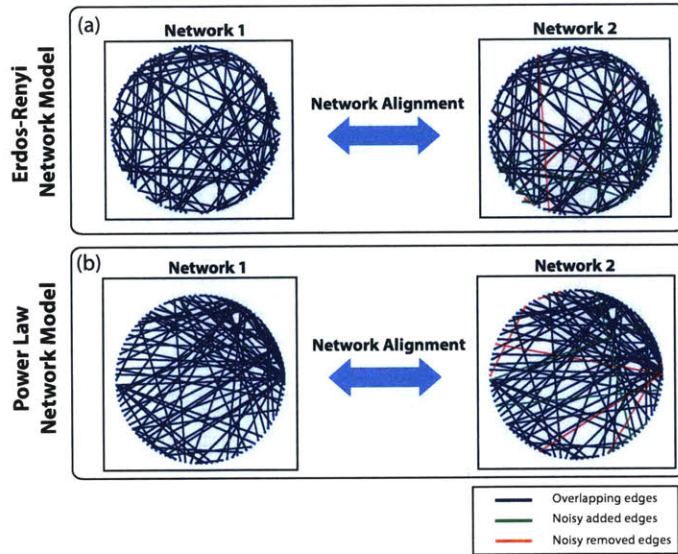


Figure 5-5: Examples of (a) Erdős-Rényi, and (b) power law networks used in Section 5.5.1.

In the EigenAlign+SDP algorithm 5.24, errors can also happen in other steps of the algorithm including the network partitioning and semidefinite relaxation steps. In Section 5.5.3, we empirically analyze the performance of the proposed algorithm. Our results suggest that the proposed method significantly outperforms existent network alignment techniques in aligning modular networks, while it has significantly lower computational complexity compared to standard convex-relaxation methods.

5.5 Performance Evaluation Over Synthetic Networks

We now compare the performance of the EigenAlign algorithm against other network alignment methods including IsoRank [161], NetAlign [188], Klau linearization [204] as well as an SDP-based method [185] through simulations. IsoRank is a global network alignment method which uses an iterative approach to align nodes across two networks based on their neighborhood similarities. NetAlign formulates the alignment problem in a quadratic optimization framework and uses message passing to approximately solve it. Klau linearization uses Lagrange multipliers to relax the underlying quadratic assignment problem. The SDP-based method [185] uses a

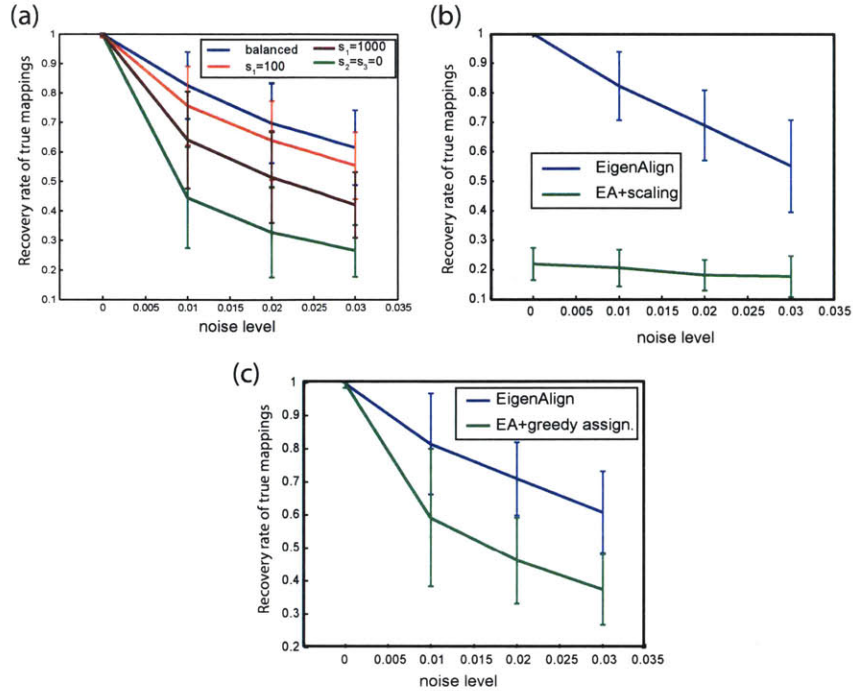


Figure 5-6: EigenAlign (EA) performance over power law graphs (a) with different alignment scores, (b) with alignment network scaling, and (c) with the greedy assignment method. At each point, simulations have been repeated 20 times.

convex relaxation of the underlying QAP based on matrix splitting. In our simulations, we use default parameters of these methods.

5.5.1 Erdős-Rényi and Power-Law graphs

Here, we test alignment methods over two types of input networks:

- *Erdős-Rényi graphs [199]*: In this case, G_1 is a symmetric random graph where, $Pr[(i, j) = 1] = p$ (see an example in Figure 5-5-a).
- *Power law graphs [221]*: We construct G_1 as follows; we start with a random subgraph with 5 nodes. At each iteration, a node is added to the network connecting to θ existent nodes with probabilities proportional to their degrees. This process is repeated till the number of nodes in the network is equal to n (see an example in Figure 5-5-b).

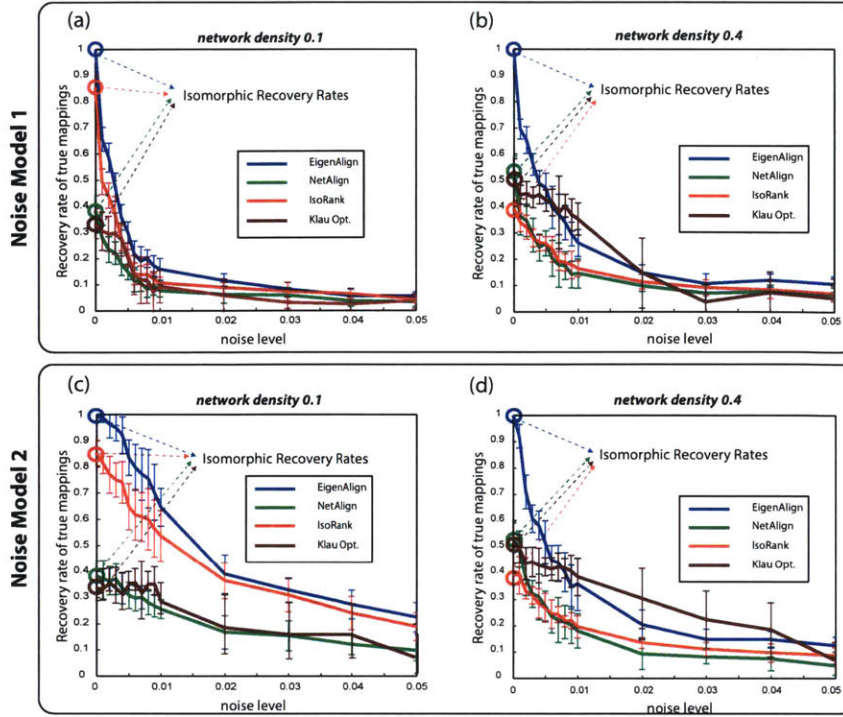


Figure 5-7: Performance evaluation of alignment methods over Erdős-Rényi graphs. EigenAlign outperforms IsoRank, NetAlign and Klau optimization in low and high network densities, in all considered noise levels, in both noise models. At each point, simulations have been repeated 10 times.

G_2 is then formed according to (5.17). In our experiments, we consider two noise models (5.14) and (5.15). In the case of the power law network model, we use the density of G_1 as parameter p in the noise model II of (5.15). We denote p_e as the *noise level* in both models.

In (5.17), P can be an arbitrary permutation matrix. In our simulations, we choose P such that, $P(i, n - i + 1) = 1$, for all $1 \leq i \leq n$, where n is the number of nodes in the network. Let \tilde{P} represent the solution of a network alignment method in aligning networks G_1 and G_2 . The recovery rate of true mappings is defined as,

$$1 - \frac{1}{2n} \|P - \tilde{P}\|. \quad (5.43)$$

As we explain in Proposition 5.16, this recovery rate is directly related to the number of inferred matches and mismatches in the network alignment optimization (5.5).

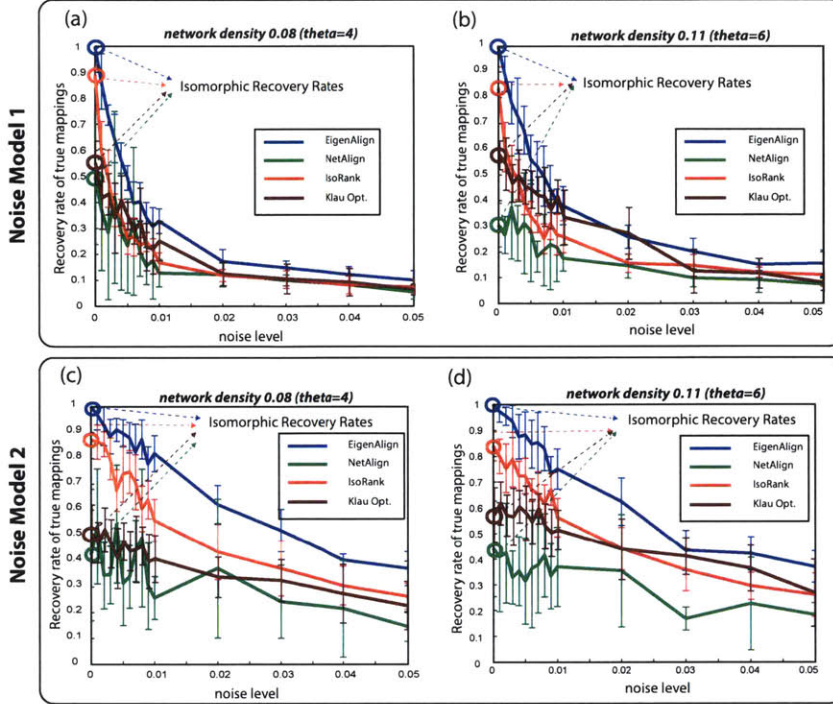


Figure 5-8: Performance evaluation of alignment methods over power law graphs. EigenAlign outperforms IsoRank, NetAlign, and Klau optimization in low and high network densities, in all considered noise levels, in both noise models. At each point, simulations have been repeated 10 times.

Thus, we illustrate performance evaluations using this metric, noting that using other metrics such as the number of inferred matches/mismatches results in similar performance trends and conclusions. Moreover, in our simulations, all mappings across networks are considered to be possible (i.e., an unrestricted network alignment case).

Figure 5-6 illustrates individual contributions of different components of the EigenAlign algorithm by isolating multiple aspects of the algorithm. In these experiments, we use the power-law network structure with 50 nodes, $\theta = 3$, and noise model II. Performance trends and conclusions are similar for other cases as well. In EigenAlign Algorithm 5.9, we use $s_1 = \alpha + \epsilon$, $s_2 = 1 + \epsilon$, and $s_3 = \epsilon$, where $\epsilon = 0.001$ and,

$$\alpha = 1 + \frac{\# \text{ of mismatches}}{\# \text{ of matches}}. \quad (5.44)$$

This choice of α satisfies the required condition $\alpha > 1$. Moreover, by this selection of

α , we have,

$$|s_1 - 1 - \epsilon|(\# \text{ of matches}) = |s_3 - 1 - \epsilon|(\# \text{ of mismatches}), \quad (5.45)$$

which makes a balance between matched and mismatched interaction scores across networks and leads to an improved performance of the method (Figure 5-6-a). Note that, ignoring mismatches ($s_2 = s_3 = 0$) leads to a significant decrease in the EigenAlign performance. Choosing other values of α can provide a way to adjust relative importance of matching and non-matching interactions in different applications. In general, this parameter can be tuned in different applications using standard machine learning techniques such as cross validations [222]. Moreover, Figure 5-6-b illustrates the performance of the EigenAlign algorithm if we scale the alignment network by its node degrees (we replace $A_{i,j}$ by $A_{i,j}/d_i d_j$, where d_i is the degree of node i). As illustrated in this figure, this scaling removes our global match-mismatch score assignments and leads to a significant decrease in the performance. It is worth noting that, other scaling/scoring methods may be deployed leading to high network alignment performance in different applications. Finally, Figure 5-6-c illustrates the performance of the EigenAlign algorithm if we replace the linear optimization of maximum weight bipartite matching with the greedy approach of Remark 5.14. As illustrated in this figure, this variation of the EigenAlign method decreases the performance, although it also decreases the computational complexity.

Next, we compare the performance of the EigenAlign method with the one of other network alignment methods over synthetic networks with $n = 100$ nodes. For large networks, the SDP-based method of [185] has high computational complexity, thus we excluded it from these experiments. Later, we will investigate the performance of the SDP-based method over small networks ($n \leq 50$).

Figure 5-7 shows the recovery rate of true mappings obtained by EigenAlign, IsoRank, NetAlign and Klau relaxation, over Erdős-Rényi graphs, and for both noise models. In the noiseless case (i.e., $p_e = 0$), the network alignment problem is simplified to the problem of graph isomorphism as explained in Section 5.2.3. In this case, the

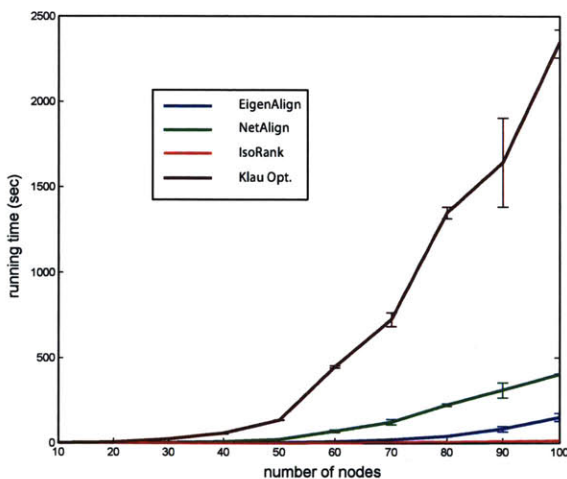


Figure 5-9: Running time of network alignment methods over Erdős-Rényi graphs.

EigenAlign algorithm finds true mappings without error in both low and high network densities (i.e., the recovery rate of (5.43) is one.). In noisy cases and in both low and high network densities, EigenAlign shows consistently high performance compared to other methods.

Figure 5-8 shows the recovery rate of true mappings obtained by three considered methods over power law graphs, and for both noise models. Similarly to the case of Erdős-Rényi graphs, EigenAlign outperforms other methods in all considered noise levels and network densities. Notably, in noiseless cases, EigenAlign finds true isomorphic mappings across networks without error.

Figure 5-9 depicts the average running time of these methods over an Erdős-Rényi graph with density $p = 0.4$, and under the noise model I. All methods have been run on the same machine and each experiment has been repeated 10 times. In these experiments, EigenAlign has the second best runtime, after IsoRank, outperforming NetAlign and Klau relaxation methods. Notably, the running time of Klau optimization method increases drastically as network size grows.

Next, we consider small networks with $n = 50$ nodes so that we are able to include a computationally expensive SDP-based method of [185] in our experiments. Note that, the method proposed in [185] only provides a bound on the value of the underlying QAP objective function and it does not provide a feasible solution. In

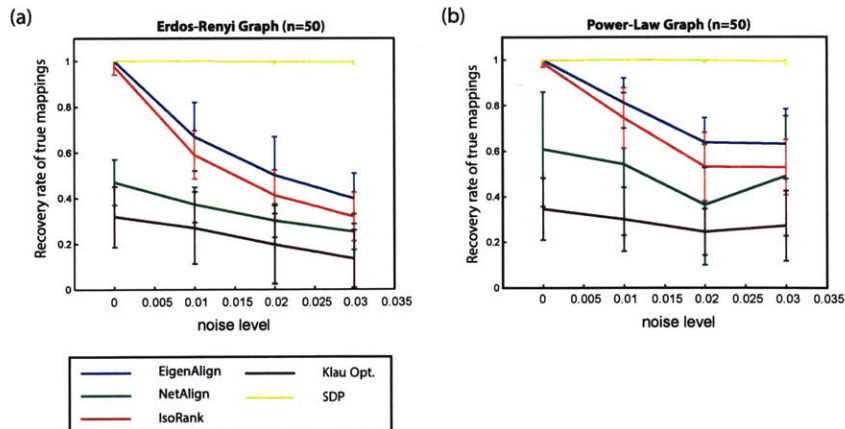


Figure 5-10: Performance evaluation of network alignment methods over (a) Erdős-Rényi ($p = 0.2$), and power law ($\theta = 3$) graphs, with $n = 50$. The computationally expensive SDP-based method outperforms other methods, while EigenAlign has the second best performance. At each point, simulations have been repeated 20 times.

order to obtain a feasible solution for the network alignment optimization, we use the SDP-based solution in a maximum weight bipartite matching optimization.

Figures 5-10-a,b illustrate the recovery rate of true mappings obtained by five considered methods over both Erdős-Rényi and power law graphs, under the noise model II. Performance trends under the noise model I is similar. As illustrated in these figures, the computationally expensive SDP-based method outperforms other network alignment methods significantly, while it has the highest running time (Figure 5-11). Notably, in these experiments, EigenAlign outperforms other methods except the SDP-based one consistent with the cases illustrated in Figures 5-7 and 5-8. These results inspired us to propose the EigenAlign+SDP method to both have high recovery rate and low computational complexity (see Section 5.4).

5.5.2 Low Rank Alignment of Networks

In this section, we evaluate the performance of LowRank Align algorithm introduced in Section 5.3. Unlike EigenAlign, LowRank Align uses top k eigenvectors of input graphs according to Optimization (5.31). Here we use $k = 4$ as larger values of k did not have a significant effect on the results. We compare the performance of

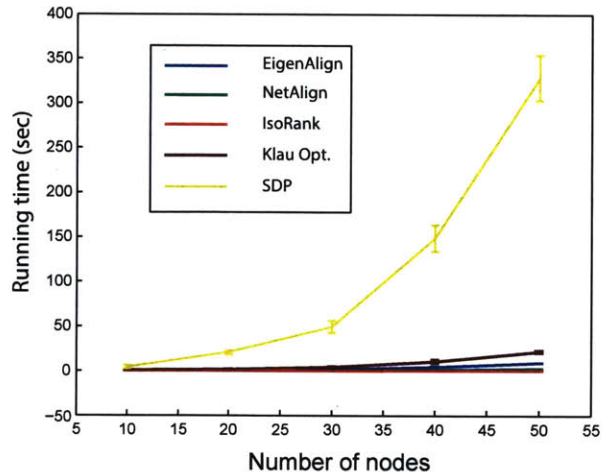


Figure 5-11: Running time of network alignment methods over Erdős-Rényi graphs with $n = 50$ and $p = 0.2$. The SDP-based method has significantly higher computational complexity compared to other network alignment techniques, which prohibits its use over large graphs. At each point, simulations have been repeated 20 times.

LowRank Align with EigenAlign, NetAlign [188], IsoRank [161] and an SDP-based network alignment method [185]. Besides these methods, we consider a network alignment solution that projects an optimal orthogonal relaxation of the underlying QAP over permutation matrices according to Theorem 5.19 and Optimization (5.22) (we refer to this method as OR, for orthogonal relaxation). As we explain in Section 5.3, there may be exponentially many optimal solutions for orthogonal relaxation of the underlying QAP. We also consider the case where all eigenvectors are oriented to have more positive components than negative ones (we refer to this method as ORP, for orthogonal relaxation with positive orientation of eigenvectors). Moreover, we consider a version of ORP where eigenvectors are scaled with eigenvalues similarly to (5.28) (we refer to this method as ORPS, for ORP with eigenvalue scaling).

Figure 5-12 shows the recovery rate, the number of inferred matches and mismatches and the mismatch-match ratio for different methods over Erdős-Rényi graphs with $n = 50$ nodes and edge density of $p = 0.2$. Similarly to our results in Section 5.5.1, the SDP-based network alignment outperforms other methods. However, this method has significantly higher computational complexity compared to other methods which prohibits its use in the alignment of large graphs. As it is illustrated

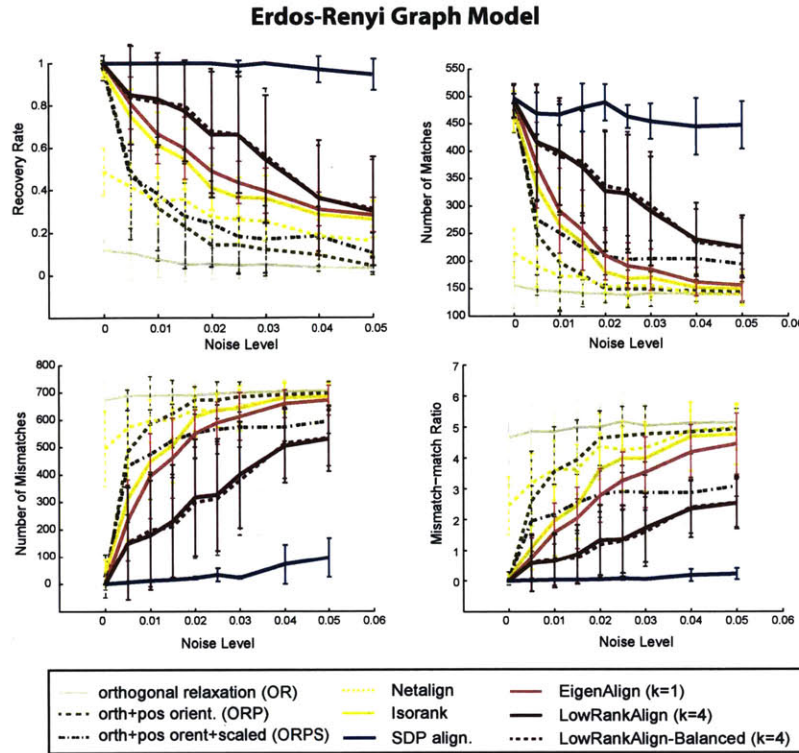


Figure 5-12: Performance evaluation of different network alignment methods including LowRank Align over Erdős-Rényi graphs with $n = 50$ nodes and edge density of $p = 0.2$. Experiments have been repeated 100 times for all methods except for the SDP-based alignment. For the SDP-based alignment, experiments have been repeated five times owing to its high computational complexity.

in this figure, LowRank Align outperforms all other methods (except the expensive SDP-based method) in different noise levels. This illustrates that using higher-order eigenvectors in aligning graph structures leads to a tighter approximation of the underlying QAP. Notably, ORPS and ORP both outperform the baseline orthogonal relaxation (OR) method, highlighting the importance of considering eigenvector orientations and scaling in the network alignment optimization described in Section 5.3. Similar results hold for the alignment of power-law graph structures.

Figure 5-13 shows performance evaluation of different methods over regular graph structures with $n = 50$ nodes and edge density of $p = 0.2$. The alignment of regular graph structures is one of the most difficult network alignment cases because of the homogeneity of node degrees. As illustrated in Figure 5-13, in this case, LowRank

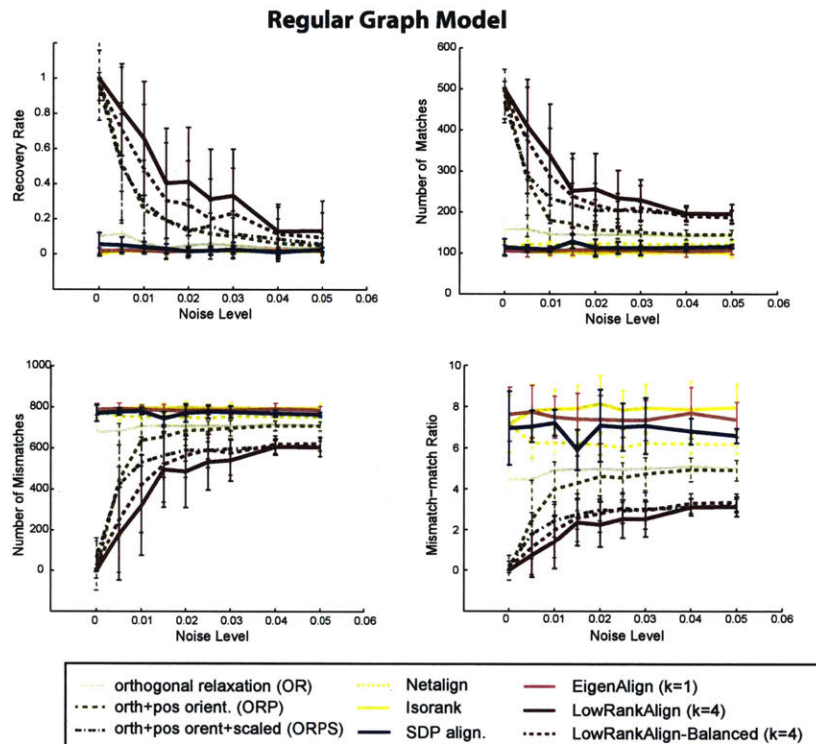


Figure 5-13: Performance evaluation of different network alignment methods including LowRank Align over regular graphs with $n = 50$ nodes and edge density of $p = 0.2$. Experiments have been repeated 100 times for all methods except for the SDP-based alignment. For the SDP-based alignment, experiments have been repeated five times owing to its high computational complexity.

Align outperforms all other methods including the expensive SDP-based network alignment. This illustrates the effectiveness of using higher order eigenvectors in aligning homogenous graph structures.

5.5.3 Modular Network Structures

We then assess the performance of the proposed EigenAlign+SDP method in Section 5.4 in aligning networks with modular structures. To be able to compare the performance of the proposed method with the one of the SDP-based method [185], we only consider small modular networks with 50 nodes, having two equal-size modules. As discussed in Section 5.5.1, the SDP-based method of [185] has high computation-

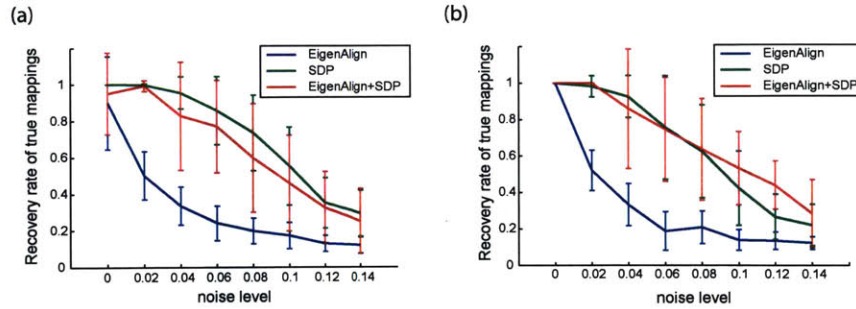


Figure 5-14: Performance evaluation of network alignment methods over modular network structures with $n = 20$, $m = 2$, $p = 0.2$, and (a) $q = 0.05$, (b) $q = 0.01$. The recovery rate of the EigenAlign+SDP method is close to the one of the SDP-based method, while it has significantly lower computational complexity.

ally complexity and using it for large networks is not practically feasible. The key idea of the EigenAlign+SDP method is to use the EigenAlign solution to split the large QAP into smaller sub-problems, enabling the use of the SDP-based relaxation method over each sub-problem. Moreover, the SDP-based method can be run in parallel over each sub-problem.

Here, we consider network and noise models described in Section 5.4. In our evaluations, we consider EigenAlign, SDP, and EigenAlign+SDP methods. Moreover, we use the performance evaluation metric introduced in Section 5.5.1. Figures 5-14-a,b illustrates the recovery rate of true mappings of different methods, in various noise levels, and for different set of parameters p and q (the density of edges within and across modules, respectively). The average running time of the methods is depicted in Figure 5-15. As it is illustrated in these figures, while the recovery rate of the EigenAlign+SDP method is close to the one of the SDP-based one, it has significantly lower computational complexity, enabling its use for large complex networks.

5.6 Performance Evaluation Over Gene Networks

Having illustrated the effectiveness of our algorithm both theoretically and through simulations, we apply it to compare gene regulatory networks across human, fly and worm species. We use regulatory networks which we inferred in Chapter 3 by inte-

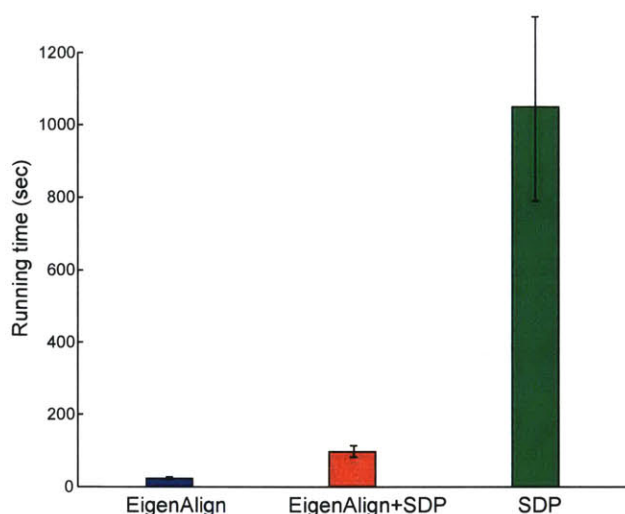


Figure 5-15: Average running time of network alignment methods over modular network structures with $n = 20$, $m = 2$, $p = 0.2$, and $q = 0.05$. The average running time of the SDP+EigenAlign algorithm is significantly lower than the one of the SDP-based method. Simulations have been repeated 20 times.

grating genome-wide functional and physical genomics datasets from ENCODE and modENCODE consortia. Comparative network analysis in evolutionary studies often requires having a one-to-one mapping across genes of two or multiple networks. We use homolog mappings across genes of human, fly and worm provided by the ENCODE and modENCODE consortia [223]. Note that human homologs refer to homologous genes in fly and worm, while fly/worm homologs are defined solely in relation to human. Homolog mappings across species are based on gene sequence similarities over corresponding gene trees [223]. However, since human, fly and worm are distant species and as a result, many gene families have undergone extensive duplications and losses, we observe non-bijective homolog mappings across their genes [223]. For example, one gene in human can be homologous to multiple genes in fly and vice versa (see an example in Figure 5-16). To infer bijective mappings as a subset of homolog genes across species, we use network alignment methods.

We assess the performance of three network alignment methods (EigenAlign, NetAlign, IsoRank) in comparative analysis across gene regulatory networks of human, fly and worm. We excluded the network alignment method based on Klau optimiza-

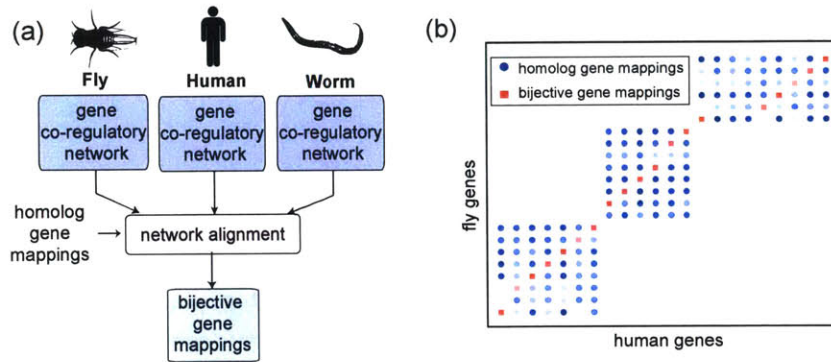


Figure 5-16: (a) A comparative analysis framework across human, fly and worm regulatory networks. (b) An application example of a network alignment method (EigenAlign) over three gene families across human-fly.

tion [204] from our analysis in this section owing to its high computational complexity, and its low performance over synthetic networks of Section 8.3. Moreover, as we discussed in Section 8.3, the SDP-based method of [185] has high-computational complexity which prohibits its use in aligning large regulatory networks. We use EigenAlign, IsoRank and NetAlign methods with the setup and parameters similarly to Section 8.3. Moreover, unlike EigenAlign, IsoRank and NetAlign methods do not take into account the directionality of edges in their network alignment setup. Thus, to have fair performance assessments of considered network alignment methods, we create un-directed co-regulatory networks using inferred regulatory networks of Chapter 3, by connecting genes when their parent TFs have an overlap larger than 25%. This results in undirected binary networks in human, fly, and worm, with 19,221, 13,642, and 19,296 nodes, and 13.9%, 3.5%, and 4.2% edge densities, respectively.

By application of network alignment methods, we infer bjective mappings across human-fly and human-worm networks. There are numerous mismatched interactions across networks of distal species partially owing to extensive gene functional divergence due to processes such as gene duplication and loss. In this situation, it is critical to consider both matches and mismatches in the network alignment optimization. Unlike existent network alignment algorithms which completely ignore mismatches across networks, EigenAlign considers both matches and mismatches across networks

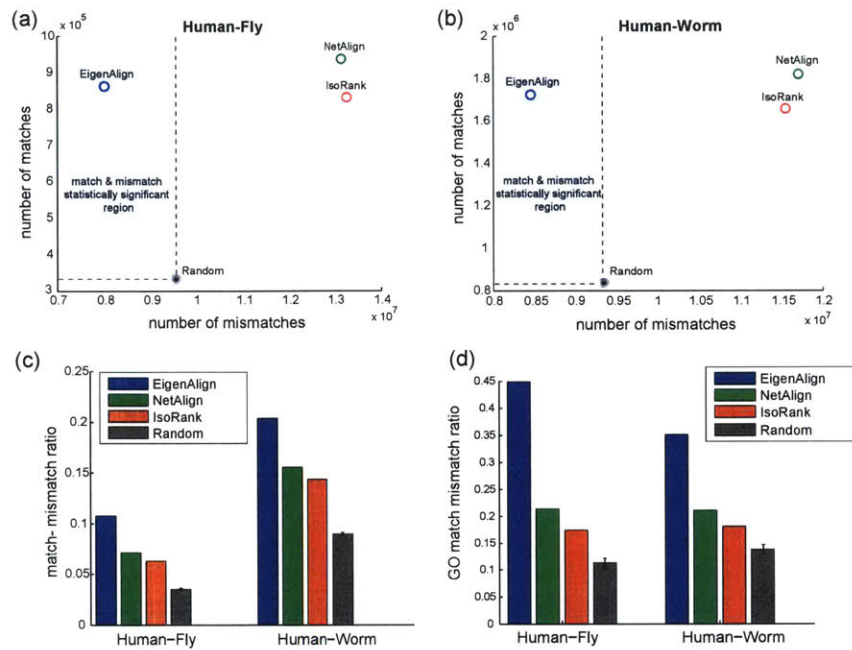


Figure 5-17: The number of matched and mismatched interactions inferred using different network alignment methods across (a) human-fly and (b) human-worm networks. (c) An average match-mismatch ratio for different network alignment methods, in both human-fly and human-worm comparisons. (d) This panel shows average match-mismatch ratios over pairs of biological processes across both human-fly and human-worm networks. EigenAlign mappings preserve regulatory patterns of biological processes across both human-fly and human-worm networks significantly better than the ones of NetAlign and IsoRank methods.

in a balanced way by setting alignment scores according to (5.44). Figures 5-17-a,b illustrate the number of matched (conserved) and mismatched (error) interactions across human-fly and human-worm networks, inferred using different methods. As it is illustrated in these figures, in both human-fly and human-worm comparisons, EigenAlign significantly outperforms other methods in terms of causing fewer number of mismatches (errors) across networks, while its performance is close to the best one in terms of the number of matches (overlaps). Figure 5-17-c illustrates the match-mismatch ratio for different network alignment methods, in both human-fly and human-worm comparisons. As it is illustrated in this figure, EigenAlign outperforms other methods significantly indicating that, it finds mappings across networks which maximize the number of matches and minimize the number of mismatches

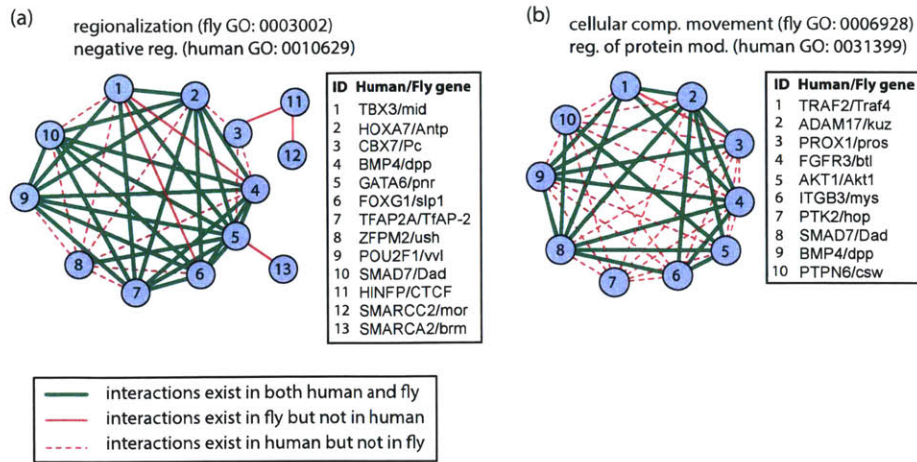


Figure 5-18: Examples of conserved biological processes (p -values < 0.05) across human and fly inferred by EigenAlign.

simultaneously. In all cases, the number of conserved interactions is statistically significant (p -value < 0.01), suggesting regulatory pathways are conserved between distal species and therefore highlighting their potential role in basic cellular processes. p -values are computed using a rank test by randomly selecting bijective mappings from gene homolog mappings while keeping network structures the same.

We examine enrichment of different biological processes over conserved subgraphs inferred by different alignment methods, across human-fly and human-worm. We expect many basic cellular processes such as transcription elongation, cellular growth, etc to be conserved across these species despite long evolutionary distances. In our comparative analysis, we consider genome ontology (GO) processes that are experimentally verified and have sizes in the range of 100 and 500 genes. We consider pairs of processes across species with more than 50 homolog gene pairs.

Figure 5-17-d shows average match-mismatch ratios over these pairs of biological processes across both human-fly and human-worm networks. As illustrated in this figure, mappings inferred by EigenAlign preserve regulatory patterns of biological processes significantly better than the ones of NetAlign and IsoRank methods, across both human-fly and human-worm networks. This is consistent with the overall match-mismatch ratio gain of the EigenAlign method across human-fly and human-worm

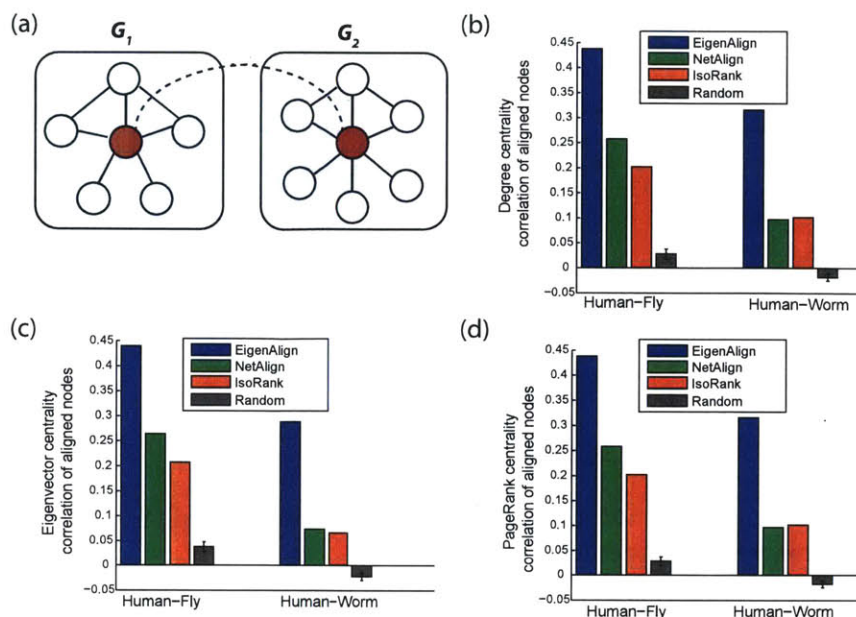


Figure 5-19: (a) A suitable alignment should map central nodes in one network to central nodes in the other network. Correlation of (b) degree centrality, (c) eigenvector centrality, and (d) page-rank centrality measures of aligned genes using different network alignment methods, in both human-fly and human-worm comparisons.

networks, demonstrated in Figure 5-17-c. Note that the average match-mismatch ratios across GO processes (Figure 5-17-d) are significantly higher than the ones across entire networks (Figure 5-17-c). This may indicate that connectivity patterns among genes involved in similar processes are better preserved compared to the rest of the regulatory network. However, note that, this enrichment of GO processes across species may be partially owing to biases in the study of these genes in the literature [224]. Thus, additional experimental validations are required which is beyond the scope of this work.

Next, we compute p -values for matches and mismatches for each pair of processes using a rank test by randomly selecting bijective mappings from gene homolog mappings. Then, we take the worst match/mismatch p -value for each pair and perform Benjamini and Hochberg multiple hypothesis correction [88]. Using EigenAlign mappings, we find five conserved pairs of regulatory processes across human-fly networks (p -values < 0.05) including regionalization/negative regulation of gene

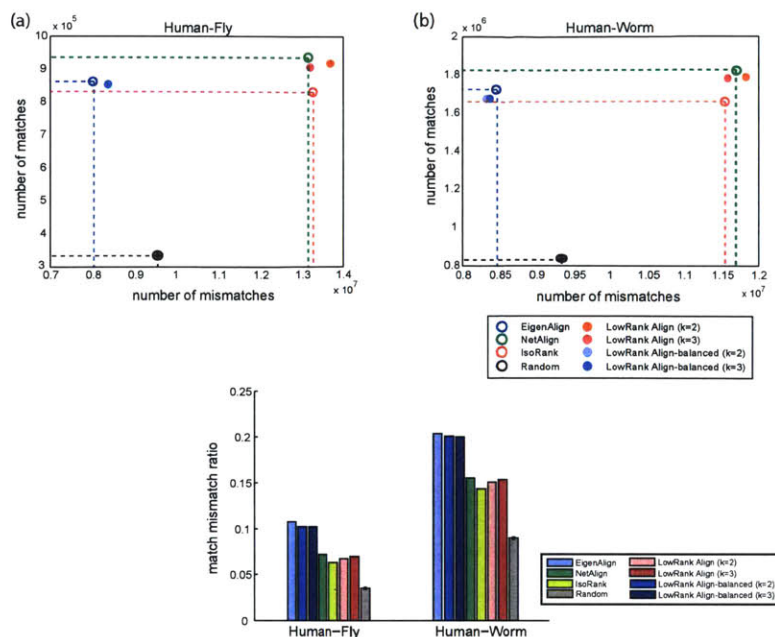


Figure 5-20: Performance evaluation of different network alignment methods including LowRank Align over co-regulatory networks. EigenAlign and LowRank Align both lead to similar solutions.

expression (Figure 5-18-a), and cellular component movement/regulation of protein modification process (Figure 5-18-b). Using NetAlign and IsoRank mappings, no significantly enriched processes across human and fly are found. Moreover, across human and worm, none of the inferred mappings leads to the inference of significantly conserved processes. Similarly to the centrality conservation analysis, this can be partially owe to a large evolutionary distance between human and worm. Finally, we note that, the solution provided by the EigenAlign method can be used in designing future experiments to extend GO catalogs by using gene regulatory networks and cross-species information.

We assess conservation of gene centralities across different networks. Centrality measures over regulatory networks provide useful insights on functional roles of genes in regulatory processes [225–227]. Therefore, comparative analysis of gene centrality measures across species can shed light on how evolution has preserved or changed regulatory patterns and functions across distal species. A suitable alignment should

map central nodes in one network to central nodes in the other network (Figure 5-19-a). Here, we consider three centrality measures: degree centrality, eigenvector centrality, and the page-rank centrality with the damping factor 0.85. In Figure 5-19, we illustrate centrality correlation of aligned genes using different network alignment methods, across both human-fly and human-worm networks. As it is illustrated in this figure, the bijective mapping inferred by EigenAlign preserves different node centralities across networks significantly better than other tested methods. That is, EigenAlign is more likely to map central nodes in one network to central nodes in the other network. Finally, we note that, the centrality conservation across human-fly networks is significantly larger than the one across human-worm networks, partially owing to a larger evolutionary distance between human and worm compared to the one between human and fly.

Note that in the alignment of regulatory networks across human, fly and worm, both EigenAlign and its higher eigenvector extension (LowRank Align) lead to similar solutions, suggesting that in our particular example of aligning regulatory networks of distal species, the leading eigenvectors yield a reasonable approximation of the underlying QAP (Figure 5-20).

5.7 User De-anonymization over twitter follower sub-graphs

Data anonymization is the process of either encrypting or removing personally identifiable information from a dataset. As more and more personal information is available over social networks such as Facebook and twitter, data anonymization is becoming critical in privacy protection of people [171, 228]. Moreover, because the size and resolution of medical and genetic datasets have increased significantly, serious challenges have arisen with respect to ensuring the privacy of patients and samples [229]. In this circumstance, it is critical to measure the vulnerability of data anonymization methods by aiming to extract personal information from an

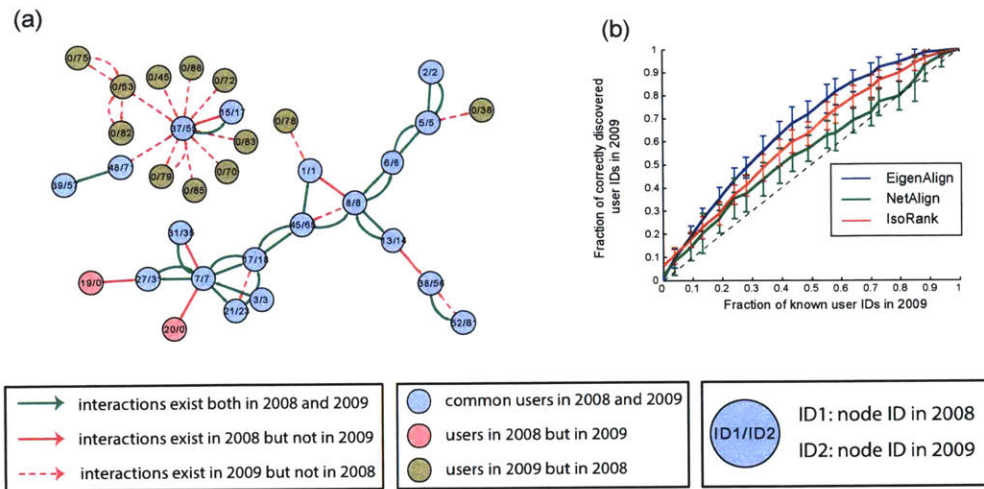


Figure 5-21: De-anonymization of user IDs over twitter subgraphs. We consider small subgraphs of the twitter follower network sampled in 2008 and 2009 years, with 53 and 88 nodes, and 58 and 91 edges, respectively. Panel (a) illustrates two large components of these networks. Panel (b) illustrates the fraction of correctly de-anonymized users using different network alignment methods. The EigenAlign algorithm is more efficient than both NetAlign and IsoRank methods in user de-anonymization over these graphs.

anonymized dataset to highlight its privacy loopholes.

Here, we examine the power of network alignment techniques in de-anonymizing user IDs using connectivity structures of networks. We consider small subgraphs of the Twitter follower network sampled in years 2008 and 2009. These networks have 53 and 88 nodes, and 58 and 91 edges, respectively (Figure 5-21-a).

Here, we wish to de-anonymize user IDs in the 2009 network using user IDs in the 2008 network by only knowing a fraction of user IDs in 2009. To do that, we apply different network alignment methods to infer an alignment across nodes of two networks. Figure 5-21-b illustrates the fraction of correctly de-anonymized users using different network alignment methods. The EigenAlign algorithm is more efficient than both NetAlign and IsoRank methods in user de-anonymization over these graphs. This application illustrates the extent of personal information that can be retrieved from network structure alone, and raises additional considerations that need to be addressed in different privacy-related applications. Applications of network alignment techniques over the entire twitter graph will require reducing their computational

complexity significantly using further relaxations and approximations that we explain in Section 5.3 (Remark 5.22).

5.8 Proofs

In this section, we present proofs of the main results of the chapter.

5.8.1 Proofs of Lemmas 5.6 and 5.7

First we prove Lemma 5.6: Let A be the alignment network of graphs G_1 and G_2 . Suppose \tilde{P} is a permutation matrix where $\rho \triangleq \frac{1}{2n}\|P - \tilde{P}\| > 0$. Let \mathbf{y} and $\tilde{\mathbf{y}}$ be vectorized versions of permutation matrices P and \tilde{P} , respectively. Then, we have,

$$\begin{aligned}
 \frac{1}{n^2}\mathbb{E}[\tilde{\mathbf{y}}_2^T A \tilde{\mathbf{y}}_2] &= (1 - \rho)[ps_1 + (1 - p)s_2] + \rho[p^2s_1 + (1 - p)^2s_2 + 2p(1 - p)s_3] \quad (5.46) \\
 &< (1 - \rho)[ps_1 + (1 - p)s_2] + \rho[p^2s_1 + (1 - p)^2s_2 + p(1 - p)(s_1 + s_2)] \\
 &= (1 - \rho)[ps_1 + (1 - p)s_2] + \rho[ps_1 + (1 - p)s_2] \\
 &= ps_1 + (1 - p)s_2 \\
 &= \frac{1}{n^2}\mathbb{E}[\mathbf{y}_1^T A \mathbf{y}_1].
 \end{aligned}$$

Now we prove Lemma 5.7: Similarly to the proof of Lemma 5.6, let A be the alignment network of graphs G_1 and G_2 and suppose \tilde{P} is a permutation matrix where $\rho \triangleq \frac{1}{2n}\|P - \tilde{P}\| > 0$. Let \mathbf{y} and $\tilde{\mathbf{y}}$ be vectorized versions of permutation matrices P and \tilde{P} , respectively. Define a' and b' as follows:

$$\begin{aligned}
 a' &\triangleq p(1 - q)s_1 + (1 - p)(1 - q)s_2 + (pq + (1 - p)q)s_3, \quad (5.47) \\
 b' &\triangleq (p^2(1 - q) + pq(1 - p))s_1 \\
 &\quad + ((1 - p)^2(1 - q) + pq(1 - p))s_2 \\
 &\quad + (2p(1 - p)(1 - q) + 2p^2q)s_3.
 \end{aligned}$$

Thus,

$$a' - b' = p(1 - p)(1 - 2q)(s_1 + s_2 - 2s_3) + q(1 - 2p)s_3. \quad (5.48)$$

Because $s_1 > s_2 > s_3$, we have, $s_1 + s_2 - 2s_3 > 0$. Because $0 < p < 1/2$ and $0 \leq q < 1/2$, we have $(1 - 2p) > 0$ and $(1 - 2q) > 0$. Therefore, according to (5.48), $a' > b'$. Thus we have,

$$\frac{1}{n^2} \mathbb{E}[\tilde{\mathbf{y}}^T A \tilde{\mathbf{y}}] = (1 - \rho)a' + \rho b' < a' = \frac{1}{n^2} \mathbb{E}[\mathbf{y}^T A \mathbf{y}].$$

5.8.2 Proof Of Corollary 5.17

Without loss of generality and to simplify notations, we assume the permutation matrix P is equal to the identity matrix I , i.e., the isomorphic mapping across G_1 and G_2 is $\{1 \leftrightarrow 1', 2 \leftrightarrow 2', \dots, n \leftrightarrow n'\}$ (otherwise, one can relabel nodes in either G_1 or G_2 to have P equal to the identity matrix). Therefore, $G_1(i, j) = G_2(i', j')$ for all $1 \leq i, j \leq n$. Recall that \mathbf{Y} is a vector of length kn which has weights for all possible mappings $(i, j') \in \mathcal{R}$. To simplify notations and without loss of generality, we re-order indices of vector \mathbf{y} as follows:

- The first n indices of \mathbf{y} correspond to correct mappings, i.e., $y(1) = y_{1,1'}, y(2) = y_{2,2'}, \dots, y(n) = y_{n,n'}$.
- The remaining $(k - 1)n$ indices of \mathbf{y} correspond to incorrect mappings. e.g., $y(n + 1) = y_{1,2'}, y(n + 2) = y_{1,3'}, \dots, y(kn) = y_{r,s'}$ ($r \neq s$).

Therefore, we can write,

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix},$$

where \mathbf{y}_1 and \mathbf{y}_2 are vectors of length n and $(k - 1)n$, respectively.

We re-order rows and columns of the alignment matrix A accordingly. Define the following notations: $\mathcal{S}_1 = \{1, 2, \dots, n\}$ and $\mathcal{S}_2 = \{n + 1, n + 2, \dots, kn\}$. The alignment

matrix A for graphs G_1 and G_2 can be characterized using equation (5.3) as follows:

$$A(t_1, t_2) = \begin{cases} (\alpha + 1)G_1(i, j)G_2(i', j') - G_1(i, j) - G_2(i', j') + 1 + \epsilon, & \text{if } t_1 \sim (i, i'), t_2 \sim (j, j'), \\ & t_1 \text{ and } t_2 \in \mathcal{S}_1, t_1 \neq t_2. \\ (\alpha + 1)G_1(i, j)G_2(r', s') - G_1(i, j) - G_2(r', s') + 1 + \epsilon, & \text{if } t_1 \sim (i, r'), t_2 \sim (j, s'), \\ & t_1 \text{ or } t_2 \in \mathcal{S}_2, t_1 \neq t_2. \\ 1 + \epsilon, & \text{if } t_1 = t_2, \end{cases} \quad (5.49)$$

where notation $t_1 \sim (i, r')$ means that, row (and column) index t_1 of the alignment matrix A corresponds to the mapping (i, r') . Since G_1 and G_2 are isomorphic with permutation matrix $P = I$, we have $G_1(i, j) = G_2(i', j')$. Therefore, equation (5.49) can be written as,

$$A(t_1, t_2) = \begin{cases} (\alpha + 1)G_1(i, j)^2 - 2G_1(i, j) + 1 + \epsilon, & \text{if } t_1 \sim (i, i'), t_2 \sim (j, j'), \\ & t_1 \text{ and } t_2 \in \mathcal{S}_1, t_1 \neq t_2. \\ (\alpha + 1)G_1(i, j)G_1(r, s) - G_1(i, j) - G_1(r, s) + 1 + \epsilon, & \text{if } t_1 \sim (i, r'), t_2 \sim (j, s'), \\ & t_1 \text{ or } t_2 \in \mathcal{S}_2, t_1 \neq t_2. \\ 1 + \epsilon, & \text{if } t_1 = t_2. \end{cases} \quad (5.50)$$

Let \bar{A} be the expected alignment matrix, where $\bar{A}(t_1, t_2) = \mathbb{E}[A(t_1, t_2)]$, the expected value of $A(t_1, t_2)$.

Lemma 5.29. Let \mathbf{v} be the eigenvector of the expected alignment matrix \bar{A} corresponding to the largest eigenvalue. Suppose

$$\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix},$$

where \mathbf{v}_1 and \mathbf{v}_2 are vectors of length n and $(k-1)n$, respectively. Then,

$$\begin{aligned} v_{1,1} &= v_{1,2} = \dots = v_{1,n} = v_1^*, \\ v_{2,1} &= v_{2,2} = \dots = v_{2,(k-1)n} = v_2^*, \end{aligned}$$

Moreover, if $n \rightarrow \infty$, then,

$$\frac{v_1^*}{v_2^*} > \beta, \quad (5.51)$$

where $\beta = 1 + \Delta$, and $0 < \Delta k < \frac{(\alpha-1)p+1+\epsilon}{(\alpha+1)p^2-2p+1+\epsilon} - 1$.

Proof. Since $G_1(i, j)$ is a Bernoulli random variable which is one with probability p , equation (5.50) leads to:

$$\bar{A}(t_1, t_2) = \begin{cases} (\alpha-1)p+1+\epsilon, & \text{if } t_1 \text{ and } t_2 \in \mathcal{S}_1, t_1 \neq t_2, \\ (\alpha+1)p^2-2p+1+\epsilon, & \text{if } t_1 \text{ or } t_2 \in \mathcal{S}_2, t_1 \neq t_2, \\ 1+\epsilon, & \text{if } t_1 = t_2. \end{cases} \quad (5.52)$$

Define $a \triangleq (\alpha-1)p+1+\epsilon$ and $b \triangleq (\alpha+1)p^2-2p+1+\epsilon$.

Since $b\mathbf{v}$ is an eigenvector of \bar{A} , we have,

$$\bar{A}\mathbf{v} = \lambda\mathbf{v}, \quad (5.53)$$

where λ is the corresponding eigenvalue of $b\mathbf{v}$. Therefore,

$$\bar{A}\mathbf{v} = \begin{bmatrix} a \sum_i v_{1,i} + b \sum_j v_{2,j} + (1+\epsilon-a)v_{1,1} \\ \vdots \\ a \sum_i v_{1,i} + b \sum_j v_{2,j} + (1+\epsilon-a)v_{1,n} \\ b \sum_i v_{1,j} + b \sum_j v_{2,j} + (1+\epsilon-a)v_{2,1} \\ \vdots \\ b \sum_i v_{1,i} + b \sum_j v_{2,j} + (1+\epsilon-a)v_{2,(k-1)n} \end{bmatrix} = \lambda \begin{bmatrix} v_{1,1} \\ \vdots \\ v_{1,n} \\ v_{2,1} \\ \vdots \\ v_{2,(k-1)n} \end{bmatrix}. \quad (5.54)$$

Therefore,

$$\begin{aligned} a \sum_i v_{1,i} + b \sum_j v_{2,j} &= v_{1,r}(\lambda + a - 1 - \epsilon), & \forall 1 \leq r \leq n, \\ b \sum_i v_{1,i} + b \sum_j v_{2,j} &= v_{2,s}(\lambda + b - 1 - \epsilon), & \forall 1 \leq s \leq (k-1)n. \end{aligned} \quad (5.55)$$

We choose ϵ so that $\lambda + a - 1 - \epsilon \neq 0$ and $\lambda + b - 1 - \epsilon \neq 0$. We will show later in this section that any sufficiently small value of ϵ satisfies these inequalities. Therefore, equation (5.55) leads to,

$$\begin{aligned} v_{1,1} &= v_{1,2} = \dots = v_{1,n} = v_1^*, \\ v_{2,1} &= v_{2,2} = \dots = v_{2,(k-1)n} = v_2^*. \end{aligned} \quad (5.56)$$

Using equations (5.55) and (5.56), we have,

$$\begin{cases} anv_1^* + b(k-1)nv_2^* = v_1^*(\lambda + a - 1 - \epsilon) \\ bnv_1^* + b(k-1)nv_2^* = v_2^*(\lambda + b - 1 - \epsilon). \end{cases} \quad (5.57)$$

We choose ϵ so that $\lambda + b(1 - (k-1)n) - 1 - \epsilon \neq 0$. We will show later in this section that any sufficiently small value of ϵ satisfies this inequality. Further, according to PerronFrobenius Theorem 5.10, $v_{1,i} > 0$ and $v_{2,j} > 0$, for all i and j . Under these conditions, solving equation (5.57) leads to:

$$(\lambda - \lambda_a)(\lambda - \lambda_b) = b^2(k-1)n^2, \quad (5.58)$$

where,

$$\begin{cases} \lambda_a = (n-1)a + 1 + \epsilon, \\ \lambda_b = ((k-1)n-1)b + 1 + \epsilon. \end{cases} \quad (5.59)$$

Equation (5.58) has two solutions for λ . However, since λ is the largest eigenvalue of the expected alignment matrix \bar{A} , we choose the largest of the roots. Note that, since $b^2(k-1)n^2 > 0$, we have $\lambda > \max(\lambda_a, \lambda_b)$. This guarantees conditions that we

put on ϵ in the early steps of the proof.

By solving equations (5.58) and (5.59), we have,

$$\lambda = \frac{\lambda_a + \lambda_b + \sqrt{(\lambda_a - \lambda_b)^2 + 4(k-1)b^2n^2}}{2}.$$

First, we show $v_1^* > v_2^*$:

As $n \rightarrow \infty$, equation (5.57) implies,

$$\frac{v_1^*}{v_2^*} = \frac{\lambda}{bn} - k + 1, \quad (5.60)$$

where λ is the largest root of equation (5.58). For sufficiently large n ,

$$\frac{v_1^*}{v_2^*} = \frac{1}{2} \left[\left(\frac{a}{b} - k + 1 \right) + \sqrt{\left(\frac{a}{b} - k + 1 \right)^2 + 4k - 4} \right]. \quad (5.61)$$

If $p \neq 0, 1$, we always have $a > b$. Therefore, there exists $\Delta > 0$ such that $\frac{a}{b} > 1 + \Delta k$.

Thus, we have,

$$\frac{a}{b} > 1 + \Delta k > 1 + \Delta \left(1 + \frac{k-1}{1+\Delta} \right) = 1 + \Delta + \frac{\Delta}{\Delta+1} (k-1). \quad (5.62)$$

Using inequality (5.62) in (5.61), we have,

$$\begin{aligned} \frac{v_1^*}{v_2^*} &> \frac{1}{2} \left[\frac{(1+\Delta)^2 - k + 1}{1+\Delta} + \frac{\sqrt{((1+\Delta)^2 - k + 1)^2 + 4(k-1)(1+\Delta)^2}}{1+\Delta} \right] \\ &= 1 + \Delta. \end{aligned} \quad (5.63)$$

This completes the proof of Lemma 5.29. □

Remark 5.30. In Lemma 5.29, if $kp = c_1$, where $c_1 \lambda k$ and $k\Gamma 1$, choosing $\alpha = c_2 k^2$ results in $\Delta \approx ck$ where $c = \frac{c_1 c_2}{c_1^2 c_2 + 1} > 1$.

If \mathbf{v} is an eigenvector with unit norm, we have:

$$\|\mathbf{v}\| = 1 \Rightarrow nv_1^* + (k-1)nv_2^* = 1. \quad (5.64)$$

By using this equation and Lemma 5.29, for sufficiently large n , we have:

$$\begin{aligned} v_1^* &= \frac{\beta}{\sqrt{\beta^2 + k}} \frac{1}{\sqrt{n}}, \\ v_1^* &= \frac{1}{\sqrt{\beta^2 + k}} \frac{1}{\sqrt{n}}. \end{aligned} \tag{5.65}$$

Since edges of network G_1 are random, the alignment network is a also random matrix. Let Ω be a set of outcomes of edge variables of network G_1 , and let $A(\Omega)$ be the corresponding alignment network. Let $\mathbf{v}(\Omega)$ represent an eigenvector of $A(\Omega)$ with the largest eigenvalue.

Lemma 5.31. Let $\mathbf{v}(\Omega)$ be a unite norm eigenvector with the largest eigenvalue of the alignment matrix $A(\Omega)$. Suppose

$$\mathbf{v}(\Omega) = \begin{bmatrix} \mathbf{v}_1(\Omega) \\ \mathbf{v}_2(\Omega) \end{bmatrix},$$

where $\mathbf{v}_1(\Omega)$ and $\mathbf{v}_2(\Omega)$ are vectors of length n and kn , respectively. Let $\bar{\mathbf{v}} = \mathbb{E}[\mathbf{v}(\Omega)]$.

Then,

- (i) For all $1 \leq i \leq n$, $\mathbb{E}[v_{1,i}(\Omega)] = \bar{\mathbf{v}}_1$, and $\sigma_{v_{1,i}(\Omega)}^2 = \sigma_1^2$.
- (ii) For all $1 \leq j \leq (k-1)n$, $\mathbb{E}[v_{2,j}(\Omega)] = \bar{\mathbf{v}}_2$, and $\sigma_{v_{2,j}(\Omega)}^2 = \sigma_2^2$.
- (iii) Let \mathbf{v} be the eigenvector with the largest eigenvalue of the expected alignment matrix \bar{A} . Then, $\mathbf{v}(\Omega) \cdot \mathbf{v} \geq 1 - \gamma$, where γ is a small positive number, w.h.p.
- (iv) $1 - \gamma \leq \|\bar{\mathbf{v}}\| \leq 1$.
- (v) $\sigma_1^2 \leq \frac{2\gamma}{n}$ and, $\sigma_2^2 \leq \frac{2\gamma}{(k-1)n}$.

- (vi) If $v_1^*/v_2^* = \beta$, for sufficiently large n and w.h.p., $\bar{v}_1 = \frac{a_1}{\sqrt{n}}$ and $\bar{v}_2 = \frac{a_2}{\sqrt{n}}$, where $a_1 > a_2$ and,

$$\begin{aligned} a_1 &\approx \frac{\beta(1-\gamma)}{\sqrt{\beta^2 + k - 1}}, \\ a_2 &\approx \frac{(1-\gamma)}{\sqrt{\beta^2 + k - 1}}. \end{aligned} \quad (5.66)$$

Proof. Owing to symmetry of the problem, for all $1 \leq i \leq n$, random variables $v_{1,i}(\Omega)$ have the same probability distribution. Therefore, they have the same mean and variance. The same argument holds for random variables $v_{2,j}(\Omega)$ for all $1 \leq j \leq (k-1)n$. This proves parts (i) and (ii).

To prove part (iii), first we show that, with high probability and for sufficiently large n , $\|E\| \triangleq \|A(\Omega) - \bar{A}\| < \delta n$, where δ is a small positive number, and $\|\cdot\|$ is the spectral norm operator. Since \bar{A} and $A(\Omega)$ have the same diagonal elements, all diagonal elements of E are zero.

Theorem 5.32 (Gershgorin Circle Theorem). Let E be a complex matrix with entries $e_{i,j}$. Let $R_i = \sum_{j \neq i} |e_{i,j}|$ be the sum of absolute values of the off-diagonal entries in the row i . Let $D(e_{i,i}, R_i)$ be the closed disc centered at $e_{i,i}$ with radius R_i . Every eigenvalue of E lies within at least one of the Gershgorin discs $D(e_{ii}, R_i)$.

Proof. See reference [211]. □

We use Gershgorin Circle Theorem to show that, $\|E\| < \delta n$ for sufficiently large n , w.h.p. First, we show that, for sufficiently large n and w.h.p., $R_i < \delta n$, where δ is a small positive number. To simplify notations, we write R_t for different t as follows:

If $i \leq n$, then,

$$\begin{aligned} R_t &= \sum_{i=2}^n (\alpha + 1)G_i^2 - 2G_i + 1 + \sum_{(i,j) \in \mathcal{B}} (\alpha + 1)G_i G_j - G_i - G_j + 1 \\ &\quad - (n-1)((\alpha + 1)p - 2p + 1) - (k-1)n((\alpha + 1)p^2 - 2p + 1), \end{aligned} \quad (5.67)$$

where G_i is an iid Bernoulli variable with $Pr[G_i = 1] = p$, and $\mathcal{B} = \mathcal{R} - \{(i, i') :$

$1 \leq i \leq n$. Similarly, if $n < t \leq kn$, we can write,

$$R_t = \sum_{(i,j) \in \mathcal{B}} ((\alpha + 1)G_i G_j - G_i - G_j + 1) - (k - 1)n((\alpha + 1)p^2 - 2p + 1). \quad (5.68)$$

Using Chernoff bound, for sufficiently large n , for a given $\delta_1 > 0$, there exists $\epsilon_1 > 0$ so that,

$$\begin{aligned} Pr \left[\left| \frac{1}{n} \sum_{i=1}^n G_i - p \right| > \delta_1 \right] &\leq e^{-n\epsilon_1}, \\ Pr \left[\left| \frac{1}{n} \sum_{i=1}^n G_i^2 - p \right| > \delta_1 \right] &\leq e^{-n\epsilon_1}, \\ Pr \left[\left| \frac{1}{(k-1)n} \sum_{(i,j) \in \mathcal{B}} G_i G_j - p^2 \right| > \delta_1 \right] &\leq e^{-n\epsilon_1}. \end{aligned} \quad (5.69)$$

Proposition 5.33. Let U_1 and U_2 be two random variables such that,

$$\begin{aligned} Pr[U_1 \in (-\delta_1, \delta_1)] &> 1 - e^{-n\epsilon_1}, \\ Pr[U_2 \in (-\delta_2, \delta_2)] &> 1 - e^{-n\epsilon_2}. \end{aligned}$$

Then, w.h.p.,

$$\begin{aligned} Pr[U_1 + U_2 \geq \delta_1 + \delta_2] &\leq e^{-n \min(\epsilon_1, \epsilon_2)}, \\ Pr[U_1 U_2 \geq \delta_1 \delta_2] &\leq e^{-n \min(\epsilon_1, \epsilon_2)}. \end{aligned}$$

Proof. Let T be a random variable representing the event $U_1 \in (-\delta_1, \delta_1)$, and T^c be its complement. Then,

$$\begin{aligned} Pr[U_1 + U_2 \geq \delta_1 + \delta_2] &= Pr[U_1 + U_2 \geq \delta_1 + \delta_2 | T] Pr[T] \\ &\quad + Pr[U_1 + U_2 \geq \delta_1 + \delta_2 | T^c] Pr[T^c] \\ &\leq Pr[U_2 \geq \delta_2] + Pr[T^c] \\ &\leq e^{-n \min(\epsilon_1, \epsilon_2)}. \end{aligned}$$

A similar argument can be made for the case of U_1U_2 . This completes the proof of Proposition 5.33. \square

According to equations (5.67) and (5.68), for all $1 \leq t \leq kn$, $E[R_t] = 0$. Using Proposition 5.33 and equation (5.69), there exists $\delta > 0$ such that, $|R_t| \leq \delta n$ for sufficiently large n , w.h.p. Thus, using Gershgorin circle Theorem 5.32, the maximum eigenvalue of matrix E is smaller than δn , for sufficiently large n , w.h.p., which indicates that $\|E\| \leq \delta n$, for sufficiently large n , w.h.p.

Theorem 5.34 (Wedin Sin Theorem). Let \mathbf{v} and $\mathbf{v}(\Omega)$ be eigenvectors of matrices \bar{A} and $A(\Omega)$ corresponding to their largest eigenvalues, respectively. Let λ_1 and λ_2 be the largest and second largest eigenvalues of matrix \bar{A} . Then, there exists a positive constant μ such that,

$$|\sin \angle(\mathbf{v}, \mathbf{v}(\Omega))| \leq \mu \frac{\|A - A(\Omega)\|}{\lambda_1 - \lambda_2}. \quad (5.70)$$

Proof. See reference [210]. \square

According to equation (5.55), as $n \rightarrow \infty$, the two largest eigenvalues of the matrix \bar{A} are the roots of equation (5.58) because other eigenvalues are equal to $-a + 1 + \epsilon$ or $-b + 1 + \epsilon$. Solving equation (5.58) for sufficiently large n , we have,

$$\begin{aligned} |\lambda_1 - \lambda_2| &= \sqrt{(\lambda_a - \lambda_b)^2 + 4(k-1)b^2n^2} \\ &= nb \sqrt{\left(\frac{a}{b} - k + 1\right)^2 + 4(k-1)} \\ &> nb \sqrt{\left(1 + \Delta \frac{k-1}{1+\Delta}\right) + 4(k-1)} \\ &= \frac{b(\beta^2 + k - 1)}{\beta} n. \end{aligned} \quad (5.71)$$

Therefore, using equation (5.71) in Wedin sin Theorem 5.34, for any small positive δ , we have,

$$|\sin \angle(\mathbf{v}, \mathbf{v}(\Omega))| \leq \frac{\mu\beta}{b(\beta^2 + k - 1)} \delta. \quad (5.72)$$

This completes the proof of part (iii).

To prove part (iv), first we use Jensen's inequality. Since norm is a convex function, we have,

$$\|\bar{\mathbf{v}}\| = \|\mathbb{E}[\mathbf{v}(\Omega)]\| \leq \mathbb{E}[\|\mathbf{v}(\Omega)\|] = 1. \quad (5.73)$$

From part (iii), we have,

$$\mathbf{v}(\Omega) \cdot \mathbf{v} \geq 1 - \gamma \Rightarrow \bar{\mathbf{v}} \cdot \mathbf{v} \geq 1 - \gamma. \quad (5.74)$$

Then, using Cauchy Schwarz inequality, we have,

$$\|\bar{\mathbf{v}}\| = \|\bar{\mathbf{v}}\| \|\mathbf{v}\| \geq \bar{\mathbf{v}} \cdot \mathbf{v} \geq 1 - \gamma. \quad (5.75)$$

This completes the proof of part (iv).

To prove part (v), we can write,

$$\begin{aligned} \mathbb{E}[\|\mathbf{v}(\Omega) - \bar{\mathbf{v}}\|^2] &= \mathbb{E}[1 + \|\bar{\mathbf{v}}\|^2 - 2\mathbf{v}(\Omega) \cdot \bar{\mathbf{v}}] = 1 - \|\bar{\mathbf{v}}\|^2 \\ &\leq 1 - (1 - \gamma)^2 = \gamma(2 - \gamma) \leq 2\gamma. \end{aligned} \quad (5.76)$$

On the other hand,

$$\mathbb{E}[\|\mathbf{v}(\Omega) - \bar{\mathbf{v}}\|^2] = n\sigma_1^2 + (k - 1)n\sigma_2^2. \quad (5.77)$$

Therefore,

$$\begin{aligned} n\sigma_1^2 \leq 2\gamma &\Rightarrow \sigma_1^2 \leq \frac{2\gamma}{n}, \\ (k - 1)n\sigma_2^2 \leq 2\gamma &\Rightarrow \sigma_2^2 \leq \frac{2\gamma}{(k - 1)n}. \end{aligned} \quad (5.78)$$

This completes the proof of part (v).

To prove part (vi), we assume that in the worst case, vectors \mathbf{v} and $\bar{\mathbf{v}}$ has inner

product value of $1 - \gamma$. Therefore, using part (iii) and (iv), we have,

$$\begin{cases} nv_1^* \bar{v}_1 + (k-1)nv_2^* \bar{v}_2 = 1 - \gamma, \\ n\bar{v}_1^2 + (k-1)n\bar{v}_2^2 = c, \end{cases} \quad (5.79)$$

where $1 - \gamma \leq c = \|\bar{v}\| \leq 1$. Solving equation (5.79) using equation (5.65) for sufficiently large n , we have,

$$\begin{aligned} a_1 &= \frac{\beta(1-\gamma) + \sqrt{(k-1)(c - (\gamma-1)^2)}}{\sqrt{\beta^2 + k-1}}, \\ a_2 &= \frac{(1-\gamma)\sqrt{k-1} - \beta\sqrt{c - (\gamma-1)^2}}{\sqrt{(k-1)(\beta^2 + k-1)}}. \end{aligned} \quad (5.80)$$

Choosing δ sufficiently small and n sufficiently large, equation (5.80) leads to equation (5.79). This completes the proof of part (vi). \square

Now, we have all technology to prove Theorem 5.17. Recall notations $\mathcal{S}_1 = \{1, 2, \dots, n\}$ and $\mathcal{S}_2 = \{n+1, n+2, \dots, kn\}$. Let $\mathcal{S}_1(m) \subseteq \mathcal{S}_1$ and $\mathcal{S}_2(m) \subseteq \mathcal{S}_2$ where $|\mathcal{S}_1(m)| = |\mathcal{S}_2(m)| = m$. Define following variables:

$$\begin{aligned} U_1 &\triangleq \frac{1}{n} \sum_{i \in \mathcal{S}_1} v_i(\Omega), \\ U_2 &\triangleq \frac{1}{n} \left[\sum_{i \in \mathcal{S}_1(n-m)} v_i(\Omega) + \sum_{j \in \mathcal{S}_2(m)} v_j(\Omega) \right]. \end{aligned} \quad (5.81)$$

Let $\rho = m/n \neq 0$. According to Lemma 5.31, $\mathbb{E}[U_1] = \bar{v}_1$ and $\mathbb{E}[U_2] = (1-\rho)\bar{v}_1 + \rho\bar{v}_2$. Moreover, $\sigma_{U_1}^2 \leq \sigma_1^2$ and $\sigma_{U_2}^2 \leq (1-\rho)\sigma_1^2 + \rho\sigma_2^2 < \sigma_1^2$. Define,

$$d \triangleq \frac{|\mathbb{E}[U_1] - \mathbb{E}[U_2]|}{2} = \frac{c_1(1-\gamma)}{\sqrt{n}}, \quad (5.82)$$

where $c_1 = \frac{\rho\Delta}{2\sqrt{(1+\Delta)^2 + k-1}}$.

Using Chebyshev's inequality, we have,

$$\begin{aligned} \Pr[U_1 \leq \mathbb{E}[U_1] - d] &\leq \frac{2}{c_1^2} \frac{\gamma}{(1-\gamma)^2} \leq c_2 \gamma = \epsilon_1, \\ \Pr[U_2 \geq \mathbb{E}[U_2] + d] &\leq \frac{2}{c_1^2} \frac{\gamma}{(1-\gamma)^2} \leq c_2 \gamma = \epsilon_1. \end{aligned} \quad (5.83)$$

Therefore,

$$\Pr[U_1 < U_2] < \epsilon_1, \quad (5.84)$$

where ϵ_1 can be arbitrarily small for sufficiently large n . This completes the proof of Theorem 5.17.

5.8.3 Proof of Theorem 5.15

Without loss of generality and similarly to the proof of Theorem 5.17, let $P = I$. Let A be the alignment network of graphs G_1 and G_2 defined according to equation (5.3). Similarly to the proof of Theorem 5.17, re-order row (and column) indices of matrix A so that the first n indices correspond to the true mappings $\{(i, i') : i \in \mathcal{V}_1, i' \in \mathcal{V}_2\}$. Define the expected alignment network \bar{A} as $\bar{A}(t_1, t_2) = \mathbb{E}[A(t_1, t_2)]$, where t_1 and t_2 are two possible mappings across networks. Recall notations $\mathcal{S}_1 = \{1, 2, \dots, n\}$ and $\mathcal{S}_2 = \{n+1, n+2, \dots, kn\}$.

First, we consider the noise model I (5.14):

Define,

$$\begin{aligned} a' &\triangleq p(1-p_e)(\alpha + \epsilon) + (1-p)(1-p_e)(1 + \epsilon) + (pp_e + (1-p)p_e)\epsilon \\ b' &\triangleq (p^2(1-p_e) + pp_e(1-p))(\alpha + \epsilon) \\ &\quad + ((1-p)^2(1-p_e) + pp_e(1-p))(1 + \epsilon) \\ &\quad + (2p(1-p)(1-p_e) + 2p^2p_e)\epsilon. \end{aligned} \quad (5.85)$$

Since $G_1(i, j)$ and $Q(i, j)$ are Bernoulli random variables with parameters p and p_e ,

respectively, the expected alignment network can be simplified as follows:

$$\bar{A}(t_1, t_2) = \begin{cases} a', & \text{if } t_1 \text{ and } t_2 \in \mathcal{S}_1, t_1 \neq t_2, \\ b', & \text{if } t_1 \text{ or } t_2 \in \mathcal{S}_2, t_1 \neq t_2, \\ 1 + \epsilon, & \text{if } t_1 = t_2. \end{cases} \quad (5.86)$$

We have,

$$a' - b' = (\alpha + 1)(2p_e - 1)p(p - 1) + p_e(1 - 2p)\epsilon. \quad (5.87)$$

Thus, if $p \neq 0, 1$ and $p_e < 1/2$, for small enough ϵ , $a' > b' > 0$. Therefore, there exists a positive Δ such that $\frac{a'}{b'} = 1 + \Delta$. The rest of the proof is similar to the one of Theorem 5.17.

The proof for the noise model II of (5.15) is similar. To simplify notation and illustrate the main idea, here we assume ϵ is sufficiently small with negligible effects.

Define,

$$\begin{aligned} a'' &\triangleq p(1 - p_e)(\alpha) + (1 - p)(1 - p_{e_2}) = 1 - p(1 + \alpha(p_e - 1) + p_e) \\ b'' &\triangleq p^2(1 - p_e)\alpha + (1 - p)^2(1 - p_{e_2}) + 2p(1 - p)p_{e_2}(1 + \alpha) \\ &= 1 - p(2 + p_e) + p^2(1 + \alpha + 2p_e). \end{aligned} \quad (5.88)$$

The expected alignment network in this case is:

$$\bar{A}(t_1, t_2) = \begin{cases} a'', & \text{if } t_1 \text{ and } t_2 \in \mathcal{S}_1, t_1 \neq t_2, \\ b'', & \text{if } t_1 \text{ or } t_2 \in \mathcal{S}_2, t_1 \neq t_2, \\ 1 + \epsilon, & \text{if } t_1 = t_2. \end{cases} \quad (5.89)$$

Moreover, we have,

$$a'' - b'' = p((1 - p - p_e)(1 + \alpha) + p_e(1 - 2p)). \quad (5.90)$$

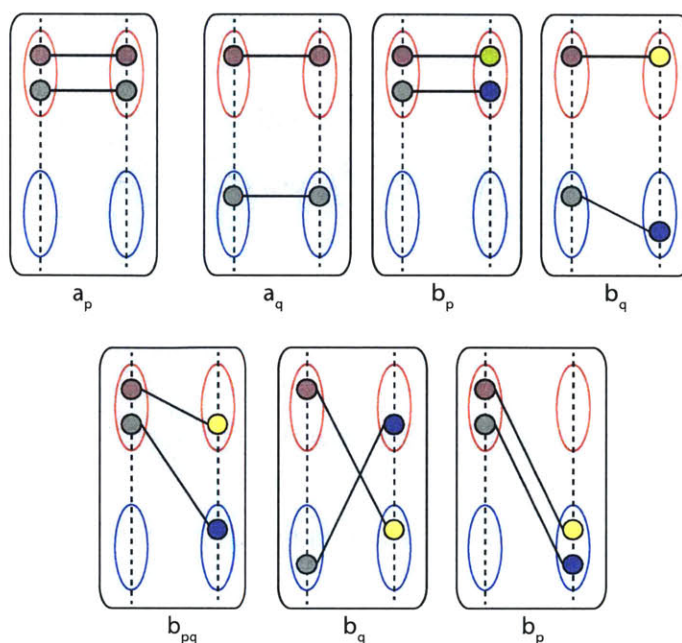


Figure 5-22: Expected alignment scores of different mapping combinations. Ellipses with the same color represent corresponding modules across two networks. Nodes with the same color represent true mappings across two networks.

If $p < 1/2$ and $p_e < 1/2$, then $a'' - b'' > 0$. The rest of the proof is similar to the previous case.

5.8.4 Proof Of Theorem 5.18

By writing Taylor's expansion of $Tr(G_1 X G_2 X^T)$ around the point X_0 , we have,

$$Tr(G_1 X G_2 X^T) = Tr(G_1 X_0 G_2 X_0^T) + 2Tr(G_1 X_0 G_2 (X - X_0)^T) + Tr(G_1 (X - X_0) G_2 (X - X_0)^T). \quad (5.91)$$

Let $\Delta \triangleq X - X_0$. Thus, we have,

$$f(X) = \tilde{f}(X) + Tr(G_1 \Delta G_2 \Delta^T). \quad (5.92)$$

Let $\sigma_i(G)$ be the i -th largest singular value of matrix G .

Theorem 5.35 (Von Neumann's trace inequality). Suppose A and B are two $n \times n$

complex matrices. We have,

$$|Tr(AB)| \leq \sum_{i=1}^n \sigma_i(A)\sigma_i(B). \quad (5.93)$$

Using Theorem 5.35, we have,

$$|Tr(G_1\Delta G_2\Delta^T)| \leq \sum_{i=1}^n \sigma_i(G_1\Delta)\sigma_i(G_2\Delta^T). \quad (5.94)$$

Moreover, using Allahverdiev's characterization of singular values, we have,

$$\begin{aligned} \sigma_i(G_1\Delta) &\leq \min\{\sigma_k(G_1)\sigma_{i+1-k}(\Delta) : 1 \leq k \leq i\}, \\ &\leq \sigma_i(G_1)\sigma_1(\Delta). \end{aligned} \quad (5.95)$$

Moreover,

$$\begin{aligned} \sigma_i(G_1\Delta) &\geq \max\{\sigma_k(G_1)\sigma_{n+i-k}(\Delta) : i \leq k \leq n\}, \\ &\geq \sigma_i(G_1)\sigma_n(\Delta). \end{aligned} \quad (5.96)$$

Moreover, since $\|\Delta\|_2 \leq \epsilon$, we have $|\sigma_i(\Delta)| \leq \epsilon$, for $1 \leq i \leq n$. Using (5.95) and (5.96), we have,

$$|\sigma_i(G_1\Delta)| \leq \epsilon\sigma_i(G_1). \quad (5.97)$$

Similarly, we have,

$$|\sigma_i(G_2\Delta^T)| \leq \epsilon\sigma_i(G_2). \quad (5.98)$$

Thus, using (5.94) and (5.98), we have,

$$|Tr(G_1\Delta G_2\Delta^T)| \leq \epsilon^2 \sum_{i=1}^n \sigma_i(G_1)\sigma_i(G_2). \quad (5.99)$$

Using (5.92) and (5.99), we have,

$$|f(X) - \tilde{f}(X)| \leq \epsilon^2 \sum_{i=1}^n \sigma_i(G_1) \sigma_i(G_2). \quad (5.100)$$

Moreover, since X^* and X_{lin}^* are optimal solutions of optimizations (5.19) and (5.23), respectively, we have,

$$\begin{aligned} f(X^*) &\leq \tilde{f}(X^*) + \epsilon^2 \sum_{i=1}^n \sigma_i(G_1) \sigma_i(G_2) \\ &\leq \tilde{f}(X_{lin}^*) + \epsilon^2 \sum_{i=1}^n \sigma_i(G_1) \sigma_i(G_2) \\ f(X^*) &\geq f(X_{lin}^*) \geq \tilde{f}(X_{lin}^*) - \epsilon^2 \sum_{i=1}^n \sigma_i(G_1) \sigma_i(G_2). \end{aligned} \quad (5.101)$$

This completes the proof.

5.8.5 Proof of Theorem 5.25

Without loss of generality, let $P = I$ in (5.17). Let A be the alignment network of graphs G_1 and G_2 defined according to (5.3). G_1 and G_2 are stochastic block networks with m modules, each with n nodes. Modules of G_1 are represented by $\{V_1^1, V_1^2, \dots, V_1^m\}$ where $V_1^a = \{i_a^1, i_a^2, \dots, i_a^n\}$. Similarly, modules of G_2 are represented by $\{V_2^1, V_2^2, \dots, V_2^m\}$ where $V_2^a = \{j_a^1, j_a^2, \dots, j_a^n\}$. The alignment network A has $n^2 m^2$ nodes. We re-order row (and column) indices of matrix A so that the first mn^2 indices are within module mappings:

$$(i_1^1, j_1^1), (i_1^2, j_1^2), \dots, (i_1^n, j_1^n), (i_1^1, j_1^2), (i_1^1, j_1^3), \dots, (i_1^n, j_1^{n-1}).$$

The remaining $(m^2 - m)n^2$ indices correspond to across module mappings:

$$(i_1^1, j_2^1), (i_1^1, j_2^2), \dots, (i_m^n, j_{m-1}^{n-1}), (i_m^n, j_{m-1}^n).$$

Define,

$$\begin{aligned}
 a_p &\triangleq (\alpha - 1)p + 1, & (5.102) \\
 b_p &\triangleq (\alpha + 1)p^2 - 2p + 1, \\
 a_q &\triangleq (\alpha - 1)q + 1, \\
 b_q &\triangleq (\alpha + 1)q^2 - 2q + 1, \\
 b_{pq} &\triangleq \alpha pq + (1 - p)(1 - q).
 \end{aligned}$$

To form the expected alignment network \bar{A} , we consider seven types of mapping-pairs illustrated in Figure 5-22. Thus, the expected network alignment matrix has the following structure:

$$\bar{A} = \left[\begin{array}{cccc|cccc}
 A_1 & A_2 & A_2 & \cdots & A_2 & A_3 & \cdots & A_3 \\
 A_2 & A_1 & A_2 & \cdots & A_2 & & & \\
 \vdots & & \ddots & & \vdots & \vdots & \ddots & \vdots \\
 A_2 & \cdots & & & A_1 & A_3 & \cdots & A_3 \\
 \hline
 A_3 & \cdots & & & A_3 & A_4 & A_5 & A_5 & \cdots & A_5 \\
 & & & & & A_5 & A_4 & A_5 & \cdots & A_5 \\
 \vdots & & \ddots & & \vdots & \vdots & & \ddots & & \vdots \\
 A_3 & \cdots & & & A_3 & A_5 & \cdots & & & A_4
 \end{array} \right], \quad (5.103)$$

where the upper and lower blocks have mn^2 and $(m^2 - m)n^2$ rows, respectively. A_i

is a matrix of size $n^2 \times n^2$, defined as follows:

$$\begin{aligned}
 A_1 &= \left[\begin{array}{ccc|ccc} a_p & \cdots & a_p & b_p & \cdots & b_p \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ a_p & \cdots & a_p & b_p & \cdots & b_p \\ \hline b_p & \cdots & b_p & b_p & \cdots & b_p \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ b_p & \cdots & b_p & b_p & \cdots & b_p \end{array} \right], & (5.104) \\
 A_2 &= \left[\begin{array}{ccc|ccc} a_q & \cdots & a_q & b_q & \cdots & b_q \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ a_q & \cdots & a_q & b_q & \cdots & b_q \\ \hline b_q & \cdots & b_q & b_q & \cdots & b_q \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ b_q & \cdots & b_q & b_q & \cdots & b_q \end{array} \right],
 \end{aligned}$$

$$A_3 = b_{pq} \mathbf{1}_{n^2},$$

$$A_4 = b_p \mathbf{1}_{n^2},$$

$$A_5 = b_q \mathbf{1}_{n^2},$$

where $\mathbf{1}_{n^2}$ is a $n^2 \times n^2$ matrix whose elements are ones.

Suppose v is the leading eigenvector of the expected alignment matrix \bar{A} . We wish to show that,

$$v(i) > v(j), \quad \forall 1 \leq i \leq mn^2, \forall mn^2 < j \leq m^2n^2. \quad (5.105)$$

Lemma 5.36. Let A be a positive matrix. If one or any number of entries of row i are increased and all the other rows remain fixed, and if the i -th entry of the Perron vector is held a fixed constant equal to 1, then the remaining entries of the Perron vector strictly decrease.

Proof. See a proof in [230]. □

Define matrix B as follows:

$$B = \left[\begin{array}{ccccc|ccc} B_1 & B_2 & B_2 & \cdots & B_2 & B_3 & \cdots & B_3 \\ B_2 & B_1 & B_2 & \cdots & B_2 & & & \\ \vdots & & \ddots & & \vdots & \vdots & \ddots & \vdots \\ B_2 & & \cdots & & B_1 & B_3 & \cdots & B_3 \\ \hline B_3 & & \cdots & & B_3 & B_1 & B_2 & B_2 & \cdots & B_2 \\ & & & & & B_2 & B_1 & B_2 & \cdots & B_2 \\ \vdots & & \ddots & & \vdots & \vdots & & \ddots & & \vdots \\ B_3 & & \cdots & & B_3 & B_2 & \cdots & & & B_1 \end{array} \right], \quad (5.106)$$

where,

$$B_1 = b_p \mathbf{1}_{n^2}, \quad (5.107)$$

$$B_2 = b_q \mathbf{1}_{n^2},$$

$$B_3 = b_{pq} \mathbf{1}_{n^2}.$$

Because $a_p > b_p$ and $a_q > b_q$, using Lemma 5.36, it is sufficient to show (5.105) for the matrix B . Suppose v_B is the leading eigenvector of the matrix B . We wish to show,

$$v_B(i) > v_B(j), \quad \forall 1 \leq i \leq mn^2, \forall mn^2 < j \leq m^2n^2. \quad (5.108)$$

Note that,

$$B = \mathbf{1}_{n^2} \otimes C, \quad (5.109)$$

where C is a $m^2 \times m^2$ matrix with a structure illustrated as follows:

$$C = \left[\begin{array}{ccccc|ccccc} b_p & b_q & b_q & \cdots & b_q & b_{pq} & & \cdots & & b_{pq} \\ b_q & b_p & b_q & \cdots & b_q & & & & & \\ \vdots & & \ddots & & \vdots & \vdots & & \ddots & & \vdots \\ b_q & & \cdots & & b_p & b_{pq} & & \cdots & & b_{pq} \\ \hline b_{pq} & & \cdots & & b_{pq} & b_p & b_q & b_q & \cdots & b_q \\ & & & & & b_q & b_p & b_q & \cdots & b_q \\ \vdots & & \ddots & & \vdots & \vdots & & \ddots & & \vdots \\ b_{pq} & & \cdots & & b_{pq} & b_q & & \cdots & & b_p \end{array} \right]. \quad (5.110)$$

Let v_C be the leading eigenvector of the matrix C . Using the property of Kronecker products, we have,

$$v_B = \mathbf{1}_{n^2} \otimes v_C. \quad (5.111)$$

Therefore, to show (5.105) and (5.108), it is sufficient to show that,

$$v_C(i) > v_C(j), \quad \forall 1 \leq i \leq m, \forall m < j \leq m^2. \quad (5.112)$$

Lemma 5.37. Consider the matrix

$$X = \left[\begin{array}{ccccc|ccccc} c & b & b & \cdots & b & a & & \cdots & & a \\ b & c & b & \cdots & b & & & & & \\ \vdots & & \ddots & & \vdots & \vdots & & \ddots & & \vdots \\ b & & \cdots & & c & a & & \cdots & & a \\ \hline a & & \cdots & & a & c & b & b & \cdots & b \\ & & & & & b & c & b & \cdots & b \\ \vdots & & \ddots & & \vdots & \vdots & & \ddots & & \vdots \\ a & & \cdots & & a & b & & \cdots & & c \end{array} \right], \quad (5.113)$$

where the size of the top and bottom blocks are n_1 and n_2 , respectively. Let v^* be

the leading eigenvector of X . Then, if $a > b > 0$, $c > 0$, and $n_2 > n_1$, we have,

$$v^*(i) > v^*(j), \quad \forall 1 \leq i \leq n_1, \forall n_1 < j \leq n_1 + n_2. \quad (5.114)$$

Proof. Suppose v is an eigenvector of the matrix X with the corresponding eigenvalue λ . Owing to the symmetric structure of the matrix X , we have,

$$\begin{cases} v(1) = \dots = v(n_1) \triangleq v_1, \\ v(n_1 + 1) = \dots = v(n_1 + n_2) \triangleq v_2. \end{cases} \quad (5.115)$$

Using (5.115) in the eigen decomposition equality, we have,

$$\begin{cases} cv_1 + (n_1 - 1)bv_1 + n_2av_2 = \lambda v_1 \\ n_1av_1 + cv_2 + (n_2 - 1)bv_2 = \lambda v_2. \end{cases} \quad (5.116)$$

Thus, we have,

$$(\lambda' - \lambda_1)(\lambda' - \lambda_2) = n_1n_2a^2, \quad (5.117)$$

where,

$$\begin{cases} \lambda' = \lambda - c, \\ \lambda_1 = (n_1 - 1)b, \\ \lambda_2 = (n_2 - 1)b. \end{cases} \quad (5.118)$$

Let λ^* be the largest root of (5.117), which corresponds to the leading eigenvector of the matrix v^* . To prove the lemma, it is sufficient to show that,

$$\lambda^* > n_1a + (n_2 - 1)b. \quad (5.119)$$

This is true if $a > b$ and $n_2 > n_1$. To show this, we need to show,

$$\lambda^* = \frac{\lambda_1 + \lambda_2 + \sqrt{(\lambda_1 - \lambda_2)^2 + 4n_1n_2a^2}}{2} > n_1a + (n_2 - 1)b, \quad (5.120)$$

which is true under the conditions of the lemma. This completes the proof. \square

If $p > q$ and $\alpha > 1/q - 1$, we have $b_{pq} > b_q$. Thus, Lemma 5.37 leads to (5.112). This completes the proof.

5.8.6 Proof of Theorem 5.26

We use the same setup considered in the proof of Theorem 5.25 to form the expected alignment network. Considering $0 < p_e^2 \lambda_1$ and $0 < p_{e_2}^2 \lambda_1$, expected scores of different mapping-pairs illustrated in Figure 5-22 can be approximated as follows:

$$\begin{aligned}
 a_p &\simeq p(1 - p_e)\alpha, & (5.121) \\
 b_p &\simeq p^2(1 - p_e)^2\alpha + 2p(1 - p)p_{e_2}\alpha, \\
 a_q &\simeq q(1 - p_e)\alpha, \\
 b_q &\simeq q^2(1 - p_e)^2\alpha + 2q(1 - q)p_{e_2}\alpha, \\
 b_{pq} &\simeq pq(1 - p_e)^2\alpha + p(1 - q)p_{e_2}\alpha + (1 - p)qp_e\alpha.
 \end{aligned}$$

The proof is similar to the one of Theorem 5.25. To use Lemma 5.36, we need to have,

$$\begin{cases} a_p > b_p, \\ a_q > b_q, \end{cases} \quad (5.122)$$

which results in following conditions:

$$\begin{aligned}
 p &< \frac{1 - pp_e}{1 + p_e^2}, & (5.123) \\
 q &< \frac{1 - pp_e}{1 + p_e^2}.
 \end{aligned}$$

Because $p_e^2 \ll 1$, we then have,

$$\begin{aligned}
 p &< \frac{1}{1 + p_e}, & (5.124) \\
 q &< \frac{1}{1 + p_e}.
 \end{aligned}$$

To use Lemma 5.37, we need to have $b_{pq} > b_q$. Using (5.121), we have:

$$b_{pq} - b_q = (1 - p)q(p - q) + p_e \underbrace{(p^2(1 + 2q) - 6pq + q(1 + 2q))}_{\text{polynomial I}}. \quad (5.125)$$

To show the non-negativity of the right-hand side of (5.125), it is sufficient to show the non-negativity of polynomial I. This polynomial has two roots at

$$p_{\text{roots}} = \frac{6q \pm \sqrt{-4q(4q - 1)(q - 1)}}{2(1 + 2q)}. \quad (5.126)$$

If $0 < q \leq 1/4$,

$$-4q(4q - 1)(q - 1) < 0. \quad (5.127)$$

Because the value of the polynomial I at $p = 0$ is positive, if $0 < q \leq 1/4$, the polynomial is always non-negative. If $q > 1/4$, we need to have,

$$p \leq \frac{6q - \sqrt{4q(4q - 1)(1 - q)}}{2(1 + 2q)}, \quad (5.128)$$

which guarantees the non-negativity of polynomial I. The rest of the proof is similar to the one of Theorem 5.25.

5.8.7 Proof of Theorem 5.27

We use the same setup considered in the proof of Theorem 5.25 to form the expected alignment network. Suppose $S(P)$ and $S(\tilde{P})$ correspond to the expected objective function of the network alignment optimization 5.5 using permutation matrices P and \tilde{P} , respectively, where,

$$\frac{1}{nm} \|P - \tilde{P}\| > 0. \quad (5.129)$$

We wish to show that, $S(P) > S(\tilde{P})$. We have,

$$\begin{aligned} S(P) &= mn^2 a_p + (m^2 - m)n^2 a_q \\ &> (mn)^2 a_q \end{aligned} \quad (5.130)$$

where a_p and a_q are defined according to (5.102). Under conditions of Theorem 5.25, we have $b_p > b_{pq} > b_q$ according to (5.102). Thus,

$$S(\tilde{P}) \leq (mn)^2 b_p. \quad (5.131)$$

Using (5.130) and (5.131), we need to show that $a_q > b_p$. We have,

$$b_p - a_q = (\alpha + 1)p^2 - 2p - (\alpha - 1)q. \quad (5.132)$$

This polynomial have two roots at

$$p_{roots} = \frac{1 \pm \sqrt{1 + (\alpha^2 - 1)q}}{\alpha + 1}. \quad (5.133)$$

Because $\alpha > 1$, the minimum root is always negative. Moreover, at $p = 0$, the polynomial value is negative. Thus, (5.132) is negative if

$$0 < p \leq \frac{1 + \sqrt{1 + (\alpha^2 - 1)q}}{1 + \alpha}. \quad (5.134)$$

This completes the proof.

5.8.8 Proof of Theorem 5.28

We use the same setup considered in the proof of Theorem 5.25 to form the expected alignment network. Similarly to the proof of Theorem 5.27, we need to show $a_q > b_p$ according to (5.121). We have,

$$a_q - b_p = \alpha(q(1 - p_e) - p^2(1 + p_e^2)) \quad (5.135)$$

which is positive if

$$p^2 \leq \frac{q(1 - p_e)}{1 + p_e^2}. \quad (5.136)$$

This completes the proof.

Chapter 6

Network Infusion

Information from a single node (entity) can reach other nodes (entities) by propagation over network connections. For instance, a virus infection (either computer or biological) can propagate to different nodes in a network and become an epidemic [231], while rumors can spread in a social network through social interactions [232]. Even a financial failure of an institution can have cascading effects on other financial entities and may lead to a financial crisis [233]. As a final example, in some human diseases, abnormal activities of few genes can cause their target genes and therefore some essential biological processes to fail to operate normally in the cell [234, 235].

In order to gain insight into these processes, mathematical models have been developed, primarily focusing on application to the study of virus propagation in networks ([236, 237]). A well-established continuous-time diffusion model for viral epidemics is known as the susceptible-infected (SI) model [238], where infected nodes spread the virus to their neighbors probabilistically. For that diffusion model, [231, 239–241] explore the relationship between network structure, infection rate, and the size of the epidemics, while [242–244] consider learning SI model parameters. Other diffusion methods use random walks to model information spread and label propagation in networks [245–247]. These references study the forward problem of signal diffusion.

Source inference is the inverse problem. It aims to infer *source nodes* in a network by merely knowing the network structure and observing the information spread at single or multiple snapshots (Figure 6-1). Even within the context of the well-studied diffusion kernels, source inference is a difficult problem in great part owing to the presence of path multiplicity in the network [248]. Recently, the inverse problem of a diffusion process in a network under a discrete time memoryless diffusion model, and when time steps are known, has been considered [249], while the problem of identifying seed nodes (effectors) of a partially activated network in the steady-state of an Independent-Cascade model is investigated in [246]. Moreover reference [250] has considered the source inference problem using incomplete diffusion traces by maximizing the likelihood of the trace under the learned model. The problem setup and diffusion models considered in those works are different than the continuous-time diffusion

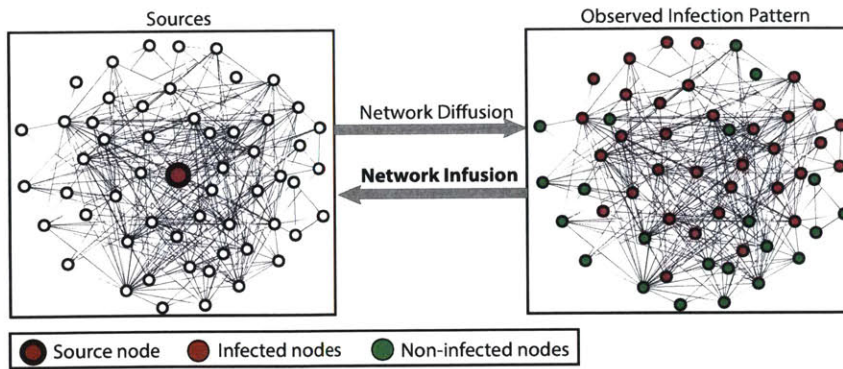


Figure 6-1: Network Infusion (NI) aims to identify source node(s) by reversing information propagation in the network. NI is based on a path-based network diffusion process that closely approximates the observed diffusion pattern, while leading to a tractable source inference method for large complex networks. The displayed network and infection pattern are parts of the Digg social news network.

setup considered in the present paper. The work in [251] uses the Minimum Description Length principle to identify source nodes using a heuristic cost function which combines the model cost and the data cost. Moreover, for the case of having a single source node in the network, some methods infer the source node based on distance centrality [76], or degree centrality [208] measures of the infected subgraph. These methods can be efficiently applied to large networks, but, amongst other drawbacks, their performance lacks provable guarantees in general. For tree structures under a homogeneous SI diffusion model, [252] computes a maximum likelihood solution for the source inference problem and provides provable guarantees for its performance. Over tree structures, their solution is equivalent to the distance centrality of the infected subgraph. The problem of inferring multiple sources in the network has an additional combinatorial complexity compared to the single-source case (see Remark 6.6). Reference [253] has considered this problem under an Independent Cascade (IC) model and introduced a polynomial time algorithm based on dynamic-programming for some special cases. We study this problem for incoherent sources and under the SI model in Section 6.2.3. We review this prior work more extensively in Section 6.1.3.

Source inference in real-world networks is made more challenging as the true propagation dynamics are typically unknown [254], and when they have been directly

measured, they rarely conform to the assumptions of well-studied canonical diffusion models such as SI, partially owing to heterogeneous information diffusion over network edges, latent sources of information, noisy network connections, non-memoryless transport and so forth [255]. As an example, we consider news spread over the Digg social news networks [256] for more than 3,500 news stories. We find that in approximately 65% of cases, nodes who have received the news at a time t , did not have any neighbors who had already received the news by that time violating the most basic conditional independency assumption of the SI model. Furthermore, the empirical distribution of remaining news propagation times over edges of the Digg social news network cannot be approximated closely by a single distribution of a homogenous SI diffusion model, even by fitting a general Weibull distribution to the observed data (Figure 6-13).

Owing to high computational complexity of solving the source inference problem under the well-studied SI diffusion models and considering the fact that those kernels are unlikely to match precisely a real-world diffusion, our key idea to solve the inverse problem is to identify a diffusion process that closely approximates the observed diffusion pattern, but also leads to a tractable source inference method by design. Thus, we develop a diffusion kernel that is distinct from the standard SI diffusion models, but its order of diffusion well-approximates many of them in various setups. We shall show that this kernel leads to an efficient source inference method that can be computed efficiently for large complex networks and shall provide theoretical performance guarantees under general conditions. The key original observation, from both a theoretical and practical perspective, is that in order to solve the inverse problem one does not need to know the full dynamics of the diffusion, , instead to solve the inversion one can do so from statistics that are consistent across many diffusion models.

Instead of the full network, our proposed continuous-time network diffusion model considers k edge-disjoint shortest paths among pairs of nodes, neglecting other paths in the network. We call this kernel a *path-based network diffusion kernel*. Propagation times in the kernel are stochastically independent, which leads to efficient kernel com-

putation even for large complex networks. Using the path-based network diffusion kernel, we propose a computationally tractable general method for source inference called *Network Infusion* (NI), by maximizing the likelihood (NI-ML) or minimizing the prediction error (NI-ME). The NI-ME algorithm is based on an asymmetric Hamming premetric function, and unlike NI-ML, it can be tuned to balance between false positive and false negative error types. Our approach can be computed efficiently for large complex networks, similarly to the distance and degree centrality methods. However, unlike those methods, we provide provable performance guarantees under a continuous-time dynamic SI diffusion setup. We prove that under the SI diffusion model, (i) the maximum-likelihood NI algorithm is mean-field optimal for tree structures, and (ii) the minimum-error NI algorithm is mean-field optimal for regular tree structures. All proofs are presented in the Section 6.6.

Most existing source inference methods consider the case of a single-source homogeneous diffusion setup, partially owing to additional combinatorial complexity of more general cases. We show that the proposed NI framework can be used efficiently for multi-source and heterogeneous diffusion setups under general conditions. Particularly, we prove that the proposed multi-source NI algorithm, which is based on localized likelihood optimization, is mean-field optimal in the regular tree structure for sufficiently-distant sources. Moreover, for a heterogeneous diffusion setup, we show that the network diffusion kernel can be characterized using the phase-type distribution of a Markov chain absorbing time, leading to efficient source inference methods. We also extend our NI algorithms to the cases with unknown or partially known diffusion model parameters such as observation times, by introducing techniques to learn these parameters from observed sample values.

We apply NI to several synthetic networks considering an underlying standard SI diffusion model and compare its performance to existing source inference methods. Having verified the effectiveness of NI both theoretically and through simulations, we then apply it to a real data application to identify the news sources for over 3,500 stories in the Digg social news network. We demonstrate the superior performance of NI compared to existing methods by validating the results based on annotated

information, that was not provided to the algorithms. Finally, we propose an integrative source inference framework which combines source prediction ranks of network infusion and distance centrality and leads to a robust performance in cases where the underlying dynamics are unknown.

6.1 Problem Setup, Computational Difficulties, and Prior Work

In this section, we present the source inference problem and explain its underlying challenges. We also review prior work and present notation used in the rest of the paper.

6.1.1 Source Inference Problem Setup

Let $G = (V, E)$ be a binary, possibly directed, graph representing a network with n nodes, where $G(i, j) = 1$ means that there is an edge from node i to node j (i.e., $(i, j) \in E$) and $G(i, j) = 0$ means there is none. Let $\mathcal{N}(i)$ represent the set of neighbors of node i in the network. If G is directed, $\mathcal{N}(i)$ represents the set of parents of node i . For the sake of description, we illustrate the problem setup and notation in the context of a virus infection spread in the network with the understanding that our framework can be used to solve a more general source inference problem. Let $\mathcal{S} \subset V$ be the set of source nodes in the network. When a node gets infected, it spreads infection to its neighbors, which causes the propagation of infection in the network. Let $\mathcal{T}_{(i,j)}$ be the non-negative, continuous random variable representing the virus traveling time over the edge $(i, j) \in E$. The $\mathcal{T}_{(i,j)}$ variables are assumed to be mutually independent. Let $\mathcal{P}_{i \rightarrow j}$ denote a path, i.e. an ordered set of edges, connecting node i to node j in the network. We define $\mathcal{T}_{\mathcal{P}_{i \rightarrow j}}$ as a random variable representing the virus traveling time over the path $\mathcal{P}_{i \rightarrow j}$, with the following cumulative density

function,

$$F_{\mathcal{P}_{i \rightarrow j}}(t) \triangleq Pr[\mathcal{T}_{\mathcal{P}_{i \rightarrow j}} \leq t]. \quad (6.1)$$

Let $\mathbf{y}(t) \in \{0, 1\}^n$ be the node infection vector at time t , where $y_i(t) = 1$ means that node i is infected at time t . Suppose \mathcal{T}_i is a random variable representing the time that node i gets infected. We assume that if a node gets infected it remains infected (i.e., there is no recovery). Suppose τ_i is a realization of the random variable \mathcal{T}_i . Thus, $y_i(t) = 1$ if $t \geq \tau_i$, otherwise $y_i(t) = 0$. If i is a source node, $\mathcal{T}_i = 0$ and $y_i(t) = 1$ for all $t \geq 0$. The set $V^t = \{i : y_i(t) = 1\}$ represents all nodes that are infected at time t .

Definition 6.1. In a dynamic Susceptible-Infected (SI) diffusion setup, we have

$$\mathcal{T}_i = \min_{j \in \mathcal{N}(i)} (\mathcal{T}_j + \mathcal{T}_{(j,i)}). \quad (6.2)$$

Let $\{\mathbf{y}(t) : t \in (0, \infty)\}$ represent a continuous-time stationary stochastic process of diffusion in the network G . In the source inference problem, given the sample values at times $\{t_1, \dots, t_z\}$ (i.e., $\{\mathbf{y}(t_1), \dots, \mathbf{y}(t_z)\}$), as well as the underlying graph structure G , we wish to infer source nodes that started the infection at time 0. We assume that the number of sources to be inferred (i.e., m) and the observation time stamps (i.e., $\{t_1, \dots, t_z\}$) are also given. We discuss the cases with unknown or partially known parameters in Section 6.2.5.

One way to formulate the source inference problem is to use a standard maximum a posteriori (MAP) estimation.

Definition 6.2 (MAP Source Inference). The MAP source inference solves the following optimization:

$$\begin{aligned} \arg \max \quad & Pr(\mathbf{y}(0) | \mathbf{y}(t_1), \dots, \mathbf{y}(t_z)), \\ & \|\mathbf{y}(0)\|_{l_0} = m, \end{aligned} \quad (6.3)$$

where m is the number of source nodes in the network, and $\|\cdot\|_{l_0}$ represents the l_0 norm of a vector.

In some applications, there may be nonuniform prior probabilities for different candidate source nodes. The MAP source inference optimization takes into account these prior probabilities as well. If there is no informative prior probabilities for candidate source nodes, the MAP source Optimization (6.3) can be simplified to the following maximum likelihood (ML) source estimation:

Definition 6.3 (ML Source Inference). The ML source inference solves the following optimization:

$$\begin{aligned} \arg \max \quad & Pr(\mathbf{y}(t_1), \dots, \mathbf{y}(t_z) | \mathbf{y}(0)), \\ & \|\mathbf{y}(0)\|_{l_0} = m, \end{aligned} \tag{6.4}$$

where its objective function is an ML function (score) of source candidates.

An alternative formulation for the source inference problem is based on minimizing the prediction error. In Section 6.2.2, we propose a minimum prediction error formulation that uses an asymmetric Hamming pre-metric function and can balance between false positive and false negative error types by tuning a parameter.

6.1.2 Computational Difficulties of the Source Inference Problem

In this section, we explain underlying challenges of the source inference problem.

Remark 6.4. Suppose the underlying network G has 4 nodes and 3 undirected edges as depicted in Figure 6-2. Suppose the underlying diffusion is according to the SI model of Definition 6.1. Let the edge holding time variables $\mathcal{T}_{(i,j)}$ be mutually independent and be distributed according to an exponential distribution with parameter

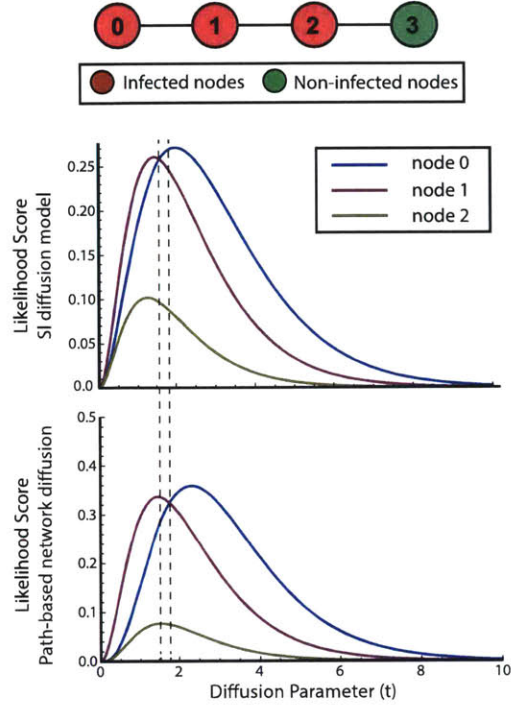


Figure 6-2: Likelihood scores based on the SI diffusion model and the path-based network diffusion kernel for a depicted line network. Although the underlying diffusion model is SI, NI optimization leads to a tight approximation of the exact solution for a wide range of parameter t (except the range between two vertical dashed lines). Unlike SI, NI computation can be efficiently extended to large complex networks, owing to decoupling of its likelihood terms.

λ :

$$f_{i,j}(\tau_{i,j}) = \lambda e^{-\lambda \tau_{i,j}}, \quad \forall (i,j) \in E. \quad (6.5)$$

Without loss of generality, let $\lambda = 1$. Suppose there is a single source in the network (i.e., $m = 1$), and we observe the infection pattern at a single snapshot at time t . Let the observed infection pattern at time t be $\mathbf{y}(t) = (1, 1, 1, 0)$, implying that nodes $\{0, 1, 2\}$ are infected at time t , while node 3 is not infected at that time. Our goal is to find the most likely source node, according to the ML Optimization 6.4. Let τ_i be a realization of the random variable \mathcal{T}_i (i.e., the time that virus arrives at node i). If node 0 was the source node (i.e., $\tau_0 = 0$), we would have $\tau_1 \leq \tau_2 \leq t \leq \tau_3$ because of

the underlying network structure. Thus,

$$\begin{aligned}
 Pr(\mathbf{y}(t) = (1, 1, 1, 0) | \mathbf{y}(0) = (1, 0, 0, 0)) & \quad (6.6) \\
 &= \int_{\tau_1=0}^t \int_{\tau_2=\tau_1}^t \int_{\tau_3=t}^{\infty} e^{-\tau_1} e^{-(\tau_2-\tau_1)} e^{-(\tau_3-\tau_2)} d\tau_1 d\tau_2 d\tau_3 \\
 &= \int_{\tau_1=0}^t \int_{\tau_2=\tau_1}^t \int_{\tau_3=t}^{\infty} e^{-\tau_3} d\tau_1 d\tau_2 d\tau_3 \\
 &= \frac{1}{2} t^2 e^{-t}.
 \end{aligned}$$

Similarly, we have,

$$\begin{aligned}
 Pr(\mathbf{y}(t) = (1, 1, 1, 0) | \mathbf{y}(0) = (0, 1, 0, 0)) &= t(1 - e^{-t})e^{-t}, \quad (6.7) \\
 Pr(\mathbf{y}(t) = (1, 1, 1, 0) | \mathbf{y}(0) = (0, 0, 1, 0)) &= e^{-t} - (1 + te^{-2t}).
 \end{aligned}$$

These likelihood functions are plotted in 6-2. For a given observation time stamp t , the ML source estimator selects the node with the maximum likelihood score as the source node, according to Optimization 6.4. Note that an optimal source solution depends on the observation time parameter t (i.e., for $t \lesssim 1.6$, node 1 and for $t \gtrsim 1.6$, node 0 are ML optimal source nodes.)

Remark 6.5. Suppose G is a network with 5 nodes and 5 edges as shown in Figure 6-3-a. Consider the same diffusion setup as the one of Remark 6.4. Let $\mathbf{y}(t) = (1, 1, 1, 1, 0)$; i.e., nodes $\{0, 1, 2, 3\}$ are infected at time t while node 4 is not infected at that time. Similarly to Remark 6.4, let τ_i be a realization of the random variable \mathcal{T}_i , a variable representing the time that virus arrives at node i . If node 0 was the source node (i.e., $\tau_0 = 0$), we would have $\tau_1 \leq \min(\tau_2, \tau_3) \leq \tau_4$, $\max(\tau_2, \tau_3) \leq t$, and $\tau_4 > t$. Thus,

$$\begin{aligned}
 Pr(\mathbf{y}(t) = (1, 1, 1, 1, 0) | \mathbf{y}(0) = (1, 0, 0, 0, 0)) & \quad (6.8) \\
 &= 2e^{-t} - e^{-2t}(1 + (1+t)^2).
 \end{aligned}$$

In this case, likelihood computation is more complicated than the case of Remark 6.4, because both variables \mathcal{T}_2 and \mathcal{T}_3 depend on \mathcal{T}_1 , and therefore, consecutive terms do not cancel as in (6.6). Moreover, note that there are two paths from node 0 to node 4 that overlap at edge $(0, 1)$. As we have mentioned earlier, such overlaps are a source of difficulty in the source inference problem, which is illustrated by this simplest example, because the variable \mathcal{T}_4 depends on both variables \mathcal{T}_2 and \mathcal{T}_3 through a $\min(., .)$ function which makes computation of the likelihood integral complicated.

Remark 6.6. Remarks 6.4 and 6.5 explain underlying source inference challenges for a single source node in the network. The case of having multiple source nodes has additional complexity because likelihood scores of Optimization (6.4) should be computed for all possible subsets of infected nodes. For the case of having m sources in the network, there are $\binom{|V^t|}{m}$ candidate source sets where for each of them a likelihood score should be computed. If there are a significant number of infected nodes in the network (i.e., $V^t = \mathcal{O}(n)$), there would be $\mathcal{O}(n^m)$ source candidate sets. This makes the multi-source inference problem computationally expensive for large networks, even for small values of m .

Moreover, in Remarks 6.4 and 6.5, we assume that the edge holding time distribution is known and follows an exponential distribution with the same parameter for all edges. This is the standard diffusion model used in most epidemic studies [239], because the exponential distribution has a single parameter and is memoryless. However, in some practical applications, the edge holding time distribution may be unknown and/or may vary for different edges. We discuss this case in Section 6.2.4.

6.1.3 Prior Work

While our approach considers a general network diffusion setup and its inverse problem, most of the literature considers application to specific problems. The most common ones focus on studying different models of virus propagation in population networks. A standard information diffusion model in this setup is known as the susceptible-infected-recovered (SIR) model [238]. At any time, nodes have three types

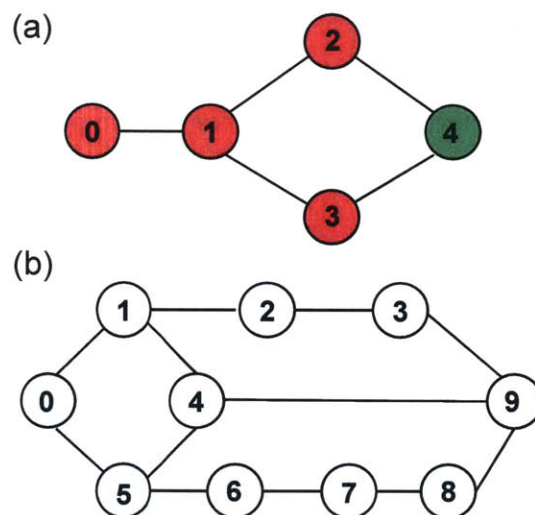


Figure 6-3: (a) A graph considered in Remark 6.5. (b) An example graph with overlapping shortest paths between nodes 0 and 9.

in this model: susceptible nodes which are capable of getting infected, infected nodes that spread virus in the network, and recovered nodes that are cured and can no longer become infected. Under the SIR diffusion model, infection spreads from sources to susceptible nodes probabilistically. References [231, 239–241] discuss the relationship among network structure, infection rate, and the size of the epidemics under this diffusion model. Learning different diffusion parameters of this model have been considered in references [242–244]. Some other diffusion methods use random walks to model information spread and label propagation in networks [245–247]. In these models, a random walker goes to a neighbor node with a probability inversely related to node degrees. Therefore, high degree nodes may be less influential in information spread in the network which may be counter-intuitive in some applications.

Although there are several works on understanding mechanisms of diffusion processes in different networks, there is somehow little work on studying the inverse diffusion problem to infer information sources, in great part owing to the presence of path multiplicity in the network [248], described in Remarks 6.4, 6.5, and 6.6. Recently, [249] considers the inverse problem of a diffusion process in a network under a discrete time memoryless diffusion model when time steps are known, while ref-

erence [246] investigates the problem of identify seed nodes (effectors) of a partially activated network in the steady-state of an Independent-Cascade model. Moreover reference [250] has studied this problem using incomplete diffusion traces by maximizing the likelihood of the trace under the learned model. The diffusion models and the problem setup considered in these works are different from the continuous-time dynamic diffusion setup considered in this paper. Furthermore, reference [251] uses the Minimum Description Length principle to identify source nodes in the network. Their method uses a heuristic cost function which combines the model cost and the data cost and lacks theoretical performance characterization.

For the case of having a single source node in the network, some methods infer the source node based on distance centrality [76], or degree centrality [208] measures of the infected subgraph. These methods are efficient to apply to large networks. However, because they do not assume any particular diffusion model, their performance lacks provable guarantees in general. For tree structures and under a homogeneous SI diffusion model, reference [252] computes a maximum likelihood solution for the source inference problem and provides provable guarantees for its performance. Over tree structures, their solution is in fact equivalent to the distance centrality of the infected subgraph.

The multi-source inference problem has an additional combinatorial complexity compared to the single-source case (see Remark 6.6). Reference [253] has studied the complexity of this problem for different types of graphs under an Independent Cascade (IC) model. It shows that for arbitrary graphs this problem is NP-hard to even approximate. However, for some special cases, they introduce a polynomial time algorithm based on dynamic-programming to solve the multi-source problem under the IC model. We study this problem under the SI model for incoherent sources in Section 6.2.3.

6.2 Main Results

In this section, first, we introduce a path-based network diffusion kernel which is used in proposed source inference methods. Then we present algorithms to infer single and multiple sources in the network. Finally, we present NI algorithms for heterogeneous diffusion and non-parametric cases. For the sake of description, we shall, as before, have a recurrent example of the virus infection spread in the network with the understanding that our framework can be used to solve a more general source inference problem.

6.2.1 The path-based Network Diffusion Kernel

Consider the case when there exists a single source node in the network. The multi-source generalization treated in Section 6.2.3. Suppose the network structure $G = (V, E)$ is known, and a snapshot $\mathbf{y}(t)$ from some *real* diffusion dynamics is observed at time t . In general, there may be several diffusion processes that lead to the observed infection snapshot in the network, either exactly or approximately. Suppose $\hat{\mathbf{y}}(t')$ is the sample generated at time t' using a certain diffusion model. One way to characterize the error of this diffusion model to explain the observed diffusion sample is to use an asymmetric Hamming premetric function as follows:

$$\min_{t'} h_{\alpha}(\mathbf{y}(t), \hat{\mathbf{y}}(t')) \triangleq (1 - \alpha) \sum_{i: y_i(t)=1} \mathbf{1}_{\hat{y}_i=0} + \alpha \sum_{i: y_i(t)=0} \mathbf{1}_{\hat{y}_i(t')=1}, \quad (6.9)$$

where $0 \leq \alpha \leq 1$. This error metric assigns weight α to false positive and weight $1 - \alpha$ to false negatives error types. To solve the inverse problem, one may select a diffusion process which approximates the observed diffusion pattern closely and also leads to a tractable source inference method. Although the SI diffusion model may be well-suited to model the forward problem of information diffusion in the network, solving the inverse problem (the source inference problem) under this model is challenging in general, in great part owing to the presence of path multiplicity in the network,

as we explain in Remarks 6.4, 6.5 and 6.6. Here, we present a path-based network diffusion kernel that is distinct from the standard SI diffusion models, but its order of diffusion approximates well many of them. We will show that this kernel leads to an efficient source inference method with theoretical performance guarantees, under some general conditions.

In our diffusion model, instead of the full network, we consider up to k edge-disjoint shortest paths among pairs of nodes, neglecting other paths in the network. Suppose $\mathcal{P}_{i \rightarrow j}^1, \mathcal{P}_{i \rightarrow j}^2, \dots$ represent different paths between nodes i and j in the network. The length of a path $\mathcal{P}_{i \rightarrow j}^r$ is denoted by $|\mathcal{P}_{i \rightarrow j}^r|$. Let $E_{i \rightarrow j}^r$ be the set of edges of the path $\mathcal{P}_{i \rightarrow j}^r$. We say two paths are edge-disjoint if the set of their edges do not overlap. Let $\{\mathcal{P}_{i \rightarrow j}^1, \mathcal{P}_{i \rightarrow j}^2, \dots, \mathcal{P}_{i \rightarrow j}^k\}$ represent k disjoint shortest paths between nodes i and j . We choose these paths iteratively so that

- $|\mathcal{P}_{i \rightarrow j}^1| \leq |\mathcal{P}_{i \rightarrow j}^2| \leq \dots \leq |\mathcal{P}_{i \rightarrow j}^k|$,
- paths are disjoint. I.e., for $1 < r \leq k$, $E_{i \rightarrow j}^r \cap (\bigcup_{a=1}^{r-1} E_{i \rightarrow j}^a) = \emptyset$,
- $\mathcal{P}_{i \rightarrow j}^r$ is a shortest path between nodes i and j in the network $G' = (V, E - \bigcup_{a=1}^{r-1} E_{i \rightarrow j}^a)$.

In some cases, the shortest path solutions may not be unique. That is, there are at least two shortest paths connecting nodes i to j in the network. If these shortest paths do not overlap, the resulting path length vector $(|\mathcal{P}_{i \rightarrow j}^1|, \dots, |\mathcal{P}_{i \rightarrow j}^k|)$ is the same irrespective of the selection order. Thus, the tie breaking can be done randomly. However, in the case of having overlapping shortest paths, one way to break the tie among these paths is to choose the one which leads to a shorter path in the next step. For example, consider the network depicted in Figure 6-3-b. There are two paths of length 3 between nodes 0 and 9. Choosing the path $0 - 5 - 4 - 9$ leads to the next independent path $0 - 1 - 2 - 3 - 9$ with length 4, while choosing the path $0 - 1 - 4 - 9$ leads to the next path $0 - 5 - 6 - 7 - 8 - 9$ of length 5. Thus, the algorithm chooses the path $0 - 5 - 4 - 9$. If next paths have the same length, tie would be broken considering more future steps. In practice, this case has negligible effect in the

performance of the source inference method. Methods based on message-passing or dynamic programming can be used to select optimal k shortest paths in the network as well [257, 258]. In this paper, we break ties randomly among paths with the same length.

Recall that $\mathcal{T}_{\mathcal{P}_{i \rightarrow j}^r}$ represents the virus traveling time over the path $\mathcal{P}_{i \rightarrow j}^r$ whose cumulative density function is denoted by $F_{\mathcal{P}_{i \rightarrow j}^r}(\cdot)$ according to Equation (6.1).

Definition 6.7 (Path-based network diffusion kernel). Let $p_{i,j}(t)$ be the probability of node j being infected at time t if node i is the source node. Thus

$$p_{i,j}(t) = Pr[y_j(t) = 1 | y_i(0) = 1] \quad (6.10)$$

$$\triangleq 1 - \prod_{r=1}^k (1 - F_{\mathcal{P}_{i \rightarrow j}^r}(t)),$$

where k is the number of independent shortest paths between nodes i and j . $P(t) = [p_{i,j}(t)]$ is called a path-based network diffusion kernel.

In the path-based network diffusion kernel node j gets infected at time t if the infection reaches to it over at least one of the k independent shortest paths connecting that node to the source node. The path-based network diffusion kernel provides a non-dynamic diffusion basis for the network and is based on two important assumptions: that edge holding time variables $\mathcal{T}_{(i,j)}$ are mutually independent; and the paths are disjoint. A path-based network diffusion kernel with $k = 1$ only considers the shortest paths in the network and has the least computational complexity among other path-based network diffusion kernels. Considering more paths among nodes in the network (i.e., $k > 1$) provides a better characterization of network diffusion processes with the cost of increased kernel computational complexity (Proposition 6.8). Computation of the path-based network diffusion kernel compared to the SI diffusion model is provably efficient even for large and complex networks.

In the following, we highlight properties and relaxations of the path-based network diffusion kernel:

The path-based network diffusion kernel provides a non-dynamic diffusion model, where nodes become infected independently based on their distances (path lengths) to source nodes. Unlike the dynamic SI model, in the path network diffusion model, it is possible (though unlikely) to have $y_i(t) = 1$ while $y_j(t) = 0$, for all neighbors of node i (i.e., $j \in \mathcal{N}(i)$). The key idea is that to infer the source node in the network, full characterization of diffusion dynamics, in many cases, may not be necessary as long as the diffusion model approximates the observed samples closely (e.g., according to an error metric of (6.9)). For instance, consider the setup of Remark 6.4 where the underlying diffusion model is according to a SI model. In that example, we compute source likelihood scores in (6.6) and (6.7) by integrating likelihood conditional density functions. The likelihood computation under this model becomes challenging for complex networks. However, according to the path-based network diffusion model of Definition 6.7, these likelihood scores are decoupled to separate terms and can be computed efficiently as follows:

$$\begin{aligned}
 Pr(\mathbf{y}(t) = (1, 1, 1, 0) | \mathbf{y}(0) = (1, 0, 0, 0)) &= F(1, t)F(2, t)\bar{F}(3, t), & (6.11) \\
 Pr(\mathbf{y}(t) = (1, 1, 1, 0) | \mathbf{y}(0) = (0, 1, 0, 0)) &= F(1, t)^2\bar{F}(2, t), \\
 Pr(\mathbf{y}(t) = (1, 1, 1, 0) | \mathbf{y}(0) = (0, 0, 1, 0)) &= F(1, t)F(2, t)\bar{F}(1, t),
 \end{aligned}$$

where $F(l, t)$ is the Erlang cumulative distribution function over a path of length l , that we shall show in (6.14). Figure 6-2 shows likelihood scores of infected nodes computed according to (6.11). This example illustrates that for a wide range of parameter t , both models lead to the same optimal solution. Moreover, unlike the SI model, likelihood functions can be computed efficiently using the path-based network diffusion kernel, even for large complex networks.

The path-based diffusion kernel considers only the top k shortest paths among nodes, neglecting other paths in the networks. The effect of long paths is dominated by the one of short ones leading to a tight approximation. Suppose $\mathcal{P}_{i \rightarrow j}^1$ and $\mathcal{P}_{i \rightarrow j}^2$ represents two paths between nodes i and j where $|\mathcal{P}_{i \rightarrow j}^1| \lambda |\mathcal{P}_{i \rightarrow j}^2|$ (i.e., the path $\mathcal{P}_{i \rightarrow j}^2$

is much longer than the path $\mathcal{P}_{i \rightarrow j}^1$). Thus, for a wide range of parameter t , we have $F_{\mathcal{P}_{i \rightarrow j}^1}(t) \Gamma F_{\mathcal{P}_{i \rightarrow j}^2}(t)$, and therefore

$$(1 - F_{\mathcal{P}_{i \rightarrow j}^1}(t))(1 - F_{\mathcal{P}_{i \rightarrow j}^2}(t)) \approx 1 - F_{\mathcal{P}_{i \rightarrow j}^1}(t). \quad (6.12)$$

Note that for very small or large t values (i.e., $t \rightarrow 0$ or $t \rightarrow \infty$), both $F_{\mathcal{P}_{i \rightarrow j}^1}(\cdot)$ and $F_{\mathcal{P}_{i \rightarrow j}^2}(\cdot)$ go to 0 and 1, respectively, and thus the approximation (6.12) remains tight. For an example network depicted in Figure 6-4-a, we illustrate the tightness of the first order approximation (i.e., $k = 1$) for different lengths of the path $\mathcal{P}_{i \rightarrow j}^2$. In general for large k values the gap between the approximate and the exact kernels becomes small with the cost of increased kernel computational complexity (see Proposition 6.8). The same approximation holds for overlapping paths with different lengths as it is illustrated in Figures 6-4-b.

Finally, path-based network diffusion kernel only considers independent shortest paths among nodes and therefore ignores the effects of non-disjoint paths in the network. This is a critical relaxation because, as we explain in Remark 6.5, overlapping paths and dependent variables make the likelihood computation and therefore source inference challenging. In general, if there are many overlapping shortest paths among nodes in the network, this approximation might not be tight. However, in network structures whose paths do not overlap significantly (for example tree structures), this approximation is tight. In Section 6.2.1, we show that in a common model for sparse random networks [199], shortest paths among nodes that are the dominating factors in the path-based network diffusion kernel computation, are extremely unlikely to overlap with each other, leading to a tight kernel approximation.

One main advantage of using the path-based network diffusion kernel compared to other diffusion models such as the SI diffusion model is its efficient computation even for large and complex networks:

Proposition 6.8 (Computational complexity of path-based network diffusion kernel). Let $G = (V, E)$ be a directed network with n nodes and $|E|$ edges. Then, computation

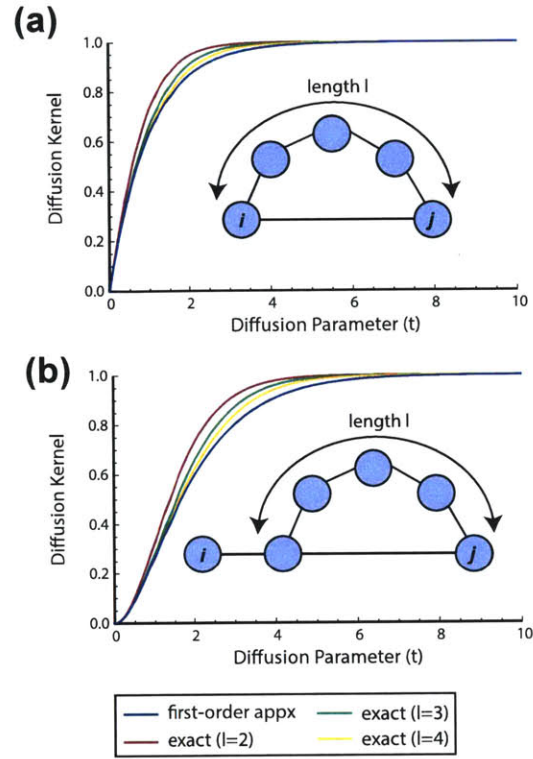


Figure 6-4: Tightness of the first order approximation of the path-based network diffusion kernel over example networks.

of the k -path network diffusion kernel of Definition 6.7 has a worst case computational complexity $\mathcal{O}(k|E|n + kn^2 \log(n))$.

Remark 6.9. To solve the source inference problem, one only needs to compute rows of the path-based network diffusion kernel which correspond to infected nodes ($i \in V^t$). Thus, time complexity of kernel computation can be reduced to $\mathcal{O}(k|E||V^t| + k|V^t|n \log(n))$, where $|V^t|$ is the number of observed infected nodes in the network at time t .

Computation of the path-based network diffusion kernel depends on edge holding time distributions. If the virus traveling time variables $\mathcal{T}_{(i,j)}$ are identically distributed for all edges in the network, the underlying diffusion process is called *homogeneous*. On the other hand, if holding time distributions differ among edges in the network, the resulting diffusion process is *heterogeneous*. In this section, we consider a homogeneous diffusion setup, where the holding time distribution is an exponential

distribution with parameter λ for all edges. Without loss of generality, we assume that $\lambda = 1$. It is important to note that the proposed framework is not limited to use of either the exponential or a single delay distribution. Other edge holding time distributions can be used in the framework as well. For example, the case of heterogeneous diffusion is considered in Section 6.2.4.

Under the setup considered in this section, the virus traveling time over each path in the network has an Erlang distribution, because it is the sum of independent exponential variables. Thus, we have that

$$F_{\mathcal{P}_{i \rightarrow j}^r}(t) = Pr[\mathcal{T}_{\mathcal{P}_{i \rightarrow j}^r} \leq t] = \frac{\gamma(|\mathcal{P}_{i \rightarrow j}^r|, \lambda t)}{(|\mathcal{P}_{i \rightarrow j}^r| - 1)!}, \quad (6.13)$$

where $\gamma(\cdot)$ is the lower incomplete gamma function [259]. $|\mathcal{P}_{i \rightarrow j}^r|$ (the path length connecting node i to j) is also called the Erlang's shape parameter. Because $F_{\mathcal{P}_{i \rightarrow j}^r}(t)$ is only a function of the path length and parameter t , to simplify notation, we define

$$F(l, t) \triangleq F_{\mathcal{P}_{i \rightarrow j}^r}(t), \quad (6.14)$$

where $l = |\mathcal{P}_{i \rightarrow j}^r|$. The k -path network diffusion kernel of Definition 6.7 using the Erlang distribution is called a path-based Erlang network diffusion kernel. If only one shortest path among nodes is considered (i.e., $k = 1$), the diffusion kernel of Definition 6.7 is called the shortest path network diffusion kernel.

Definition 6.10 (Shortest path Erlang diffusion kernel). The shortest path Erlang network diffusion kernel is defined as follows:

$$p_{i,j}(t) = Pr[y_j(t) = 1 | y_i(0) = 1] = F(d_{i,j}, t), \quad (6.15)$$

where $d_{i,j}$ is the length of the shortest path connecting node i to node j , and $F(d_{i,j}, t)$ represents the Erlang cumulative distribution function of (6.14).

The shortest path Erlang diffusion kernel can be viewed as the first order approxima-

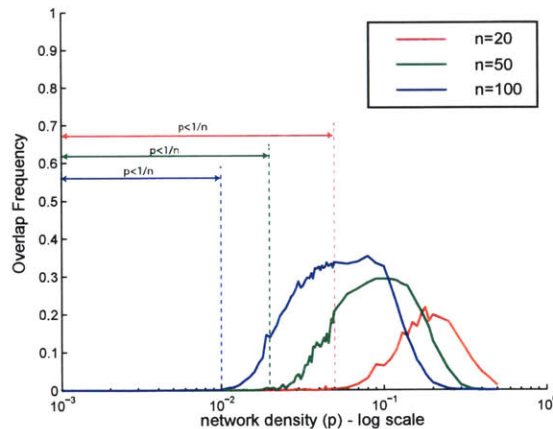


Figure 6-5: The frequency of having overlapping shortest paths between two randomly selected nodes over an Erdős-Rényi graph with parameter p . In sparse graph ($p \leq \frac{1}{n}$), the overlap frequency is small. Experiments are repeated 20,000 times for each case.

tion of the underlying diffusion process. It has the least computational complexity among other path-based network diffusion kernels which makes it suitable to be used over large and complex networks. Moreover, this kernel has a single parameter t which can be learned reliably using the observed samples (see Section 6.2.5).

Tightness of the Path-based Network Diffusion Kernel for Sparse Erdős-Rényi Graphs

The path-based network diffusion kernel of Definition 6.10 does not consider effect of overlapping paths among nodes. This can cause the SI approximation to be loose particularly if there are several overlapping shortest paths among nodes in the network since shortest paths are dominant factors in kernel computation. In the following proposition we show that in a sufficiently sparse Erdős-Rényi graph, the probability of having overlapping shortest paths among nodes is small.

Proposition 6.11. Let $G = (V, E)$ be an undirected Erdős-Rényi graph with n nodes where $Pr[(i, j) \in E] = p$. Consider two nodes i and j where $d(i, j) \leq l_0$. If $p < \frac{c}{n}$ where $c = 1/n^{2l_0}$, the probability of having overlapping shortest paths between nodes i and j goes to zero asymptotically.

Figure 6-5 illustrates Proposition 6.11 for Erdős-Rényi graphs with different num-

ber of nodes and different parameters p . As illustrated in this figure, shortest paths are less likely to overlap in sparse networks. Moreover, in very dense networks, which practically might be less interesting, the shortest path overlap probability decreases as well because most node pairs are connected by one-hop or two-hop paths.

6.2.2 Global Source Inference

In this section, we describe a source inference method called Network Infusion (NI) which aims to solve the inverse diffusion problem over a given network using observed infection patterns. The method described in this section finds a single node as the source of the *global* information spread in the network. In Section 6.2.3, we consider the case when more than one source node exists, where each source causes a *local* infection propagation in the network. In this section, we also assume that the infection pattern is observed at a single snapshot at time t (i.e., $\mathbf{y}(t)$ is given). The case of having multiple snapshots is considered in Section 6.3.

Recall that V^t is the set of observed infected nodes at time t , and $P(t) = [p_{i,j}(t)]$ represents the path-based network diffusion kernel according to Definition 6.7. Under the path-based network diffusion kernel, the joint diffusion probability distribution can be decoupled into individual marginal distributions, which leads to a tractable ML Optimization, even for complex networks:

Algorithm 6.12. A maximum-likelihood NI algorithm (NI-ML) infers the source node by solving the following optimization:

$$\begin{aligned} & \arg \max_{i \in V^t} \mathcal{L}(i, t) & (6.16) \\ & = \arg \max_{i \in V^t} \sum_{j \in V^t} \log(p_{i,j}(t)) + \sum_{j \notin V^t} \log(1 - p_{i,j}(t)). \end{aligned}$$

In Optimization (6.16), $\mathcal{L}(i, t)$ is called the NI log-likelihood function (score) of node i at time t . In this optimization, the parameter t is the time at which the observation is made and is assumed to be known. For cases that this parameter is unknown, we introduce techniques to learn a within-diffusion-model time parameter

with provable performance guarantees in various setups (Section 6.2.5).

To have a well-defined NI log-likelihood objective function, $p_{i,j}(t)$ should be non-zero for infected nodes i and j (i.e., $p_{i,j}(t) \neq 0$ when $i, j \in V^t$). If infected nodes form a strongly connected sub-graph over the network, this condition is always satisfied. In practice, if $p_{i,j}(t) = 0$ for some $i, j \in V^t$ (i.e., i and j are disconnected in the graph), we assume that $p_{i,j}(t) = \epsilon$ (e.g., $\epsilon = 10^{-6}$). Note that for $j \notin V^t$, for any value of $t > 0$, $p_{i,j}(t) < 1$. Therefore, the second term in the summation of Optimization (6.16) is always well-defined.

NI Algorithm 6.12 aims to infer the source node by maximizing $\mathcal{L}(i, t)$. An alternative approach is to infer the source node by minimizing the expected prediction error of the observed infection pattern. We describe this approach in the following: Let $h_\alpha(\mathbf{y}, \mathbf{x})$ be a weighted Hamming premetric between two binary sequences \mathbf{x} and \mathbf{y} defined as follows:

$$h_\alpha(\mathbf{y}, \mathbf{x}) \triangleq (1 - \alpha) \sum_{i:y_i=1} \mathbf{1}_{x_i=0} + \alpha \sum_{i:y_i=0} \mathbf{1}_{x_i=1}, \quad (6.17)$$

where $0 \leq \alpha \leq 1$. If $\alpha = 1/2$, $h_\alpha(\cdot, \cdot)$ is a metric distance. If $\alpha \neq 1/2$, $h_\alpha(\cdot, \cdot)$ is a premetric (not a metric) because it does not satisfy the symmetric property of distance metrics (i.e., there exist \mathbf{x} and \mathbf{y} such that $h_\alpha(\mathbf{x}, \mathbf{y}) \neq h_\alpha(\mathbf{y}, \mathbf{x})$), and it does not satisfy the triangle inequality as well (i.e., there exist \mathbf{x} , \mathbf{y} and \mathbf{z} such that $h_\alpha(\mathbf{x}, \mathbf{y}) > h_\alpha(\mathbf{y}, \mathbf{z}) + h_\alpha(\mathbf{z}, \mathbf{x})$).

Remark 6.13. $h_\alpha(\mathbf{y}, \mathbf{x})$ generalizes Hamming distance between binary sequences \mathbf{y} and \mathbf{x} using different weights for different error types. Suppose \mathbf{x} is a prediction of the sequence \mathbf{y} . There are two types of possible errors: (1) if $y_i = 1$ and the prediction is zero (i.e., $x_i = 0$, false negative error), (2) if $y_i = 0$ and the prediction is one (i.e., $x_i = 1$, false positive error). $h_\alpha(\mathbf{y}, \mathbf{x})$ combines these errors by assigning weight $1 - \alpha$ to false negatives (i.e., missing ones), and weight α to false positives (i.e., missing zeros). Having different weights for different error types can be useful, specially if the sequence \mathbf{y} is sparse. Suppose \mathbf{y} has κ ones (positives), and $n - \kappa$ zeros (negatives).

Therefore, there are κ possible type 1 errors and $n - \kappa$ type 2 errors in prediction. In this case, to have balance between the number of true negative and false positive errors, one can choose $\alpha = \kappa/n$ in calculation of $h_\alpha(\mathbf{y}, \mathbf{x})$.

Now we introduce a NI algorithm which infers the source node by minimizing the prediction error.

Algorithm 6.14. A minimum error NI algorithm (NI-ME) infers the source node by solving the following optimization:

$$\begin{aligned} \arg \min_{i \in V^t} \mathcal{H}_\alpha(i, t) &= \arg \min_{i \in V^t} \mathbb{E}[h_\alpha(\mathbf{y}(t), \mathbf{x}_i(t))] & (6.18) \\ &= \arg \min_{i \in V^t} (1 - \alpha) \sum_{j \in V^t} (1 - p_{i,j}(t)) + \alpha \sum_{j \notin V^t} p_{i,j}(t), \end{aligned}$$

where $\mathbf{x}_i(t)$ is a binary prediction vector of node i at time t with probability distribution $P_i(t)$, and $\mathcal{H}_\alpha(i, t)$ is the expected prediction error of node i at time t .

Similarly to Maximum Likelihood NI (NI-ML) Algorithm, we assume that the parameter t (the time at which observation is made) is known. We discuss the case when this parameter is unknown in Section 6.2.5.

Remark 6.15. According to Remark 6.13, to have balance between false positive and false negative error types, one can use $\alpha = |V^t|/n$ where $|V^t|$ is the number of infected nodes (positives) at time t . However, in general, this parameter can be tuned in different applications using standard machine learning techniques such as cross validations [222].

The proposed NI methods based on maximum likelihood (NI-ML, Algorithm 6.12) and minimum error (NI-ME, Algorithm 6.14) are efficient to solve even for large complex networks:

Proposition 6.16. Suppose the underlying network $G = (V, E)$ has n nodes and $|E|$ edges. Let V^t represent the set of infected nodes at time t . Then, a worst case computational complexity of NI Algorithms 6.12 and 6.14 is $\mathcal{O}(|V^t|(k|E| + kn \log(n)))$.

In the rest of this section, we analyze the performance of NI Algorithms 6.12 and 6.14 under a standard SI diffusion model of Definition 6.1.

Theorem 6.17. Let $G = (V, E)$ be an undirected tree with countably infinite nodes. Suppose node s is the source node, t is the infection observation time, and the underlying diffusion process is according to the SI model of Definition 6.1. Then, we have,

$$\mathbb{E}[\mathcal{L}(s, t)] \geq \mathbb{E}[\mathcal{L}(i, t')], \quad \forall i, \forall t', \quad (6.19)$$

where $\mathbb{E}[\mathcal{L}(i, t')]$ is the expected NI log-likelihood score of node i with parameter t' .

In the setup of Theorem 6.17, similarly to the setup of reference [252], we assume that the set of vertices is countably infinite to avoid boundary effects. Theorem 6.17 provides a mean-field (expected) optimality for Algorithm 6.12. In words, it considers the case when we have sufficient samples from independent infection spreads in the network starting from the same source node.

Note that the function $\mathcal{L}(i, t)$ provides a statistics of node i being the source node. Theorem 6.17 considers the expectation of this function under the SI diffusion model. Since both expectation and summation are linear operations, the graph structure is tree, and in $\mathcal{L}(i, t)$ nodes have been decoupled, computing this expectation under the SI model and the path-based network diffusion kernel leads to the same argument (for more details, see the proof of the Theorem). However by considering the expectation under the path-based network diffusion kernel, the argument of Theorem 6.17 holds for a general graph structure.

In (6.19), i can be equal to s (the source node), and/or t' can be equal to t as well. If i is equal to s , we have the following:

Corollary 6.18. Under the conditions of Theorem 6.17, we have,

$$\mathbb{E}[\mathcal{L}(s, t)] \geq \mathbb{E}[\mathcal{L}(s, t')], \quad \forall t', \quad (6.20)$$

where the equality holds iff $t = t'$.

Remark 6.19. In this remark, we highlight the difference between parameters t and t' in Theorem 6.17. The parameter t is the time at which we observe the infection pattern in the network. If this parameter is known, it can be used to compute likelihood scores according to Optimization (6.16). However, this parameter may be unknown and one may use an estimate of this parameter in Optimization (6.16) (i.e., using t' instead of t). Theorem 6.17 indicates that even if different parameters $t' \neq t$ are used to compute source likelihood scores for different nodes, the likelihood score obtained by the source node s and the true parameter t is optimal in expectation. This theorem and corresponding Proposition 6.18 provide a theoretical basis to estimate the underlying true parameter t by maximizing the likelihood score for each node over different values of t' (for more details, see Section 6.2.5).

In the following, we present the mean-field optimality of minimum error NI algorithm (NI-ME) over regular tree structures:

Theorem 6.20. Let $G = (V, E)$ be a regular undirected tree with countably infinite nodes. Suppose node s is the source node, t is the observation time, and the underlying diffusion process is according to the SI model of Definition 6.1. Then, for any value of $0 < \alpha < 1$ and $t' > 0$, we have,

$$\mathbb{E}[H_\alpha(s, t')] < \mathbb{E}[H_\alpha(i, t')], \quad \forall i \neq s, \forall t' > 0, \quad (6.21)$$

where $\mathbb{E}[H_\alpha(i, t')]$ is the expected prediction error of node i using parameter t' . Equality (6.21) holds iff $s = i$.

The mean field optimality of NI-ME algorithm holds for all values of $0 < \alpha < 1$ under the setup of Theorem 6.20. In practice and under more general conditions, we find that α selection according to Remarks 6.13 and 6.15 leads to a robust performance, owing to the balance between true negative and false positive errors (Sections 6.4 and 6.5).

Remark 6.21. The NI-ML mean-field optimality of Theorem 6.17 holds even if different t' values are used for different nodes. However, the mean-field optimality of the NI-ME method of Theorem 6.20 holds if the same t' parameter is used for all nodes. Interestingly, even if the parameter used in the NI-ME algorithm is different than the true observation time parameter (i.e., $t' \neq t$), the optimality argument of Theorem 6.20 holds which indicates the robustness of the method with respect to this parameter. Moreover, the NI-ME optimality of inequality (6.21) is strict, while the one of NI-ML method according to the inequality (6.19) may have multiple optimal solutions.

6.2.3 Multi-source NI

In this section, we consider the multi-source inference problem, where there exists m sources in the network. We consider this problem when sources are sufficiently distant from each other and only a single snapshot, at time t , is available (i.e., $\mathbf{y}(t)$ is given). For simplicity of the analysis, we consider $k = 1$ in the path-based network diffusion kernel (i.e., only shortest paths are considered).

Let $G = (V, E)$ be the underlying network where $d_{i,j}$ represents the length of the shortest path between node i and node j . Define $D(i, R) \triangleq \{j \in V | d_{i,j} < R\}$ as a disk with radius R centered at node i , which we refer to as the R -neighborhood of node i in the network. Similarly, the union of disks with radius R centered at nodes of the set $V_1 \subset V$ is defined as $D(V_1, R) \triangleq \{j \in V | \exists i \in V_1, d_{i,j} < R\}$. We define the following distances in the network:

$$\begin{aligned} d_0 &\triangleq \arg \max_d F(d, t) > \frac{1}{2}, \\ d_1^\epsilon &\triangleq \arg \min_d F(d, t) < \frac{\epsilon}{nm}, \end{aligned} \tag{6.22}$$

where $F(d, t)$ is defined according to (6.14).

Definition 6.22 (ϵ -Coherent Sources). Let $G = (V, E)$ be a binary network. Sources

$\mathcal{S} = \{s_1, s_2, \dots, s_m\}$ are ϵ -coherent if,

$$d(s_a, s_b) > 2(d_0 + d_1^\epsilon), \quad \forall 1 \leq a, b \leq m, a \neq b, \quad (6.23)$$

where d_0 and d_1 are defined according to (6.22).

Intuitively, sources are incoherent if they are sufficiently distant from each other in the network so that their infection effects at time t do not overlap in the network (for instance, viruses released from them, with high probability, have not visited the same nodes.). This assumption is a critical condition to solve the multi-source NI problem efficiently.

Definition 6.23 (Multi-Source Network Diffusion Kernel). Suppose $G = (V, E)$ is a possibly directed binary graph and there exist m source nodes $\mathcal{S} = \{s_1, \dots, s_m\}$ in the network that are ϵ -coherent. We say a node $j \notin \mathcal{S}$ gets infected at time t if it gets a virus from at least one of the sources. Thus, we have,

$$Pr[y_j(t) = 1] \triangleq 1 - \prod_{s \in \mathcal{S}} \bar{p}_{s,j}(t), \quad (6.24)$$

where $\bar{p}_{s,j}(t) = 1 - p_{s,j}(t)$.

Using multi-source network diffusion kernel of Definition 6.23, the log-likelihood function $\mathcal{L}(\mathcal{S}, t)$ and the Hamming error function $\mathcal{H}_\alpha(\mathcal{S}, t)$ are defined as follows:

$$\begin{aligned} \mathcal{L}(\mathcal{S}, t) &\triangleq \sum_{j \in V^t} \log(1 - \prod_{s \in \mathcal{S}} \bar{p}_{s,j}(t)) + \sum_{j \notin V^t} \log(\prod_{s \in \mathcal{S}} \bar{p}_{s,j}(t)), \\ \mathcal{H}_\alpha(\mathcal{S}, t) &\triangleq (1 - \alpha) \sum_{j \in V^t} \prod_{s \in \mathcal{S}} \bar{p}_{s,j}(t) + \alpha \sum_{j \notin V^t} (1 - \prod_{s \in \mathcal{S}} \bar{p}_{s,j}(t)). \end{aligned} \quad (6.25)$$

Similarly to Algorithms 6.12 and 6.14, NI aims to find a set of m sources which maximizes the log-likelihood score, or minimizes the weighted Hamming error. However, unlike the single source case, these optimizations are computationally costly because

all $\binom{|V^t|}{m}$ possible source combinations should be evaluated. If the number of infected nodes is significant ($|V^t| = \mathcal{O}(n)$), even for small constant number of sources, one needs to compute the likelihood or error scores for approximately $\mathcal{O}(n^m)$ possible source subsets, which may be computationally overly challenging for large networks. One way to solve this combinatorial optimization is to take an iterative approach, where, at each step, one source node is inferred. However, at each step, using single source NI methods may not lead to an appropriate approximation because single source NI methods aim to find the source node which explains the entire infection pattern in the network, while in the multi-source case, the entire infection pattern are caused by multiple sources. To avoid this problem, at each step, we use a localized version of NI methods developed in Algorithms 6.12 and 6.14, where sources explain the infection pattern only around their neighborhood in the network.

Definition 6.24. The localized likelihood function of node i in its d_0 neighborhood is defined as,

$$\mathcal{L}_{d_0}(i, t) \triangleq \sum_{\substack{j \in V^t \\ j \in D(i, d_0)}} \log(p_{i,j}(t)) + \sum_{\substack{j \notin V^t \\ j \in D(i, d_0)}} \log(1 - p_{i,j}(t)), \quad (6.26)$$

where only nodes in the d_0 neighborhood of node i is considered in likelihood computation.

A similar argument can be expressed for the localized Hamming prediction error. For large d_0 values, the localized likelihood score is similar to the global likelihood score of (6.16). Using localized likelihood function is important in the multi-source NI problem because source candidates cannot explain the infection pattern caused by other sources. In the following, we propose an efficient localized NI method to solve the multi-source inference problem by maximizing localized likelihood scores of source candidates using a greedy approach. A similar algorithm can be designed for the localized minimum error NI.

Algorithm 6.25 (Multi-source NI-ML Algorithm). Suppose \mathcal{S}_r is the set of inferred sources at iteration r . The localized NI-ML algorithm has the following steps:

- *Step 0:* $\mathcal{S}_0 = \emptyset$.

- *Step $r+1$:*

- *Likelihood computation:* compute s_{r+1} using the following optimization,

$$s_{r+1} = \arg \max_{i \in V^t - D(\mathcal{S}_k, d_1^\epsilon)} \mathcal{L}_{d_0}(i, t). \quad (6.27)$$

- *Update the source set:* add s_{r+1} to the list of inferred sources,

$$\mathcal{S}_{r+1} = \mathcal{S}_r \cup s_{r+1},$$

- *Termination:* stop if $r = m$.

In the following, we show that if sources are sufficiently incoherent (i.e., sufficiently distant from each other in the network), the solution of localized NI Algorithm 6.25 approximates the exact solution closely.

Theorem 6.26. Let $G = (V, E)$ be a regular undirected tree with countably infinite nodes. Suppose sources are ϵ -coherent according to Definition 6.22, and the underlying diffusion process is according to the SI model of Definition 6.1. Suppose \mathcal{S}_r is the set of sources inferred by localized NI Algorithm 6.25 till iteration r . If $\mathcal{S}_k \subset \mathcal{S}$, then with probability at least $1 - \epsilon$, there exists a source node that has not been inferred yet whose localized likelihood score is optimal in expectation:

$$\begin{aligned} \exists s \in \mathcal{S} - \mathcal{S}_r, \quad \mathbb{E}[\mathcal{L}_{d_0}(s, t)] &\geq \mathbb{E}[\mathcal{L}_{d_0}(i, t)] \\ \forall i \in V^t - D(\mathcal{S}_k, d_1^\epsilon). \end{aligned} \quad (6.28)$$

Proposition 6.27. A worst case computational complexity of localized NI Algorithm 6.25 is $\mathcal{O}(|V^t|(k|E| + kn \log(n) + mn))$.

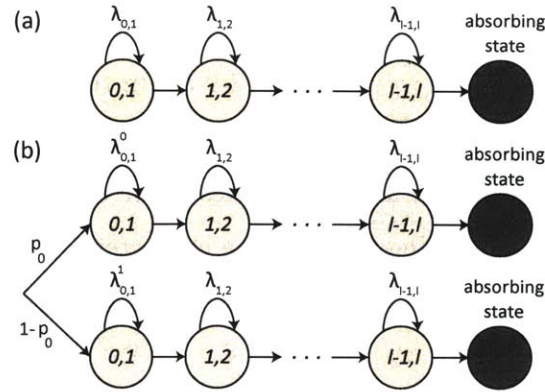


Figure 6-6: (a) A Markov chain of a hypo-exponential distribution. (b) A Markov chain of a mixed hypo-exponential distribution.

6.2.4 NI for Heterogeneous Network Diffusion

In previous sections, we have assumed that the infection spread in the network is homogeneous; i.e., virus traveling time variables $\mathcal{T}_{(i,j)}$ are i.i.d. for all edges in the network. This can be an appropriate model for the binary (unweighted) graphs. However, if edges have weights, the infection spread in the network may be heterogeneous; i.e., the infection spread is faster over strong connections compared to the one of weak edges. Suppose $G = (V, E, W)$ represents a weighted graph, where $w(i, j) > 0$ if $(i, j) \in E$, and $w(i, j) = 0$ otherwise. One way to model a heterogeneous diffusion in the network is to assume that edge holding time variables $\mathcal{T}_{(i,j)}$ are distributed independently according to an exponential distribution with parameter $\lambda_{i,j} = w(i, j)$. According to this model, the average holding time of edge (i, j) is $1/w_{i,j}$, indicating the fast spread of infection over strong connections in the network.

Recall that $\mathcal{T}_{\mathcal{P}_{i \rightarrow j}^r}$ represents the virus traveling time variable from node i to node j over the path $\mathcal{P}_{i \rightarrow j}^r$. To simplify notations and highlight the main idea, consider the path $\mathcal{P}_{0 \rightarrow l}^r = \{0 \rightarrow 1 \rightarrow 2 \dots \rightarrow l\}$. The virus traveling time from node 0 to node l over this path ($\mathcal{T}_{\mathcal{P}_{0 \rightarrow l}^r}$) is a *hypoexponential* variable whose distribution is a special case of the *phase-type distribution*. For this path, we consider a Markov chain with $l + 1$ states, where the first l states are transient, and the state $l + 1$ is an absorbing state. Each transient state of this Markov chain corresponds to

an edge (i, j) over this path whose holding time is characterized by an exponential distribution with rate $\lambda_{i,j} = w(i, j)$ (Figure 12-2-a). In this setup, the virus traveling time from node 0 to node l over the path $\mathcal{P}_{0 \rightarrow l}^r$ is equal to the time from the start of the process until reaching to the absorbing state of the corresponding Markov chain. The distribution of this absorbing time can be characterized as a special case of the phase-type distribution. A subgenerator matrix of the Markov chain of Figure 12-2-a is defined as follows:

$$\begin{bmatrix} -\lambda_{0,1} & \lambda_{0,1} & 0 & \dots & 0 & 0 \\ 0 & -\lambda_{1,2} & \lambda_{1,2} & \ddots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \ddots & -\lambda_{l-3,l-2} & \lambda_{l-3,l-2} & 0 \\ 0 & 0 & \dots & 0 & -\lambda_{l-2,l-1} & \lambda_{l-2,l-1} \\ 0 & 0 & \dots & 0 & 0 & -\lambda_{l-1,l} \end{bmatrix}. \quad (6.29)$$

For simplicity, denote the above matrix by $\Theta \equiv \Theta(\lambda_{0,1}, \dots, \lambda_{l-1,l})$. Define $\alpha = (1, 0, \dots, 0)$ as the probability of starting in each of the l states. Then, the Markov chain absorption time is distributed according to $PH(\alpha, \Theta)$, where $PH(\alpha, \Theta)$ represents a phase-type distribution. In this special case, this distribution is also called a hypoexponential distribution. A similar subgenerator matrix Θ can be defined for a general path $\mathcal{P}_{i \rightarrow j}^r$ connecting nodes i to j . Thus, we have

$$F_{\mathcal{P}_{i \rightarrow j}^r}(t) = Pr[\mathcal{T}_{\mathcal{P}_{i \rightarrow j}^r} \leq t] = 1 - \alpha e^{t\Theta} \mathbf{1}, \quad (6.30)$$

where $\mathbf{1}$ is a column vector of ones of the size $|\mathcal{P}_{i \rightarrow j}^r|$, and e^X is the matrix exponential of X . For an unweighted graph where all edges have the same rate λ , (6.32) is simplified to (6.13). For the weighted graph $G = (V, E, W)$, we compute k shortest paths among pairs of nodes over the graph $G' = (V, E, W')$, where $w'(i, j) = 1/w(i, j)$ if $(i, j) \in E$, otherwise $w'(i, j) = \infty$. Then, the path network diffusion kernel for a

weighted graph $G = (V, E, W)$ can be defined according to Definition (6.7). Using this kernel, NI algorithms introduced in Sections 6.2.2 and 6.2.3 can then be used to infer the source node under the heterogeneous diffusion in the network. Note that this framework can be extended to a more complex diffusion setup as well. We provide an example of such diffusion setup in the following:

Example 6.28. Consider the the path $\mathcal{P}_{0 \rightarrow l}^r = \{0 \rightarrow 1 \rightarrow 2 \rightarrow \dots \rightarrow l\}$. Suppose the edge $(0, 1)$ spreads the infection with rates $\lambda_{0,1}^0$ and $\lambda_{0,1}^1$ with probabilities p_0 and $1 - p_0$, respectively. Suppose other edges (i, j) of this path spread the infection with rate $\lambda_{i,j}$. Figure 12-2-b illustrates the corresponding Markov chain for this path. The subgenerator matrices of this Markov chain can be characterized as follows:

$$\Theta_i = \begin{bmatrix} -\lambda_{0,1}^i & \lambda_{0,1}^i & 0 & \dots & 0 & 0 \\ 0 & -\lambda_{1,2} & \lambda_{1,2} & \ddots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \ddots & -\lambda_{l-3,l-2} & \lambda_{l-3,l-2} & 0 \\ 0 & 0 & \dots & 0 & -\lambda_{l-2,l-1} & \lambda_{l-2,l-1} \\ 0 & 0 & \dots & 0 & 0 & -\lambda_{l-1,l} \end{bmatrix}, \quad (6.31)$$

for $i = 0, 1$. Then, for this path, we have

$$F_{\mathcal{P}_{0 \rightarrow l}^r}(t) = Pr[\mathcal{T}_{\mathcal{P}_{0 \rightarrow l}^r} \leq t] = 1 - \alpha(p_0 e^{t\Theta_0} + (1 - p_0)e^{t\Theta_1})\mathbf{1}. \quad (6.32)$$

To compute the path-based network diffusion kernel, we compute shortest paths among pairs of nodes over the graph $G' = (V, E, W')$, where $w'(0, 1) = p_0/\lambda_{0,1}^0 + (1 - p_0)/\lambda_{0,1}^1$, $w'(i, j) = 1/w(i, j)$ if $(i, j) \in E$ and $(i, j) \neq (0, 1)$, and $w'(i, j) = \infty$ otherwise. This example illustrates that NI framework can be used even under a complex heterogeneous diffusion setup.

6.2.5 Non-parametric NI - Unknown Snapshot Time

In some real-world applications, only the network structure G and infection patterns $\{\mathbf{y}(t_1), \dots, \mathbf{y}(t_z)\}$ are known and therefore to use NI algorithms, we need to learn the parameters such as observation times $\{t_1, \dots, t_z\}$, the number of sources m in the network, and the number of disjoint shortest paths considered in the diffusion kernel k . In the following, we introduce efficient techniques to learn these parameters:

Observation time parameters: In the maximum likelihood NI Algorithm 6.12, according to Remark 6.19, the true parameter t is the one that maximizes the expected source likelihood score according to Theorem 6.17. Thus, in the case of unknown parameter t , we solve the following optimization:

$$(s, t) = \arg \max_{i, t'} \mathcal{L}(i, t'), \quad (6.33)$$

where $\mathcal{L}(i, t')$ is the log-likelihood function of node i using parameter t' . One way to solve Optimization (6.33) approximately is to quantize the range of parameter t (i.e., $t \in (0, t_{max})$) to b bins and evaluate the objective function in each case. Because we assume that the λ parameter of the edge holding time distribution is equal to one, one appropriate choice for t_{max} is the diameter of the infected subgraph, defined as the longest shortest path among pairs of infected nodes. The number of quantization levels b determines the resolution of the inferred parameter t and therefore the tightness of the approximation. If t_{max} is large and the true t parameter is small, to have a tight approximation, the number of quantization levels b should be large which may be computationally costly. In this case, one approach to estimate parameter t is to use the first moment approximation of the Erlang network diffusion kernel over source neighbors. Suppose μ is the fraction of the infected neighbors of source s . Since infection probabilities of source neighbors approximately come from an exponential distribution, for a given parameter t , $\mu \approx 1 - e^{-t}$. Therefore,

$$t \approx -\ln(1 - \mu). \quad (6.34)$$

In the minimum error NI Algorithm 6.14, according to Remark 6.21, the prediction error of all infected nodes should be computed using the same parameter t . In the setup of Theorem 6.20, any value of parameter t leads to an optimal solution in expectation. In general, we suggest the following approach to choose this parameter: First, for each node, we minimize the prediction error for different values of the parameter t as follows,

$$t_i^* = \arg \max_{t'} \mathcal{H}_\alpha(i, t'). \quad (6.35)$$

This Optimization can be solved approximately similarly to the case of maximum likelihood NI Optimization (6.33). For the small t values, we use (6.34) to obtain t_i^* . Then, to obtain a fixed t parameter for all nodes, we take the median of t_i^* parameters of nodes with the minimum prediction error (in our experiments, we consider top 10 predictions with the smallest error to estimate the parameter t). In the cases of multi-source and multi-snapshot NI, one can use similar approaches to estimate time stamp parameters.

The number of sources: In NI algorithms presented in Section 6.2.3, we assume that the number of sources in the network (i.e., the parameter m) is known. In the case of unknown parameter m , if sources are sufficiently incoherent according to Definition 6.22, one can estimate m as follows: because sources are incoherent, their caused infected nodes do not overlap with each other with high probability. Thus, the number of connected components of the infected sub-graph (or the number of infected clusters in the network) can provide a reliable estimate of the number of sources in this case.

Regularization parameter of NI-ME: The minimum error NI Algorithm 6.14 has a regularization parameter α which balances between false positive and false negative error types. In the setup of Theorem 6.20, any value of $0 < \alpha < 1$ leads to an optimal expected weighted Hamming error solution of (6.21). However, in general, we choose this parameter according to Remarks 6.13 and 6.15 to have a balance between the

number of false negative and false positive errors.

Number of disjoint shortest paths considered in kernel computation: Our proposed path-based network diffusion kernel considers up to k independent shortest paths among nodes where a node j gets infected at time t if the infection reaches to it over at least one of the k disjoint shortest paths connecting that node to the source node. A path-based network diffusion kernel with $k = 1$ has the least computational complexity among other path-based network diffusion kernels. Considering more paths among nodes in the network (i.e., $k > 1$) provides a better characterization of network diffusion processes with the cost of increased kernel computational complexity (Proposition 6.8). For example, over an asymmetric grid network of Section 6.4, considering increased numbers of paths among nodes to form the path network diffusion kernel improves the performance of NI methods significantly as higher order paths partially capture the asymmetric diffusion spread in the network.

6.3 NI with Multiple Snapshots

In this section, we consider the NI problem when multiple snapshots from infection patterns are available. To simplify notation and highlight the main ideas, we consider the single source case with two samples $\mathbf{y}(t_1)$ and $\mathbf{y}(t_2)$ at times t_1 and t_2 , respectively. All arguments can be extended to a more general setup as well.

Recall that V^t denotes the set of infected nodes at time t . Let $E^t = \{(i, j) | (i, j) \in E, \{i, j\} \subset V^t\}$ represent edges among infected nodes in the network. The infection subgraph $G^t = (V^t, E^t)$ is connected if there is no infection recovery and the underlying diffusion is according to a dynamic process. We define an infection contraction operator $g(\cdot)$ as follows:

$$g(v) = \begin{cases} v & v \in V \setminus V^t \\ x & o.w. \end{cases} \quad (6.36)$$

where $x \notin V$. In other words, $g(\cdot)$ maps all infected nodes to a new node x , while it maps all other nodes to themselves.

Definition 6.29 (Infusion Contraction Graph). Suppose $G = (V, E, W)$ is a weighted graph whose infected subgraph at time t is represented as $G^t = (V^t, E^t, W^t)$. An infusion contraction graph $G_c^t = (V_c^t, E_c^t, W_c^t)$ is defined as follows:

- $(i, j) \in E_c^t$ for $i, j \neq x$ iff $(i, j) \in E$. In this case, $w_c^t(i, j) = w(i, j)$.
- $(i, x) \in E_c^t$ for $i \neq x$ iff there exists $j \in V^t$ such that $(i, j) \in E$. In this case, $w_c^t(i, x) = \sum_{\substack{j \in V^t \\ (i, j) \in E}} w(i, j)$.

Intuitively, the infusion contraction graph considers the infected subgraph as one node and adjusts weights of un-infected nodes connected to the infected ones accordingly. Now we consider the source inference problem when two snapshots at times t_1 and t_2 are given. Recall that V^{t_1} and V^{t_2} denote the set of infected nodes at times t_1 and t_2 , respectively. Without loss of generality, we assume $0 < t_1 < t_2$. Using the probability chain-rule, we can re-write the likelihood scores of Optimization (6.4) as,

$$\begin{aligned} & Pr(\mathbf{y}(t_1), \mathbf{y}(t_2) | \mathcal{S} = \{s\}) \\ &= \underbrace{Pr(\mathbf{y}(t_1) | \mathcal{S} = \{s\})}_{\text{term I}} \underbrace{Pr(\mathbf{y}(t_2) | \mathbf{y}(t_1), \mathcal{S} = \{s\})}_{\text{term II}}. \end{aligned} \tag{6.37}$$

Term (I) is the likelihood score of the single source NI Optimization (6.16). We consider different possibilities for Term II as follows:

- If $y_j(t_1) = 1$ and $y_j(t_2) = 1$, $Pr(y_j(t_2) | \mathbf{y}(t_1), \mathcal{S} = \{s\}) = 1$, because if a node gets infected, it remains infected (there is no recovery). Thus, if $y_j(t_1) = 1$ and $y_j(t_2) = 0$, $Pr(y_j(t_2) | \mathbf{y}(t_1), \mathcal{S} = \{s\}) = 0$.
- Now we consider the case $y_j(t_1) = 0$. Let $G_c^{t_1}$ be the infusion contraction graph of Definition 6.29. Suppose that all infected nodes at time t_1 are mapped to the node x . The second term can be approximated as follows:

$$Pr(y_j(t_2) = 1 | \mathbf{y}(t_1), \mathcal{S} = \{s\}) \approx p_{x,j}(t_2 - t_1), \quad (6.38)$$

where $p_{x,j}(\cdot)$ is the path-based network diffusion kernel over the graph G_c^t . In other words, all infected nodes at time t_1 can be viewed as a single source x in the infusion contraction graph. Note that this approximation is tight when the underlying edge holding time distribution is an exponential distribution which has the memory-less property. Moreover, note that to compute the diffusion kernel in this case, we use the infusion contraction graph because infected nodes are not incoherent, and therefore, the multi-source diffusion kernel of Definition 6.23 cannot be used.

Under approximation (6.38), the second term of (6.37) leads to a similar expression for all source candidates, and therefore, the optimization is simplified to a single snapshot one. In practice, one may compute average source likelihood scores using all snapshots to decrease the variance of the source likelihood scores.

6.4 NI Over Synthetic Networks

In this section, we compare the performance of proposed NI algorithms with other source inference methods over four different synthetic network structures. In our simulations, we assume that there exists a single source in the network, and the underlying diffusion is according to the SI model of Definition 6.1. In the SI model, edge holding time variables $\mathcal{T}_{(i,j)}$ are i.i.d. having an exponential distribution with parameter $\lambda = 1$. Note that to generate simulated diffusion patterns, we do not use our path-based network diffusion kernel to have a fair performance comparison of our methods with the one of other source inference techniques. We use the following methods in our performance assessment:

- **Distance Centrality:** This method infers the source node with the minimum shortest path distance from all infected nodes. Suppose G is the underlying

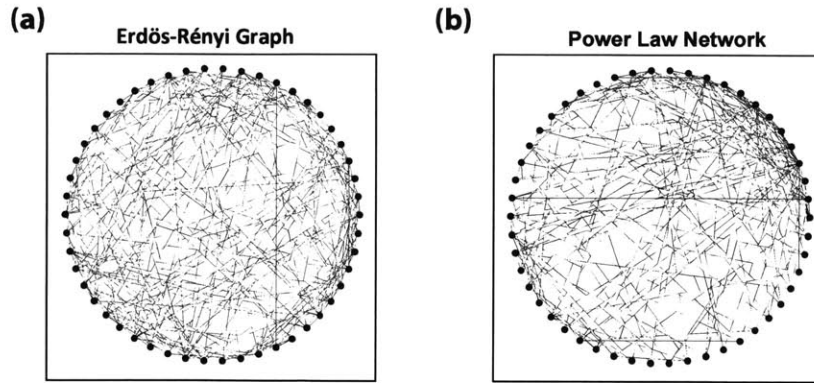


Figure 6-7: Examples of (a) Erdős-Rényi, and (b) power law networks, with 50 nodes.

network. The distance centrality of node i corresponding to infected nodes at time t is defined as follows:

$$D_t(i, G) \triangleq \sum_{j \in V^t} d(i, j), \quad (6.39)$$

where $d(i, j)$ represents the length of the shortest path between nodes i and j .

A source node is inferred using the following optimization:

$$s = \arg \min_{i \in V^t} D_t(i, G). \quad (6.40)$$

If there is no path between two nodes i and j , $d(i, j) = \infty$. This makes the distance centrality measure sensitive to noise specially in real world applications. To avoid this issue, for disconnected nodes i and j , we assign $d(i, j) = M$, where M is a large number compared to shortest path distances in the network. In our simulations, we set M as 5 times larger than the network diagonal (i.e., the longest shortest path in the network).

- **Degree Centrality:** This methods infers the source node with highest direct connections to other infected nodes. Degree centrality of node i corresponding to infected nodes V^t is defined as follows:

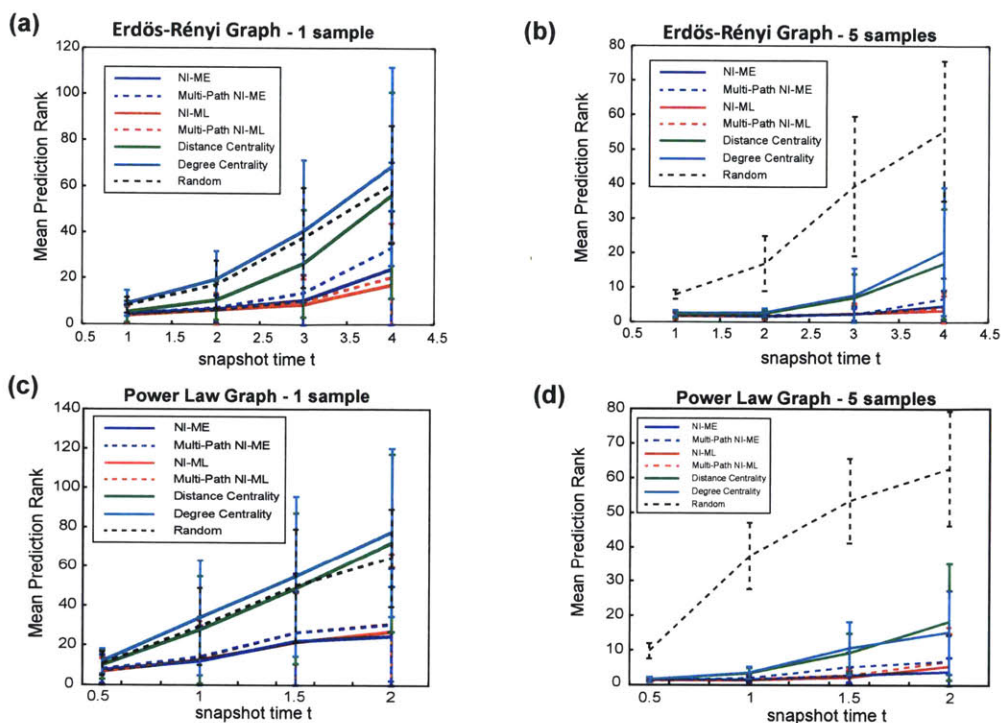


Figure 6-8: Performance of source inference methods over (a,b) Erdős-Rényi, and (c,d) power law networks graphs with 250 nodes and an average density of 1%. Error bars represent \pm standard deviations of empirical results for each case. Experiments have been repeated 100 times so that error margins are small.

$$C_t(i, G) \triangleq \sum_{j \in V^t} G(i, j). \quad (6.41)$$

Note that unlike the distance centrality method which considers both direct and indirect interactions among infected nodes, degree centrality only considers direct interactions. To infer the source node using the degree centrality approach, one needs to solve the following optimization:

$$s = \arg \min_{i \in V^t} C_t(i, G). \quad (6.42)$$

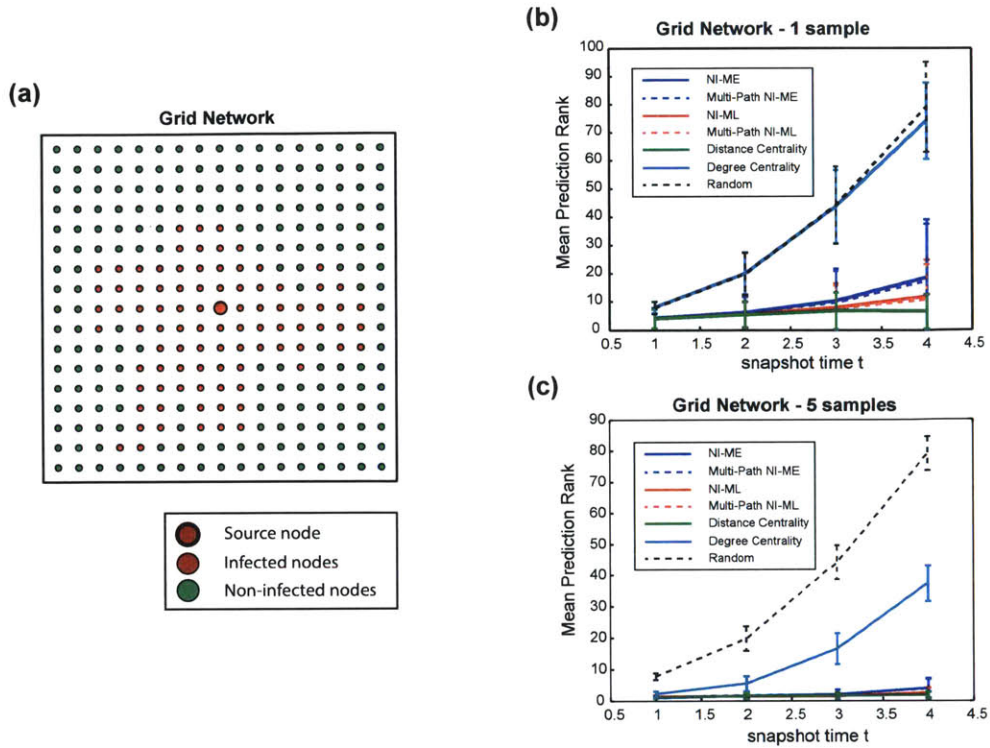


Figure 6-9: (a) An example of a grid network. Performance of source inference methods over grid networks with 250 nodes with (a) one and (b) five independent samples. Error bars represent \pm standard deviations of empirical results for each case. Experiments have been repeated 100 times so that error margins are small.

- **Network Infusion:** We use NI methods based on maximum likelihood (denoted as NI-ML) described in Algorithm 6.12, and minimum error (denoted as NI-ME) described in Algorithm 6.14. To have a fair comparison with other methods, we assume that the observation time parameter t is unknown and is estimated using the techniques presented in Section 6.2.5.

In our simulations, we use four types of input networks:

- **Erdős-Rényi graphs:** In this case, G is a symmetric random graph where, $Pr[G(i, j) = 1] = p$. Networks have 250 nodes. An example network with fewer number of nodes is shown in Figure 6-7-a). In our simulations, we use $p = 0.01$.
- **Power law graphs:** We construct G as follows [221]; we start with a random subgraph with 5 nodes. At each iteration, a node is added to the network

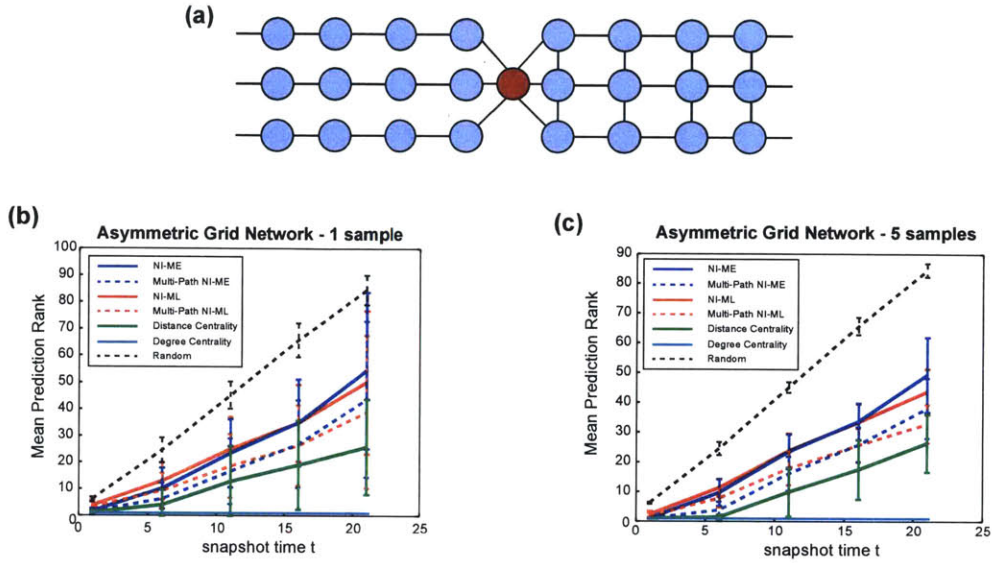


Figure 6-10: (a) An asymmetric grid network. Performance of source inference methods over asymmetric grid networks with 250 nodes with (a) one and (b) five independent samples. Node with red color is the source node. Error bars represent \pm standard deviations of empirical results for each case. Experiments have been repeated 100 times so that error margins are small.

connecting to θ existent nodes with probabilities proportional to their degrees. This process is repeated till the number of nodes in the network is equal to $n = 250$. In our simulations, we use $\theta = 2$ which results in networks with the average density approximately 0.01 (see an example with fewer number of nodes in Figure 6-7-b).

- **Grid networks:** In this case, G is an undirected square grid network with 250 nodes. Figure 6-9-a shows an example of such networks with fewer number of nodes. We assume that the source node is located at the center of the grid to avoid boundary effects.
- **Asymmetric grid networks:** In this case, G is an undirected graph with 250 nodes. It has 6 branches connected to the central node, three branches on the right with heavier connectivity among their nodes, and three branches on the left with sparse connectivity. Figure 6-10-a shows an example of such networks with fewer number of nodes. This is an adversarial example to highlight cases

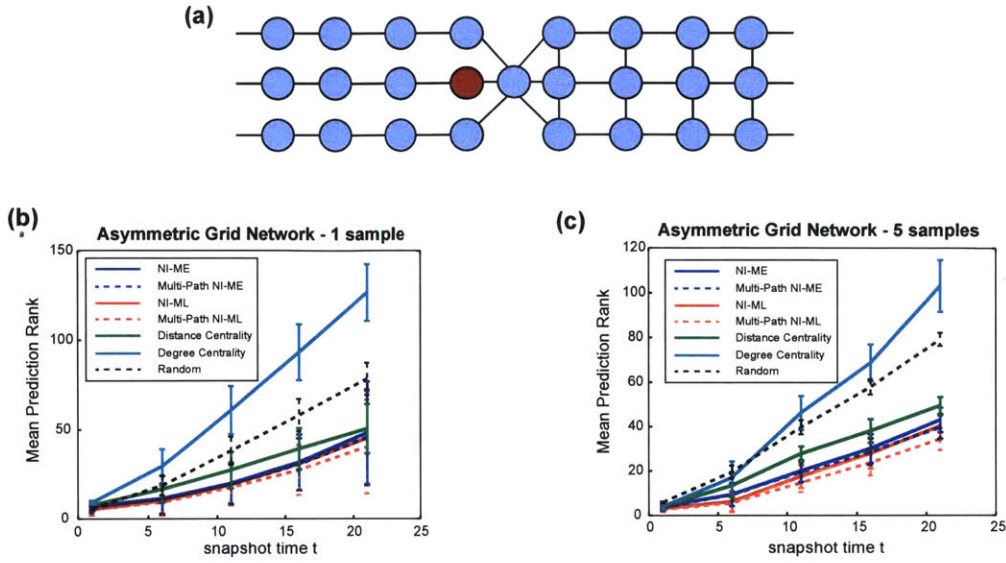


Figure 6-11: (a) An asymmetric grid network. Performance of source inference methods over asymmetric grid networks with 250 nodes with (a) one and (b) five independent samples. Node with red color is the source node. Error bars represent \pm standard deviations of empirical results for each case. Experiments have been repeated 100 times so that error margins are small.

where NI methods fail to converge to the optimal value.

In multi-path NI methods, we consider top 10 independent shortest paths to form the k -path network diffusion kernel (i.e., $k = 10$). However, in the grid network, to enhance the computation time, we consider $k = 2$ because different nodes have at most 2 independent shortest paths connected to the source node. Parameter t is the time at which we observe the infection spread, and it determines the fraction of infected nodes in the network. If t is very large compared to the graph diameter, almost all nodes in the network become infected. On the other hand, for very small values of t , the source inference problem becomes trivial. In our simulations, we consider the cases where the number of infected nodes in the network is less than 75% of the total number of nodes, and greater than at least 10 nodes. We learn this parameter using the observed infection pattern according to techniques introduced in Section 6.2.5.

For evaluation, we sort infected nodes as source candidates according to the score

obtained by different methods. High performing methods should assign the highest scores to the source node. The source node should appear on the top of the inferred source candidates. Ideally, if a method assigns the highest score to the source node, the rank of the prediction is one. We use the rank of true sources averaged over different runs of simulations. More formally, suppose $r(\mathcal{M}, s)$ is the rank of the source node s inferred by using the method \mathcal{M} . In an exact prediction, $r(\mathcal{M}, s) = 1$, while an average rank of a method based on random guessing is $r(\mathcal{M}, s) = |V^t|/2$. If $r(\mathcal{M}, s)$ is close to one, it means that the true source is among top predictions of the method \mathcal{M} . In each case, we run simulations 100 times.

Figure 6-8 compares the performance of different source inference methods over both Erdős-Rényi and power law networks, and in different ranges of the parameter t . In both network models and in all diffusion rates, NI Algorithms based on maximum likelihood (NI-ML) and minimum error (NI-ME) outperform other methods. Panels (a) and (c) illustrate the performance when only one sample from the infection pattern at time t is available, for Erdős-Rényi and power law networks, respectively. Panels (b) and (d) illustrate the performance of different methods when five independent samples from the infection pattern at time t are given, illustrating the mean-field optimality of NI methods according to Theorems 6.17 and 6.20. Because the underlying networks are sparse, according to Proposition 6.11, the performance of maximum likelihood and minimum error NI methods, both shortest-path and multi-path versions, are close to each other in both network models. Notably, unlike NI methods, the performance of other source inference methods such as distance centrality and degree centrality does not tend to converge to the optimal value even for higher sample sizes.

Figure 6-9 compares the performance of different source inference methods over grid networks in different ranges of the parameter t . In the case of having a single sample from the infection pattern, in small ranges of the parameter t , when the fraction of infected nodes is less than approximately $1/4n$, NI methods and distance centrality have approximately the same performance, significantly outperforming the degree centrality method. In higher diffusion rates, distance centrality outperforms NI-ML method and NI-ML method outperforms NI-ME method. Again in this range, the

performance of the degree centrality method is significantly worse than other methods. Having five independent samples from node infection patterns, the performance of NI methods and distance centrality converges to the optimal one.

Unlike tree and sparse Erdős-Rényi networks, there are multiple overlapping paths among nodes in the grid structures. However, as we illustrate in Figure 6-9, the performance of NI methods converge to the optimal value, similarly to the case of sparse graphs. The main reason is that grid structures are symmetric and even though paths among nodes overlap significantly, considering shortest paths among nodes approximates the true underlying diffusion based on an SI dynamics closely. In order to have an adversarial example where the performance of NI methods do not converge to the optimal one, we design an asymmetric grid structure illustrated in Figure 6-10. The three branches on the right side of the central node have strong connectivity among themselves, forming a grid, while the ones on the left is sparsely connected to each other, forming a tree. Therefore, the spread of infection will be faster on the right side compared to the left side, under a dynamic SI model. However, considering only the shortest path among nodes does not capture this asymmetric structure. Therefore, the performance of shortest path NI methods diverges from the optimal value as diffusion rate increases (see Figure 6-10). In this case, considering more paths among nodes to form the k -path network diffusion kernel according to Definition 6.7 improves the performance of NI methods significantly, because higher order paths partially capture the asymmetric diffusion spread in the network. Note that the degree centrality method has the best performance in this case, because the source node has the highest degree in the network by design. If we select another node to be the source node as illustrated in Figure 6-11-a, the performance of the degree centrality method becomes worst significantly, indicating its sensitivity to the source location. In the setup of Figure 6-11, multi-path NI methods outperform the one of other methods, although their performance do not converge to the optimal value.

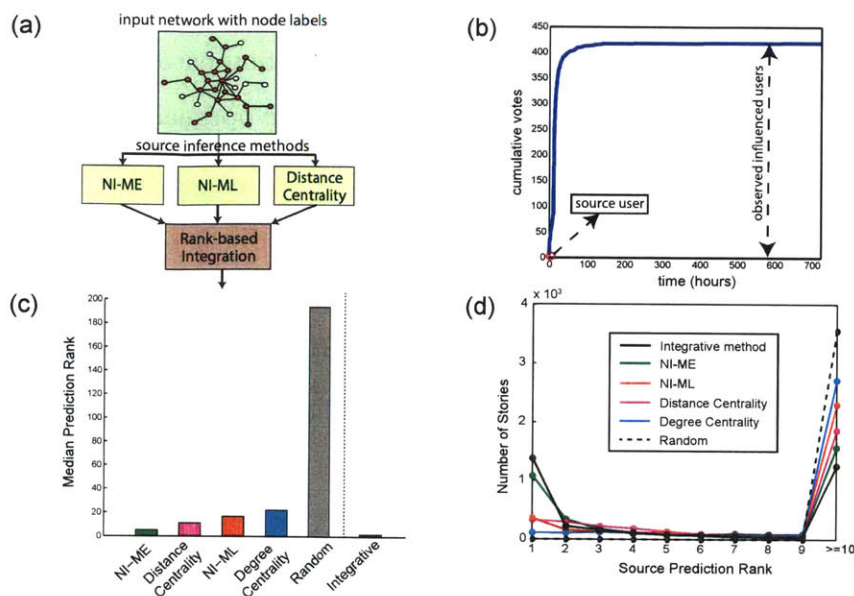


Figure 6-12: (a) An integrative source inference framework which combines prediction ranks of NI-ME, NI-ML and distance centrality methods. (b) The cumulative number of votes for a particular story over time. (c) Median prediction ranks of true sources for over 3,500 Digg stories inferred by individual and integrative source inference methods. NI-ME outperforms other individual methods with a median prediction rank of 5, while the integrative source inference leads to the best performance with a median prediction rank of 3. (d) Histogram of prediction ranks of sources inferred by different source inference methods. NI-ME and integrative source inference algorithms infer news sources optimally (in their top prediction) for approximately 30% and 38% of stories, respectively.

6.5 Inference of News Sources in a Social News Network

In this section, we evaluate the performance of NI and other source inference techniques in identifying news sources over the social news aggregator Digg (<http://digg.com>). We also propose an integrative source inference framework which combines network infusion methods (NI-ME and NI-ML) with the distance centrality method. In our integrative framework, we compute the average rank of each node to be the source node according to different source inference methods (Figure 6-12-a). We exclude the degree centrality method from our integration owing to its low performance over

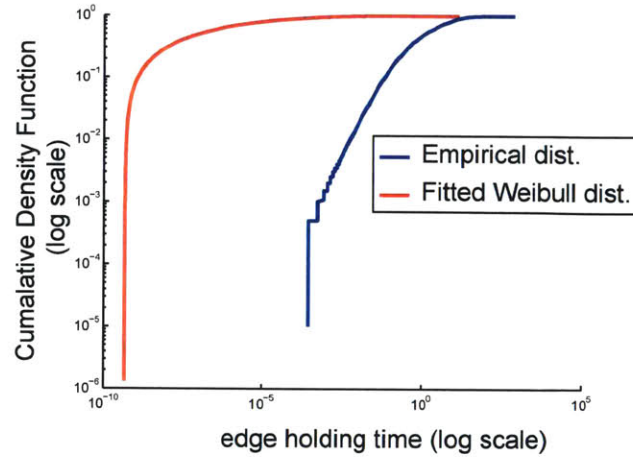


Figure 6-13: We consider news spread over the Digg social news networks for over 3,500 news stories. In approximately 65% of cases, nodes that have received the news at a time t do not have any neighbors who had already received the news by that time, violating the basic conditional independency assumption of the SI model. Furthermore, the empirical distribution of remaining news propagation times over network edges cannot be approximated closely by a single distribution of a homogenous SI diffusion model, even by fitting a general Weibull distribution to the observed data (Section 6.5).

synthetic and real networks.

Digg allows its users to submit and rate different news links. Highly rated news links are *promoted* to the front page of Digg. Digg also allows its users to become friends with other users and follow their activities over the Digg network. Digg's friendship network is asymmetric, i.e., user A can be a follower (friend) of user B but not vice versa. Reference [256] have collected voting activities and friendship connections of Digg's users over a period of a month in 2009 for 3,553 promoted stories. We use this data to form a friendship network of active Digg users with 24,219 nodes and more than 350K connections. We consider users as active if they have voted for at least 10 stories in this time period. Figure 6-1 demonstrates a small part of the Digg friendship network.

For each story, we have voting time stamps of different users where the first voting user is the source of that story. Figure 6-12-b demonstrates the number of votes of a particular story at different times. If friends of a user A vote for a specific story, it

is more likely that user A will also vote for that story and that is how information propagates over the Digg friendship network. NI aims to inverse this information propagation process to infer the news source by observing the voting pattern in a single snapshot in the steady state when no further news propagation occurs. Here, we only consider the shortest path among nodes to compute the path-based network diffusion kernel used in NI Algorithms 6.12 and 6.14.

This application provides an interesting real data framework to assess the performance and robustness of different source inference methods because the true source (i.e., the user who started the news) are known for different stories and also the underlying diffusion processes are based on real dynamics over the Digg friendship network. Moreover, not all of the voting pattern is derived by the source users and there are disconnected voting activities over the Digg friendship network. Thus, performance assessment of different source inference methods in this application can provide a measure of robustness of different methods under real-world circumstances.

Figure 6-12-c demonstrates the median prediction ranks of true sources inferred by individual and integrative source inference methods. NI-ME outperforms other individual source inference methods with a median prediction rank of 5, while the integrative source inference leads to the best performance with a median prediction rank of 3. In fact, NI-ME and integrative source inference algorithms infer news sources optimally (in their top prediction) for approximately 30% and 38% of stories, respectively (Figure 6-12-d). This illustrates the robustness of these methods in cases where the underlying dynamics are unknown.

Real dynamics over the Digg social news network is in fact significantly different from the standard SI diffusion model. Suppose for a given news propagation, a node i is infected at time τ_i . We compute a realization of random variable $\mathcal{T}_{(j,i)}$ according to the SI model 6.1. Over all considered news propagations over the Digg network, around 65% of non-source nodes do not have any infected neighbors in the network at the time of being infected, which violates the most basic conditional independency assumption of the SI model. This can be partially owe to existence of latent sources of information propagation, and/or noisy and spurious network connections. Further-

more, the empirical distribution of remaining news propagation times over edges of the Digg social news network cannot be approximated closely by a single distribution of a homogenous SI diffusion model, even by fitting a general Weibull distribution to the observed data. In Figure 6-13, we compute the best fit of Weibull distribution (which contains exponential and Rayleigh distributions) to the observed empirical distribution. Parameters of the fitted distribution are $\lambda \approx 3.13$ and $k \approx 0.51$. As it is illustrated in this figure, fitted and empirical distributions differ significantly from each other.

6.6 Proofs

In this section, we present proofs of the main results of the paper.

6.6.1 Proof of Proposition 6.8

To compute the k -path network diffusion kernel, we need to compute k -independent shortest paths among nodes. Note that ties among paths with the same length is broken randomly as explained in Section 6.2.1. Computation of these paths among one node and all other nodes using the Dijkstra's algorithm costs $\mathcal{O}(k|E| + kn \log(n))$. Thus, computational complexity of forming the entire kernel has complexity $\mathcal{O}(k|E|n + kn^2 \log(n))$.

6.6.2 Proof of Proposition 6.11

First, we compute the probability that a node $v \in V$ belongs to a cycle of length at most l . Such a cycle is determined by the $l - 1$ other vertices. By choosing them in order, there are less than n^{l-1} choices for those other vertices, while the cycle appears with probability p^l in the graph. Thus, the probability that v is involved in a cycle of length l is at most $n^{l-1}p^l \leq c^l/n$. To have an overlapping shortest path between nodes i and j , at least one of the nodes over that path should belong to a cycle of length at most $2l_0$. This happens with probability less than $l_0 c^{2l_0}/n$, which goes to

zero asymptotically if $c < 1/n^{2l_0}$.

6.6.3 Proof of Proposition 6.16

Computation of the k -path network diffusion kernel for one node has complexity $\mathcal{O}(k|E| + kn \log(n))$, according to Proposition 6.8. We need to compute the kernel for V^t nodes. Moreover, Optimizations (6.16) and (6.18) have complexity $\mathcal{O}(|V^t|n)$. Thus, the total computational complexity is $\mathcal{O}(|V^t|(k|E| + kn \log(n)))$.

6.6.4 Proof of Theorem 6.17

First, we prove the following lemma:

Lemma 6.30. Let x, y and z be positive numbers such that $0 < x, y < z$. Define

$$f(x, y) \triangleq x \log \frac{x}{y} + (z - x) \log \frac{z - x}{z - y}. \quad (6.43)$$

Then, $f(x, y) \geq 0$ where equality holds iff $x = y$.

Proof. We have,

$$\begin{aligned} \frac{\partial f}{\partial y} &= -\frac{x}{y} + \frac{z - x}{z - y}, \\ \frac{\partial^2 f}{\partial y^2} &= \frac{x}{y^2} + \frac{z - x}{(z - y)^2} > 0. \end{aligned}$$

For any $0 < x_0 < z$, define $g(y) = f(x_0, y)$. Above we show that $g(y)$ is convex, and its minimum is zero which occurs at $y = x_0$. Since this holds for any $0 < x_0 < z$, this completes the proof. \square

Recall that $\mathcal{L}(i, t')$ is the likelihood score of node i using diffusion parameter t' :

$$\mathcal{L}(i, t') = \sum_{j \in V^t} \log(p_{i,j}(t')) + \sum_{j \notin V^t} \log(1 - p_{i,j}(t')). \quad (6.44)$$

First, we prove Proposition 6.18. Let s be the source node and t be the infection observation time (t' can be different than t , see explanations of Remark 6.19). Thus, we can write,

$$\begin{aligned}
 & \mathbb{E}[\mathcal{L}(s, t)] - \mathbb{E}[\mathcal{L}(s, t')] & (6.45) \\
 &= \sum_{j \in V} p_{s,j}(t) \log(p_{s,j}(t)) + (1 - p_{s,j}(t)) \log(1 - p_{s,j}(t)) \\
 & \quad - \sum_{j \in V} p_{s,j}(t) \log(p_{s,j}(t')) + (1 - p_{s,j}(t)) \log(1 - p_{s,j}(t')) \\
 &= \sum_{j \in V} p_{s,j}(t) \log \frac{p_{s,j}(t)}{p_{s,j}(t')} + (1 - p_{s,j}(t)) \log \frac{1 - p_{s,j}(t)}{1 - p_{s,j}(t')} \\
 &\stackrel{(I)}{\geq} \sum_{j \in V} p_{s,j}(t) \log \frac{\sum_{j \in V} p_{s,j}(t)}{\sum_{j \in V} p_{s,j}(t')} + \sum_{j \in V} 1 - p_{s,j}(t) \log \frac{\sum_{j \in V} 1 - p_{s,j}(t)}{\sum_{j \in V} 1 - p_{s,j}(t')} \\
 &\stackrel{(II)}{\geq} 0.
 \end{aligned}$$

Inequality (I) follows from the log-sum inequality. Inequality (II) follows from Lemma 6.30. In particular, according to Lemma 6.30, the equality condition (II) holds iff $\sum_{j \in V} p_{s,j}(t) = \sum_{j \in V} p_{s,j}(t')$ which indicates $t = t'$. This completes the proof of Proposition 6.18. In the next step, we prove Theorem 6.17.

$$\begin{aligned}
 & \mathbb{E}[\mathcal{L}(s, t)] - \mathbb{E}[\mathcal{L}(i, t')] & (6.46) \\
 &= \sum_{j \in V} p_{s,j}(t) \log(p_{s,j}(t)) + (1 - p_{s,j}(t)) \log(1 - p_{s,j}(t)) \\
 & \quad - \sum_{j \in V} p_{s,j}(t) \log(p_{i,j}(t')) + (1 - p_{s,j}(t)) \log(1 - p_{i,j}(t')) \\
 &= \sum_{j \in V} p_{s,j}(t) \log \frac{p_{s,j}(t)}{p_{i,j}(t')} + (1 - p_{s,j}(t)) \log \frac{1 - p_{s,j}(t)}{1 - p_{i,j}(t')} \\
 &\stackrel{(III)}{\geq} \sum_{j \in V} p_{s,j}(t) \log \frac{\sum_{j \in V} p_{s,j}(t)}{\sum_{j \in V} p_{i,j}(t')} + \sum_{j \in V} 1 - p_{s,j}(t) \log \frac{\sum_{j \in V} 1 - p_{s,j}(t)}{\sum_{j \in V} 1 - p_{i,j}(t')} \\
 &\stackrel{(IV)}{\geq} 0.
 \end{aligned}$$

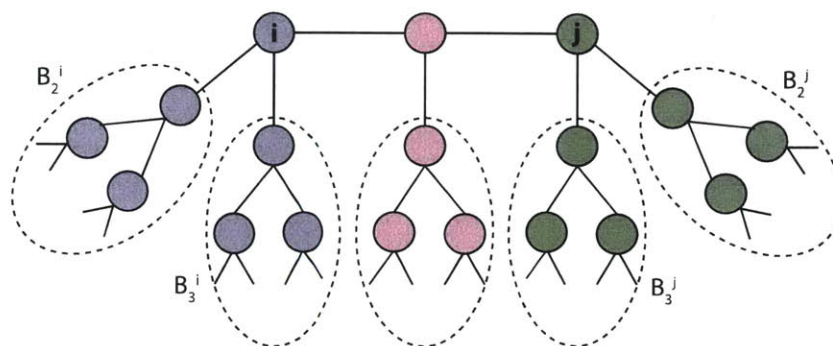


Figure 6-14: An example of a regular tree used in the proof of Theorem 6.20.

Inequality (III) follows from the log-sum inequality. Inequality (IV) follows from Lemma 6.30. This completes the proof of Theorem 6.17.

6.6.5 Proof of Theorem 6.20

To prove Theorem 6.20, first we prove regular trees are distance symmetric according to the following definition:

Definition 6.31. A graph $G = (V, E)$ is distance symmetric if for any pair of nodes $i, j \in V$, there exists a graph partition $\{V_1, V_2, V_3\}$ where,

- $\forall r \in V_1, d(i, r) = d(j, r)$. I.e., distances of nodes in V_1 from both nodes i and j are the same.
- There exists a bijective mapping function $\zeta(\cdot)$ between nodes V_2 and V_3 (i.e., $g: V_2 \rightarrow V_3$) such that for any $r \in V_2, d(i, r) = d(j, \zeta(r))$.

In the following Lemma, we show that regular trees are in fact distance symmetric:

Lemma 6.32. Let $G = (V, E)$ be a regular tree where V is countably infinite set of vertices and E is the set of edges. Then, G is distance-symmetric according to Definition 6.31.

Proof. Consider two distinctive nodes i and j . In the following, we construct graph partitions V_1, V_2 and V_3 satisfying conditions of Definition 6.31. Let degree of nodes in the regular tree be k . Thus, there are k branches connected to each nodes i and j ,

denoted by $\{B_1^i, \dots, B_k^i\}$ and $\{B_1^j, \dots, B_k^j\}$, respectively. Without loss of generality, assume $i \in B_1^j$ and $j \in B_1^i$. Add $\{B_2^i, \dots, B_k^i\}$ and $\{B_1^j, \dots, B_k^j\}$ to the partition sets V_2 and V_3 , respectively. A mapping function $\zeta(\cdot)$ of Definition 6.31 can be constructed between these sets by mapping nodes over branches B_l^i ($l \neq 1$) to nodes over branches B_l^j ($l \neq 1$) in a symmetric way (see Figure 6-14). Now consider the branch connecting nodes i and j in the graph. Let $\{1, 2, \dots, l\}$ be nodes over the shortest path connecting node i to node j . Therefore, $d(i, j) = l + 1$, the distance between nodes i and j .

- If l is odd, add non-partitioned nodes connected to nodes $\{1, \dots, l/2\}$ to the partition set V_2 . Similarly, add remaining nodes connected to nodes $\{l/2 + 1, \dots, l\}$ to the partition set V_3 .
- If l is even, add non-partitioned nodes connected to nodes $\{1, \dots, \lfloor l/2 \rfloor\}$ and $\{\lfloor l/2 \rfloor + 1, \dots, l\}$ to partition sets V_2 and V_3 , respectively. Non-partitioned nodes connected to the node $\lceil l/2 \rceil$ are assigned to the partition set V_1 .

A mapping function $\zeta(\cdot)$ of Definition 6.31 can be constructed between newly added nodes to partition sets V_2 and V_3 in a symmetric way. Moreover, nodes in the partition set V_1 have the same distance from both nodes i and j . This completes the proof. \square

Without loss of generality, suppose node 0 is the source node and we observe the infection pattern at time t according to the SI diffusion model. Thus, $Pr[y_j(t) = 1] = p_{0,j}(t)$, defined according to equation (6.13). Suppose we use parameter t' in Network Infusion Algorithm 6.14. According to equation (6.18), we have,

$$\begin{aligned}
 \mathbb{E}[\mathcal{H}_\alpha(i, t')] & \tag{6.47} \\
 &= \sum_{j \in V} (1 - \alpha) p_{0,j}(t) (1 - p_{i,j}(t')) + \alpha (1 - p_{0,j}(t)) p_{i,j}(t') \\
 &= \sum_{j \in V} p_{0,j}(t) (1 - p_{i,j}(t')) + \alpha (p_{i,j}(t') - p_{0,j}(t)).
 \end{aligned}$$

Thus, we have,

$$\begin{aligned}
 & \mathbb{E}[\mathcal{H}_\alpha(i, t')] - \mathbb{E}[\mathcal{H}_\alpha(0, t')] && (6.48) \\
 &= \sum_{j \in V} p_{0,j}(t) (p_{0,j}(t') - p_{i,j}(t')) + \alpha (p_{i,j}(t') - p_{0,j}(t')) \\
 &= \sum_{j \in V} (p_{0,j}(t') - p_{i,j}(t')) (p_{0,j}(t) - \alpha) \\
 &\stackrel{(a)}{=} \sum_{j \in V_2} (p_{0,j}(t') - p_{i,j}(t')) (p_{0,j}(t) - \alpha) \\
 &+ \sum_{j'=g(j) \in V_3} (p_{0,j'}(t') - p_{i,j'}(t')) (p_{0,j'}(t) - \alpha) \\
 &\stackrel{(b)}{=} \sum_{j \in V_2} (p_{0,j}(t') - p_{i,j}(t')) (p_{0,j}(t) - \alpha) \\
 &+ \sum_{j \in V_2} (p_{i,j}(t') - p_{0,j}(t')) (p_{i,j}(t) - \alpha) \\
 &= \sum_{j \in V_2} (p_{0,j}(t') - p_{i,j}(t')) (p_{0,j}(t) - p_{i,j}(t)).
 \end{aligned}$$

Equality (a) comes from partitioning nodes to sets V_1 , V_2 and V_3 according to Definition 6.31. The terms correspond to nodes in the partition set V_1 is equal to zero. Equality (b) comes from the fact that $d(0, j') = d(i, j)$ and $d(0, j) = d(i, j')$. Thus, $p_{0,j'}(\cdot) = p_{i,j}(\cdot)$ and $p_{i,j'}(\cdot) = p_{0,j}(\cdot)$. Therefore, if $t' = t$, we have,

$$\mathbb{E}[\mathcal{H}_\alpha(i, t')] - \mathbb{E}[\mathcal{H}_\alpha(0, t')] = \sum_{j \in V_2} (p_{0,j}(t) - p_{i,j}(t))^2, \quad (6.49)$$

which is strictly positive if $i \neq 0$. Now we consider the case that $t' \neq t$. Suppose $d_{0,j} < d_{i,j}$. Then, $p_{0,j}(t) > p_{0,j}(t')$ for any value of $t > 0$. Therefore, $(p_{0,j}(t') - p_{i,j}(t')) (p_{0,j}(t) - p_{i,j}(t)) > 0$. The same argument holds if $d_{0,j} > d_{i,j}$. If $d_{0,j} = d_{i,j}$, then $(p_{0,j}(t') - p_{i,j}(t')) (p_{0,j}(t) - p_{i,j}(t)) = 0$. This completes the proof of Theorem 6.20.

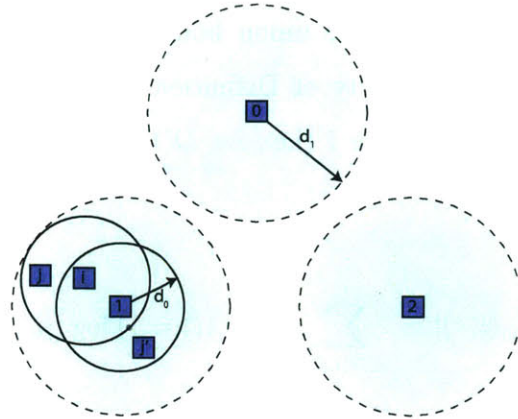


Figure 6-15: An illustrative figure of the proof of Theorem 6.26.

6.6.6 Proof of Theorem 6.26

To simplify notation, we prove this Theorem for a specific case where there are three sources in the network ($m = 3$). Also, we drop ϵ from d_1^c . All arguments can be extended to a general case. Let $\mathcal{S} = \{0, 1, 2\}$ be the sources. Suppose at the first step of the Algorithm 6.25, we have inferred the source node 0. We show that at the next step, we have,

$$\mathbb{E}[\mathcal{L}_{d_0}(s, t)] \geq \mathbb{E}[\mathcal{L}_{d_0}(i, t)], \quad (6.50)$$

for $s \in \{1, 2\}$ and for all $i \in V^t - D(0, d_1)$. Consider a node i in the d_1 -neighborhood of the source node 1 (i.e., $i \in D(1, d_1)$). Consider a node j in d_0 -neighborhood of node i (Figure 6-15). According to Equation (6.24), we have,

$$\begin{aligned} \Pr(y_j(t) = 1) &\geq p_{1,j} && (6.51) \\ \Pr(y_j(t) = 1) &\stackrel{(I)}{\leq} \sum_{s=1}^3 p_{s,j} \\ &\stackrel{(II)}{\leq} p_{1,j} + \frac{\epsilon}{n}, \end{aligned}$$

where Inequality (I) comes from the union bound of probabilities, and Inequality (II) uses incoherent source property of Definition 6.22. Consider a node i in the d_1 -neighborhood of the source node 1 (i.e., $i \in D(1, d_1)$). For this node, we have,

$$\begin{aligned}
 \mathbb{E}[\mathcal{L}_{d_0}(i, t)] &= \sum_{j \in D(i, d_0)} Pr(y_j(t) = 1) \log(p_{i,j}(t)) & (6.52) \\
 &+ Pr(y_j(t) = 0) \log(1 - p_{i,j}(t)) \\
 &\stackrel{(III)}{=} \sum_{j \in D(i, d_0)} (p_{1,j}(t) + \frac{\epsilon_j}{n}) \log(p_{i,j}(t)) \\
 &+ (1 - p_{1,j}(t) - \frac{\epsilon_j}{n}) \log(1 - p_{i,j}(t)) \\
 &= \sum_{j \in D(i, d_0)} p_{1,j}(t) \log(p_{i,j}(t)) \\
 &+ (1 - p_{1,j}(t)) \log(1 - p_{i,j}(t)) \\
 &+ \underbrace{\sum_{j \in D(i, d_0)} \frac{\epsilon_j}{n} \log \frac{p_{i,j}(t)}{1 - p_{i,j}(t)}}_{\text{term IV}} \\
 &\asymp \sum_{j \in D(i, d_0)} p_{1,j}(t) \log(p_{i,j}(t)) \\
 &+ (1 - p_{1,j}(t)) \log(1 - p_{i,j}(t)),
 \end{aligned}$$

where Equality (I) comes from Equation (6.51), and term (IV) goes to zero for sufficiently large n and a fixed t . Similarly to the proof of Theorem 6.17, we have

$$\begin{aligned}
 &\sum_{j \in D(i, d_0) \cap D(1, d_0)} p_{1,j}(t) \log(p_{1,j}(t)) & (6.53) \\
 &+ (1 - p_{1,j}(t)) \log(1 - p_{1,j}(t)) \\
 &\geq \sum_{j \in D(i, d_0) \cap D(1, d_0)} p_{1,j}(t) \log(p_{i,j}(t)) \\
 &+ (1 - p_{1,j}(t)) \log(1 - p_{i,j}(t)).
 \end{aligned}$$

Note that this inequality holds for nodes in the d_0 -neighborhood of both nodes i and 1. Now consider a node $j \in D(i, d_0) - D(1, d_0)$, and a node $j' \in D(1, d_0) - D(i, d_0)$ (see Figure 6-15). Owing to the symmetric structure of the network, similarly to Lemma 6.32, there is a one-to-one map among nodes j and j' such that $d(i, j) = d(1, j')$. For such node pairs j and j' , we have,

$$p_{1,j'}(t) \log(p_{1,j'}(t)) + (1 - p_{1,j'}(t)) \log(1 - p_{1,j'}(t)) \quad (6.54)$$

$$\begin{aligned} & - p_{1,j}(t) \log(p_{1,j}(t)) - (1 - p_{1,j}(t)) \log(1 - p_{1,j}(t)) \\ & = p_{1,j'}(t) \log(p_{1,j'}(t)) + (1 - p_{1,j'}(t)) \log(1 - p_{1,j'}(t)) \\ & - p_{1,j}(t) \log(p_{1,j}(t)) - (1 - p_{1,j}(t)) \log(1 - p_{1,j}(t)) \\ & = (p_{1,j'}(t) - p_{1,j}(t)) \log \frac{p_{1,j'}(t)}{1 - p_{1,j'}(t)} \end{aligned}$$

$$\geq 0, \quad (6.55)$$

where the inequality comes from the fact that $d(1, j') < d(1, j)$. Thus, we have,

$$\begin{aligned} & \sum_{j \in D(1, d_0) - D(1, d_0)} p_{1,j}(t) \log(p_{1,j}(t)) \quad (6.56) \\ & + (1 - p_{1,j}(t)) \log(1 - p_{1,j}(t)) \\ & \geq \sum_{j \in D(i, d_0) - D(1, d_0)} p_{1,j}(t) \log(p_{1,j}(t)) \\ & + (1 - p_{1,j}(t)) \log(1 - p_{1,j}(t)). \end{aligned}$$

Combining Inequalities (6.53) and (6.56), we have,

$$\mathbb{E}[\mathcal{L}_{d_0}(1, t)] \geq \mathbb{E}[\mathcal{L}_{d_0}(i, t)], \quad (6.57)$$

for any node i in the d_1 -neighborhood of the source node 1. The same arguments can be repeated for nodes in the d_1 -neighborhood of the source node 2. There are some remaining nodes that are not in the d_1 -neighborhood of the sources. As the last step of the proof, we show that the probability of having an infected remaining node is small. Consider node j such that $d(j, \mathcal{S}) > d_1$. according to Equation (6.24) and using the probability union bound, we have

$$\Pr(y_j(t) = 1) \leq \frac{\epsilon}{n}. \quad (6.58)$$

Let p_e denote the probability of at least one such infected node exists. We have,

$$\begin{aligned} p_e &\leq 1 - \left(1 - \frac{\epsilon}{n}\right)^n \\ &\asymp \epsilon, \end{aligned} \quad (6.59)$$

for sufficiently large n . This completes the proof of Theorem 6.25.

6.6.7 Proof of Proposition 6.27

Computation of the k -path network diffusion kernel for infected nodes has computational complexity $\mathcal{O}(|V^t|(k|E| + kn \log(n)))$ according to Proposition 6.8. Moreover, solving Optimization (6.27) for m iterations costs $\mathcal{O}(|V^t|nm)$. Thus, the total computational complexity of Algorithm 6.25 is $\mathcal{O}(|V^t|(k|E| + kn \log(n) + mn))$.

Chapter 7

Network Clustering

.

Clustering over networks is a classical problem with applications in systems biology, social sciences, and other fields [260–263]. Although most formulations of the clustering problem are NP-hard [264], several approaches have yielded useful approximate algorithms. The most well-studied approach is spectral clustering. Most spectral methods are not based on a particular generative network model; alternative, model-based approaches have also been proposed, using loopy belief propagation [265], variational Bayes [266], Gibbs sampling [267], and semidefinite programming [268, 269].

Many spectral clustering methods are derived by proposing a discrete optimization problem, and relaxing it to obtain a continuous, convex optimization whose solution is given by the eigenvectors of a normalized adjacency matrix or Laplacian. A post-processing step, typically k -means, is used to extract clusters from these eigenvectors [270, 271]. Different matrix normalizations/transformations include Modularity [260], Laplacian [272], normalized Laplacian [273, 274], and Bethe Hessian [198]. These methods are often characterized theoretically in the context of the Stochastic Block Model (SBM) [275], a simple and canonical model of community structure.

An alternative approach is to consider a network model with continuous latent positions, instead of discrete community memberships. This approach is attractive because it allows for a richer class of network models to be considered while simultaneously leading to a continuous inference problem that might be more tractable than discrete community-membership inference optimizations. A particularly interesting latent space model is the *Random Dot Product Graph* (RDPG) [276]. In the RDPG, the probability of observing an edge between two nodes is the dot product between their respective latent position vectors. Another latent space model related to the RDPG was introduced by [277]. These models are discussed further in Section 7.1. The RDPG provides a useful perspective on spectral clustering, in two ways. First, it has led to theoretical advances, including a central limit theorem for eigenvectors of an adjacency matrix [197] and a justification for k -means clustering as a post processing step for spectral clustering [195]. Second, as a natural extension of the SBM, the RDPG describes more comprehensively the types of network structures, such as variable degrees and mixed community memberships, that can be inferred using spectral

methods.

However, the spectral inference method over the RDPG is based on the heuristic that the first eigenvectors of the adjacency matrix recapitulate the non-null eigenvectors of the inner-product matrix [276]. This heuristic can fail, since the first eigenvector of the adjacency matrix often separates high-degree nodes from low-degree nodes, rather than separating the communities. Indeed, many of the normalizations that are commonly used in spectral clustering can be viewed as solutions to this problem. This problem can occur even in simple cases, such as the symmetric SBM with two clusters of equal size and density. Therefore, we were motivated to develop a more robust inference approach.

In this chapter, we consider a closely related latent space model, the *Logistic RDPG*, which uses a logistic link function mapping from latent positions to edge probabilities. Like the previously-studied RDPG, the logistic RDPG includes most SBMs as well as other types of network structure, including a variant of the degree corrected SBM. Over this model, we show that the maximum likelihood latent-position inference problem admits an asymptotically exact spectral solution. Our method is to take the top eigenvectors of the mean-centered adjacency matrix and to scale them using a logistic regression step. This result is possible because over the logistic model specifically, the likelihood function separates into a linear term that depends on the observed network and a nonlinear penalty term that does not. Because of its simplicity, the penalty term admits a Frobenius norm approximation, leading to our spectral algorithm. A similar approximation is not obtained using other link functions besides the logistic link.

We show that the likelihood of the approximate solution approaches the maximum of the likelihood function when the graph is large and the latent-position magnitudes go to zero. The asymptotic regime is not overly restrictive, as it encompasses many large SBMs at or above the detectability threshold [265]. Moreover, through simulations we show that our method performs well for finite graphs over a broad range of clustering models (in fact, it outperforms other existing spectral and SDP-based methods). Ours is the first network model over which maximum likelihood inference

has been shown to be tractable using a spectral method.

We compare the performance of our method over synthetic networks with several spectral methods, including the Modularity method [260], the Normalized Laplacian [272], and the Bethe Hessian [198]. We show that our method outperforms these methods, as well as SDP-based network clustering methods [268, 269], at both latent position inference and community recovery.

The rest of the chapter is organized as follows: In Section 7.1, we review prior work relating to spectral clustering and latent space models. In Section 7.2, we introduce the Logistic RDPG, discuss relationship between RDPGs and SBMs, and state our main result. In Section 7.3, we evaluate the performance of our method. We present proofs of our main results in Section 7.4.

7.1 Prior Work

Spectral clustering is a term that encompasses several related methods. In the simplest case of two clusters, the adjacency matrix A is normalized and an eigenvector of the normalized matrix is computed. The signs of the eigenvector entries are used to assign community memberships. When the number of clusters is greater than two, k -means is typically used to extract cluster labels from the eigenvectors [271].

For most spectral clustering methods, a heuristic objective function is associated with this technique. This connection was discovered by Hagen and Kahng [278], who showed that the unnormalized Laplacian is associated with the RatioCut objective function. The normalized Laplacian is associated with the Ncut objective function [274], and the Modularity matrix is associated with the Modularity function [157]. In each instance, an NP-hard discrete optimization problem is relaxed to obtain a tractable continuous optimization.

The performance of these methods is often characterized in the context of a theoretically tractable model, the Stochastic Blockmodel (SBM), which was introduced by Holland et al. [275]. In the SBM, each node is assigned to one community, and edges are drawn independently between nodes with probabilities determined by the

communities of the respective nodes. Theoretical bounds on the detectability of network structure for large networks have been established for this model. For a large SBM with two equally-sized communities, if the within-cluster density is p and the outside-of-cluster density is $q < p$, it was conjectured by Dece et al. [265] and proven by Mossel et al. and Massoulié [279–281] that it is possible to efficiently estimate the community labels more accurately than random guessing if and only if

$$(p - q)^2/n > 2(p + q).$$

Several spectral methods have been shown to achieve this recovery threshold [157, 198, 282]. Strong theoretical and empirical results have also been obtained using SDP-based [268, 269] and Belief Propagation-based [265] methods, which often have higher computational complexity than spectral methods. A threshold has also been discovered for *perfect* clustering, i.e. when community structure can be recovered with zero errors [268, 283].

An alternative approach to the clustering problem is to invoke a latent space model. Each node is assigned a latent position v_i , and the edges of the graph are drawn independently with probability $p_{ij} = g(v_i, v_j)$ for some $g(\cdot)$. Hoff et al. [277] considered two latent space models, the first of which was a distance model,

$$P_{i,j} = l(\|v_i - v_j\| - \mu + \beta x_{ij}), \quad l(x) = 1/(1 + \exp(-x))$$

where edge probabilities depend on the Euclidean distance between two nodes. x_{ij} is a fixed covariate term, which is not learned. Their second model is called a projection model:

$$P_{i,j} = l\left(\frac{v_i \cdot v_j}{\|v_j\|} + \beta x_{ij} - \mu\right). \quad (7.1)$$

Hoff et al. suggest to perform inference using an MCMC approach over both models. Focusing on the distance model, they have also extended their approach to allow the latent positions to themselves be drawn from a mixture distribution containing

clusters [284,285]. Efforts to improve the computational efficiency have been made in references [286,287], and with a related mixed membership blockmodel in reference [266].

Young et al. [276] introduced the Random Dot Product Graph (RDPG), in which the probability of observing an edge between two nodes is the dot product between their respective latent position vectors. The RDPG model can be written as

$$P_{i,j} = v_i \cdot v_j.$$

This model is related to the projection model of Hoff et al. (7.1).

Let V be the $n \times d$ matrix of latent positions, where n is the number of nodes and d is the dimension of the latent vectors ($d \leq n$). Sussman et al. [195] proposed an inference method over the RDPG based on the heuristic that the first eigenvectors of A will approximate the singular vectors of V , as $E(A) = VV^T$. They also characterized the distribution of the eigenvectors of A given V . Related work by Rohe et al. [288] characterized the asymptotic behavior of the eigenvectors of a normalized Laplacian, showing that they converge to the eigenvectors of a *population* Laplacian. Likewise, the inference method of Sussman et al. is shown to be consistent under some conditions.

7.2 Logistic Random Dot Product Graphs

In this section, we introduce the logistic RDPG, describe our inference method, and show that it is asymptotically equivalent to the maximum likelihood inference.

7.2.1 Definitions

Let \mathcal{A} be the set of adjacency matrices corresponding to undirected, unweighted graphs of size n .

Definition 7.1 (Stochastic Block Model). The Stochastic Block Model is a family of distributions on \mathcal{A} parameterized by (k, c, Q) , where k is the number of communities,

$c \in [k]^n$ is the community membership vector and $Q \in [0, 1]^{k \times k}$ is the matrix of community relationships. For each pair of nodes, an edge is drawn independently with probability

$$P_{i,j} := Pr(A_{ij} = 1 | c_i, c_j, Q) = Q_{c_i, c_j}.$$

Another network model that characterizes low dimensional structures is the *Random Dot Product Graph* (RDPG). This class of models includes many SBMs.

Definition 7.2 (Random Dot Product Graph). The Random Dot Product Graph with link function $g(\cdot)$ is a family of distributions on \mathcal{A} parameterized by an $n \times d$ matrix of latent positions $V \in \mathbb{R}^{n \times d}$. For each pair of nodes, an edge is drawn independently with probability

$$P_{i,j} := Pr(A_{ij} = 1 | V) = g(v_i \cdot v_j), \quad (7.2)$$

where v_i , the i -th row of the matrix V , is the latent position vector assigned to node i .

The RDPG has been formulated using a general link function $g(\cdot)$ [276]. The linear RDPG, using the identity link, has been analyzed in the literature because it leads to a spectral inference method. We will refer to this model as either the linear RDPG or as simply the RDPG. In this chapter, we consider the Logistic RDPG:

Definition 7.3 (Logistic RDPG). The logistic RDPG is the RDPG with link function:

$$g(x) = l(x - \mu), \quad l(x) := \frac{1}{1 + e^{-x}}, \quad (7.3)$$

where μ is the offset parameter of the logistic link function.

Note that this model is similar to the projection model of Hoff et al. [277] (7.1). The projection model is for a directed graph, with $P_{i,j} \neq P_{j,i}$ owing to the division by $\|v_j\|$.

Remark 7.4. The parameter μ in the logistic RDPG controls the sparsity of the network. If the latent position vector lengths are small, the density of the graph is

$$\frac{1}{n(n-1)} \sum_{i < j} \mathbb{E}(A_{ij}) \approx l(-\mu). \quad (7.4)$$

A logistic RDPG with this property is called *centered*. In Section 7.2.3, we show that asymptotically exact maximum likelihood inference of latent positions over the centered logistic RDPG can be performed using an efficient spectral algorithm.

Because we will ultimately apply our method to identify communities, it is important to understand what types of community structure are accommodated by our model. Thus, we consider which SBMs can be equivalently modeled as RDPGs. An SBM that can be equivalently described as an RDPG with some link function is called *admissible* under that model:

Definition 7.5 (Admissible SBM). A SBM with community relationship parameter Q is said to be *admissible* under the RDPG with a link function $g(\cdot)$ if $g^{-1}(Q) \succeq 0$.

In Section 7.2.2, we discuss which SBMs are and are not admissible under the linear and the logistic RDPGs.

For a general RDPG, the ML inference problem is:

Definition 7.6 (ML inference problem for the RDPG). Let $A \in \mathcal{A}$. The ML inference problem over the RDPG is:

$$\begin{aligned} \max_X \quad & \sum_{i \neq j} A_{ij} \log g(X_{i,j}) + (1 - A_{ij}) \log(1 - g(X_{i,j})), \\ & X = VV^T, \quad V \in \mathbb{R}^{n \times d}. \end{aligned} \quad (7.5)$$

Note that in the undirected case, this objective function is twice the log likelihood.

Remark 7.7. A convex semidefinite relaxation of Optimization (7.5) can be obtained

for some link functions as follows:

$$\begin{aligned} \max_X \quad & \sum_{i \neq j} A_{ij} \log g(X_{ij}) + (1 - A_{ij}) \log(1 - g(X_{ij})), \\ & X \succeq 0, \end{aligned} \tag{7.6}$$

where $X \succeq 0$ means that X is psd. For example, this optimization is convex for the logistic link function (i.e., $g(x) = l(x - \mu)$) and for the linear link function (i.e., $g(x) = x$). However, this optimization can be slow in practice, and it often leads to excessively high-rank solutions.

The ML community recovery problem for an SBM can be written similarly to Optimization (7.5) with the additional constraint that nodes in the same community are assigned to the same latent position vector:

Definition 7.8 (ML inference problem for the admissible SBM). The ML inference problem for an admissible SBM over an RDPG with the link function $g(\cdot)$ can be written as follows:

$$\begin{aligned} \max_{X, U} \quad & \sum_{i \neq j} A_{ij} \log g(X_{ij}) + (1 - A_{ij}) \log(1 - g(X_{ij})), \\ & X = VV^T, \quad v_i \in \{u_1, \dots, u_k\} \subset \mathbb{R}^{n \times d}, \quad 1 \leq i \leq n. \end{aligned} \tag{7.7}$$

RDPGs provide a natural relaxation of admissible SBMs, removing the constraint that every node must be assigned to one out of a few possible latent positions.

7.2.2 Relationship between RDPGs and SBMs

In this section we discuss the relationship between the SBM and the RDPG, both linear and logistic; in particular, we give results that characterize which SBMs are and are not admissible under these models.

First, we give a result that characterizes which SBMs are admissible under the logistic RDPG for some value of μ :

Proposition 7.9. Let Q be a matrix parameterizing a SBM. Let $Z = l^{-1}(Q) - t$, where $t = \frac{1}{n^2} \sum_{i,j} l^{-1}(Q_{i,j})$. Let $z = \mathbf{1}_n Z$ be the vector of row sums of Z . Let ν_1, \dots, ν_n and $\lambda_1, \dots, \lambda_n$ be eigenvectors and eigenvalues of Z , respectively. Then this SBM is admissible under the logistic RDPG if and only if for every eigenvalue λ_i , either:

1. λ_i is positive;
2. λ_i is zero, and $\nu_i \cdot z = 0$; or
3. λ_i is negative, and ν_i is the constant vector $\frac{1}{\sqrt{n}} \mathbf{1}_n$.

If these conditions hold, then the following function is bounded above on $C := \{u : u^T \mathbf{1}_n = 0\}$:

$$f(u) = \frac{(u \cdot z)^2}{u^T Z u}. \quad (7.8)$$

Under these conditions, suitable values of μ are $\mu \geq \mu_0$, where

$$\mu_0 = \frac{\max_{u \in C} f(u)}{n^2} - t. \quad (7.9)$$

Proof. A real symmetric matrix M is psd if and only if the quadratic form $Q(x) = x^T M x$ is nonnegative. Consider the quadratic form

$$Q(v) = v^T (Z + t + \mu) v.$$

Let $u = v - s \mathbf{1}_n$, with $s = v^T \mathbf{1}_n / n$ so that $u^T \mathbf{1}_n = 0$. Then:

$$\begin{aligned} Q(v) &= Q(u + s \mathbf{1}_n) = \sum_{i,j} (Z_{ij} + t + \mu) (u_i + s)(u_j + s) \\ &= u^T Z u + 2u \cdot z s + (\mu + t) n^2 s^2 \end{aligned} \quad (7.10)$$

The first term comes from u , Z , and again u ; the second term comes from the cross-terms with u , Z , $s \mathbf{1}_n$; the third term comes from $s \mathbf{1}_n$ and $t + \mu$. The other terms involving $t + \mu$ are zero because $|u| = 0$ and the other term involving Z is zero because $|Z| = 0$.

Equation (7.10) is a quadratic equation in s . It is nonnegative for all v and s if and only if for all $u \in C$, $u^T Z u$ is nonnegative and the discriminant is nonpositive:

$$D = (2u \cdot z)^2 - 4n^2(\mu + t)u^T Z u \leq 0.$$

If $u^T Z u$ is positive for all u , then D is nonpositive for sufficiently large μ . The sign of D does not depend on the magnitude of u . Thus without loss of generality we may assume that $\|u\|^2 = 1$, and $u^T u$ cannot be arbitrarily close to zero on $\{u \in C : \|u\|^2 = 1\}$ because this is a closed set. If $u^T u$ is zero (but not negative) for some $u \in C$, then D is nonnegative for some μ if and only if $u \cdot z = 0$ for all such u . If and only if these conditions hold, then the following rational function is bounded above on C :

$$f(u) = \frac{(u \cdot z)^2}{u^T Z u}. \quad (7.11)$$

The bound indicates that $Z + t + \mu$ is psd for $\mu \geq \mu_0$, where

$$\mu_0 = \frac{\max_{u \in C} f(u)}{n^2} - t. \quad (7.12)$$

To further simplify the conditions under which $f(\cdot)$ is bounded, diagonalize Z so that the eigenvectors are ν_1, \dots, ν_n and the eigenvalues are $\lambda_1, \dots, \lambda_n$. The quadratic form can be written as

$$v^T Z v = \sum_{i=1}^n \lambda_i (\nu_i^T v)^2.$$

If there is a negative eigenvalue whose eigenvector is not the constant vector, then there exists u such that $|u| = 0$ and $u^T Z u$ is negative, and the SBM is not admissible. Otherwise, $u^T Z u$ is nonnegative for all $|u| = 0$, but nonetheless $f(\cdot)$ may be unbounded, which occurs if and only if the kernel of Z is not orthogonal to z . \square

Proposition 7.9 gives a procedure for testing if a SBM is admissible under the logistic RDPG for a range of μ , but it gives little intuition. In the following, we present necessary conditions for admissible SBMs over both linear and logistic RDPGs; these

conditions are also sufficient conditions in the case of only two communities.

Proposition 7.10. Admissible SBMs over the linear and logistic RDPGs must satisfy

$$Q_{i,j}^2 \leq Q_{i,i}Q_{j,j}, \quad (7.13)$$

$$(g^{-1}(Q_{i,j}) + \mu)^2 \leq (g^{-1}(Q_{i,i}) + \mu)(g^{-1}(Q_{j,j}) + \mu), \quad (7.14)$$

respectively. In the case that the number of communities is equal to two, these conditions are also sufficient for a SBM to be admissible.

Proof. The condition is necessary because the dot product satisfies Cauchy-Schwartz. To see that it is sufficient in the case of two communities, observe that given numbers a, b, c with $a^2 \leq bc$, one can always find vectors u, v such that $\|u\| = b$, $\|v\| = c$, and $u \cdot v = a$. \square

In the following, we show that every SBM with two communities that is admissible under the linear RDPG is arbitrarily well-approximated by an SBM that is admissible under the logistic RDPG. The converse is not true, as we illustrate in Examples 7.12 and 7.13.

Proposition 7.11. Suppose that $\tilde{Q} \in \mathbb{R}^{2 \times 2}$ parameterizes an SBM that is admissible under the linear RDPG. For every $\epsilon > 0$ there exists Q such that $\|Q - \tilde{Q}\|_F < \epsilon$, where the SBM parameterized by Q is admissible under the logistic RDPG for some μ .

Proof. The proof is an application of the inequality of arithmetic and geometric means (AM-GM). Choose Q in the ϵ -ball around \tilde{Q} such that its entries are in $(0, 1)$ (so that $g^{-1}(Q)$ exists), and such that the inequality (7.13) strictly holds for $i \neq j$. By Proposition 7.10, for any Q , condition (7.14) can be satisfied if there exists $\mu > 0$ such that

$$\frac{(g^{-1}(Q_{ij}) + \mu)^2}{\mu} \leq \frac{(g^{-1}(Q_{ii}) + \mu)(g^{-1}(Q_{jj}) + \mu)}{\mu}, \quad (7.15)$$

which can be written as

$$\begin{aligned} & \mu + 2g^{-1}(Q_{ij}) + O(1/\mu) \\ & \leq \mu + g^{-1}(Q_{ii}) + g^{-1}(Q_{jj}) + O(1/\mu). \end{aligned} \quad (7.16)$$

This holds for sufficiently large μ if

$$\begin{aligned} & 2 \log \frac{Q_{ij}}{1 - Q_{ij}} < \log \frac{Q_{ii}}{1 - Q_{ii}} + \log \frac{Q_{jj}}{1 - Q_{jj}} \\ \iff & \left(\frac{Q_{ij}}{1 - Q_{ij}} \right)^2 < \frac{Q_{ii}Q_{jj}}{(1 - Q_{ii})(1 - Q_{jj})}. \end{aligned} \quad (7.17)$$

We have arranged so that $Q_{ij}^2 < Q_{ii}Q_{jj}$. Thus

$$\begin{aligned} \left(\frac{Q_{ij}}{1 - Q_{ij}} \right)^2 & < \left(\frac{\sqrt{Q_{ii}Q_{jj}}}{1 - \sqrt{Q_{ii}Q_{jj}}} \right)^2 \\ & = \frac{Q_{ii}Q_{jj}}{1 - 2\sqrt{Q_{ii}Q_{jj}} + Q_{ii}Q_{jj}} \end{aligned} \quad (7.18)$$

By AM-GM, $2\sqrt{Q_{ii}Q_{jj}} \leq Q_{ii} + Q_{jj}$, and

$$\begin{aligned} \left(\frac{Q_{ij}}{1 - Q_{ij}} \right)^2 & < \frac{Q_{ii}Q_{jj}}{1 - (Q_{ii} + Q_{jj}) + Q_{ii}Q_{jj}} \\ & = \frac{Q_{ii}Q_{jj}}{(1 - Q_{ii})(1 - Q_{jj})}. \end{aligned} \quad (7.19)$$

□

In the following, we give an example of a two-community SBM that is admissible under the logistic model but not under the linear model.

Example 7.12. Consider the SBM parameterized by Q with $Q_{1,1} \approx 1$, $Q_{1,2} \approx 1$, and $Q_{2,2} \approx 0$. Such a model is admissible under a Logistic RDGP with $d = 1$ by making

$$v_1 \gg 1, \quad v_2 = 1, \quad 1 \ll \mu \ll v_1.$$

However, we have arranged so that $1 \approx Q_{1,2} \gg Q_{1,1}Q_{2,2} \approx 0$. Therefore, this SBM

violates (7.13) and it is not admissible, even up to an approximation, by a linear RDPG.

An additional necessary condition on admissible SBMs can be derived from the triangle inequality. Let $d_{i,j} := \sqrt{1 - \frac{Q_{i,j}}{Q_{i,i}Q_{j,j}}}$. If the SBM parameterized by Q is admissible under the linear RDPG, with latent positions V , then $d_{i,j}$ can be written as

$$d_{i,j} = \sqrt{1 - \frac{v_i \cdot v_j}{\|v_i\| \|v_j\|}}.$$

This function is a metric (Dongen and Enright, 2012). Thus it satisfies the triangle inequality:

$$d_{i,j} \leq d_{i,k} + d_{j,k}. \tag{7.20}$$

For a fixed μ , an analogous inequality holds for the logistic RDPG with

$$d_{i,j} = \sqrt{1 - \frac{g^{-1}(Q_{i,j}) - \mu}{\sqrt{(g^{-1}(Q_{i,i}) - \mu)(g^{-1}(Q_{j,j}) - \mu)}}}.$$

As μ grows with a fixed Q , this quantity converges to $\sqrt{2}$. Thus the constraint becomes trivial: $\sqrt{2} \leq 2\sqrt{2}$.

This observation leads to the following example of a SBM that is admissible under the logistic but not the linear RDPG.

Example 7.13. Consider the SBM parameterized by

$$Q = \begin{pmatrix} .99 & .98 & .01 \\ .98 & .99 & .98 \\ .01 & .98 & .99 \end{pmatrix}$$

This matrix is described by a logistic RDPG with latent positions

$$v_1 = (a, 0), \quad v_2 = (0, b), \quad v_3 = (-a, 0)$$

where a is sufficiently large and $-a^2 \ll \mu \ll 0$. However, we have arranged so that $1 \approx d_{1,3} \gg d_{1,2} + d_{2,3} \approx 0$. Therefore, this SBM violates (7.20) and it cannot be encoded by the linear RDPG.

7.2.3 Maximum-Likelihood Inference of Latent Position Vectors

Here, we present an efficient, asymptotically exact spectral algorithm for the maximum-likelihood (ML) inference problem over the logistic RDPG, subject to mild constraints. Our algorithm is asymptotically exact in the sense that the likelihood ratio between our solution and the true maximum converges uniformly to one as the average latent position length shrinks. Importantly, the convergence is uniform over arbitrarily large graphs; therefore, this regime contains most interesting large network models, such as an SBM with large communities that cannot be perfectly recovered. Coupling this algorithm with a k -means post processing step leads to a clustering method with robust performance under different network clustering setups. We assume that the number of clusters is given (and the number of dimensions, which is typically one fewer than the number of clusters); in practice, these parameters are often set manually, but approaches have been proposed to automatically detect the number of clusters [289].

The proposed maximum-likelihood inference of latent position vectors for the logistic RDPG is described in Algorithm 5. In the following, we sketch the derivation of this algorithm in a series of lemmas. Proofs for these assertions are presented in Section 7.4.

First, we simplify the likelihood function of Optimization (7.5) using the logistic link function. Let $F(X)$ be the log-likelihood, and let the link function be $g(x) =$

Algorithm 5 ML Inference for the logistic RDPG

Require: Adjacency matrix A , number of dimensions d . (optional) number of clusters k

Form the mean-centered adjacency matrix $B := A - 1/(n(n-1))\|A\|$.

Compute d eigenvectors of B with largest eigenvalues: e_1, \dots, e_d .

Let $X_i = e_i e_i^T$ for $1 \leq i \leq d$.

Perform logistic regression of the entries of A lying above the diagonal on the corresponding entries of X_1, \dots, X_d , estimating coefficients $\lambda_1^*, \dots, \lambda_d^*$ subject to the constraint that $\lambda_i \geq 0 \quad \forall i$.

Let V be the matrix formed by concatenating $\sqrt{\lambda_1^*}e_1, \dots, \sqrt{\lambda_d^*}e_d$. Return V .

(optional) Perform k -means on V , and return the inferred clusters.

$l(x - \mu)$. Then:

$$\begin{aligned}
 F(X) &:= \sum_{i,j} A_{ij} \log l(X_{ij} - \mu) + (1 - A_{ij}) \log(1 - l(X_{ij} - \mu)) \\
 &= \sum_{i,j} A_{ij} \log \frac{l(X_{ij} - \mu)}{1 - l(X_{ij} - \mu)} + \log(1 - l(X_{ij} - \mu)) \\
 &= \sum_{i,j} A_{ij} (X_{ij} - \mu) + \log(1 - l(X_{ij} - \mu)). \tag{7.21}
 \end{aligned}$$

We have used that $l(x)/(1 - l(x)) = x$. The maximum likelihood problem takes the following form (for given μ):

$$\begin{aligned}
 \max_X \quad & Tr(AX) + \sum_{i \neq j} \log\left(1 - \frac{1}{1 + e^{-(X_{ij} - \mu)}}\right), \\
 X &= VV^T, \quad V \in \mathbb{R}^{n \times d}. \tag{7.22}
 \end{aligned}$$

The objective function has been split into a linear term that depends on the adjacency matrix A and a penalty term that does not depend on A . This simplification is what leads to tractable optimization, and it is the reason that the logistic link is needed; using e.g. a linear link, an optimization of this form is not obtained.

We define a penalty function $f(x)$ that keeps only the quadratic and higher-order terms in the penalty term of (7.22). Let

$$f(x) := -(h(x) - h(0) - h'(0)x), \quad h(x) = \log(1 - l(x - \mu)).$$

Now, $h'(0) = -l(-\mu)$. Let

$$B := A - l(-\mu)\mathbb{K}_{n \times n}$$

in order to re-write Optimization (7.22) as:

$$\begin{aligned} \max_X \quad & Tr(BX) - \sum_{i \neq j} f(X_{ij}), \\ & X = VV^T, \quad V \in \mathbb{R}^{n \times d}. \end{aligned} \tag{7.23}$$

Note that for a centered RDPG with average density \bar{A} , $\mu = l^{-1}(\bar{A})$, and B is the mean-centered adjacency matrix $A - \bar{A}\mathbb{K}_{n \times n}$. In the next step, we convert the penalty term in the objective function into a constraint:

Lemma 7.14. Suppose that X^* is the optimal solution to Optimization (7.5). Let

$$h^* := \frac{1}{n(n-1)} \sum_{i,j} f(X_{ij}^*).$$

Then X^* is also the solution to the following optimization:

$$\begin{aligned} \max_X \quad & Tr(BX) \\ & X = VV^T, \quad V \in \mathbb{R}^{n \times d} \\ & \frac{1}{n(n-1)} \sum_{i \neq j} f(X_{ij}) \leq h^*. \end{aligned} \tag{7.24}$$

In the following key lemma, we show that the inequality constraint of Optimization (7.24) can be replaced by its second order Taylor approximation.

Lemma 7.15. For any $\epsilon > 0$ and $\gamma \geq 1$, there exists $\delta > 0$ such that for any graph whose ML solution X^* satisfies

$$h^* \leq \delta \quad \text{and} \quad \max_i X_{ii}^* \leq \frac{\gamma}{n} \sum_i X_{ii}^*, \tag{7.25}$$

the following bound is satisfied. Let B be the mean centered adjacency matrix of

the chosen graph. Let $s^* \in \mathbb{R}$ be the optimal value of the following optimization, obtained at $X = X^*$:

$$\begin{aligned} \max_X \quad & \text{Tr}(BX), \\ & X = VV^T, \quad V \in \mathbb{R}^{n \times d} \\ & \frac{1}{n(n-1)} \sum_{i \neq j} f(X_{ij}) \leq h^*, \end{aligned} \tag{7.26}$$

Let \tilde{s} be the optimal value of the following optimization:

$$\begin{aligned} \max_X \quad & \text{Tr}(BX), \\ & X = VV^T, \quad V \in \mathbb{R}^{n \times d} \\ & \frac{a_2}{n(n-1)} \sum_{i \neq j} X_{ij}^2 \leq h^*. \end{aligned} \tag{7.27}$$

where $a_2 := f''(-\mu)$. Then

$$\tilde{s} \geq (1 - \epsilon)s^*.$$

The parameter h^* is related to the average length of latent-position vectors (X_{ii}^*) . If these lengths approach zero, h^* approaches zero, for a fixed γ . An implication of this constraint is that the logistic RDPG must be approximately centered. Thus, there is a natural choice for the parameter μ for the purpose of inference:

$$\hat{\mu} = -l^{-1}\left(\frac{1}{n(n-1)}\|A\|\right). \tag{7.28}$$

This estimator of μ can be viewed as the maximum-likelihood estimator specifically over the centered logistic RDPG. With this choice of μ , B , the mean-centered adjacency matrix, can be written as

$$B = A - \frac{1}{n(n-1)}\|A\|.$$

Note that changing the constant in the inequality constraint of Optimization (7.27) only changes the scale of the solution, since the shape of the feasible set does not

change. Thus, in this optimization we avoid needing to know h^* *a priori* (as long as the conditions of Lemma 7.15 are satisfied).

Next we show that that the solution to Optimization (7.27) can be recovered up to a linear transformation using spectral decomposition:

Lemma 7.16. Let \tilde{X} be the optimal solution to Optimization (7.27). Let e_1, \dots, e_d be the first d eigenvectors of B , corresponding to the largest eigenvalues. Then e_1, \dots, e_d are identical to the non-null eigenvectors of \tilde{X} , up to rotation.

Once the eigenvectors of \tilde{X} are known, it remains only to recover the corresponding eigenvalues. Instead of recovering the eigenvalues of \tilde{X} , we find the eigenvalues that maximize the likelihood, given the eigenvectors of \tilde{X} . Let $\tilde{X}_i = e_i e_i^T$. Then, the maximum-likelihood estimate of $\lambda_1, \dots, \lambda_d$ conditional on $X_1, \dots, X_d = \tilde{X}_1, \dots, \tilde{X}_d$ can be written as follows:

$$\lambda^* := \operatorname{argmax}_{\lambda=(\lambda_1, \dots, \lambda_d)} \sum_{i \neq j} \log P(A_{ij} | \tilde{X}_1, \dots, \tilde{X}_d, \lambda, \mu). \quad (7.29)$$

Lemma 7.17. Optimization (7.29) can be solved by logistic regression of the entries of A on the entries of X_1^*, \dots, X_d^* , with the constraint that the coefficients are nonnegative, and with intercept μ .

These lemmas can be used to show the asymptotic optimality of Algorithm 5:

Theorem 7.18. For all $\epsilon > 0$ and $\gamma > 1$, there exists $\delta > 0$ that satisfies the following. For any graph with size n and adjacency matrix A , suppose that X^* is the solution to the optimization

$$\begin{aligned} \max_X \quad & P(A|X), \\ & X = VV^T, \quad V \in \mathbb{R}^{n \times d}. \end{aligned}$$

Let

$$h^* := \frac{1}{n(n-1)} \sum_{i \neq j} f(X_{ij}^*).$$

If

$$h^* < \delta \quad \text{and} \quad \max_i X_{ii}^* \leq \frac{\gamma}{n} \sum_i X_{ii}^*, \quad (7.30)$$

then

$$\frac{P(A|X = X_*)}{P(A|X = X^*)} > 1 - \epsilon,$$

where X_* is the solution obtained by Algorithm 5.

This result is stronger than the statement that an approximate objective function approximates the likelihood function at the optimum of the likelihood function. Such a result, which can be obtained for many link functions (such as an affine-linear link), is not useful because it does not follow that the optimum of the approximate function lies near the optimum of the true likelihood function. Indeed, for a linear approximation, it has no optimum since the objective function is unbounded. In order to obtain the stronger statement that the likelihood at the optimum of the approximation is large, it is necessary to use a quadratic approximation. For link functions besides the logistic link, the quadratic term in the likelihood function depends on A , and a spectral optimization method cannot be obtained.

The condition in Theorem 7.18 that the lengths of optimal latent vectors are sufficiently small is not restrictive for large networks. Consider a sequence of increasingly-large SBMs with two clusters of fixed relative sizes, and a convergent sequence of admissible connectivity matrices whose average density is fixed. There are three asymptotic regimes for the community structure: (1) in which the structure of the network is too weak to detect any clusters at all; (2) in which the communities can be partially recovered, but some mis-assignments will be made; and (3) in which the communities can be recovered perfectly. The true latent position lengths go to zero in regimes (1) and (2) as well as in part of regime (3) [268]. We expect that the maximum-likelihood latent position lengths will also go to zero. Thus, if ML achieves the optimum thresholds for partial recovery and perfect recovery, then our method will as well.

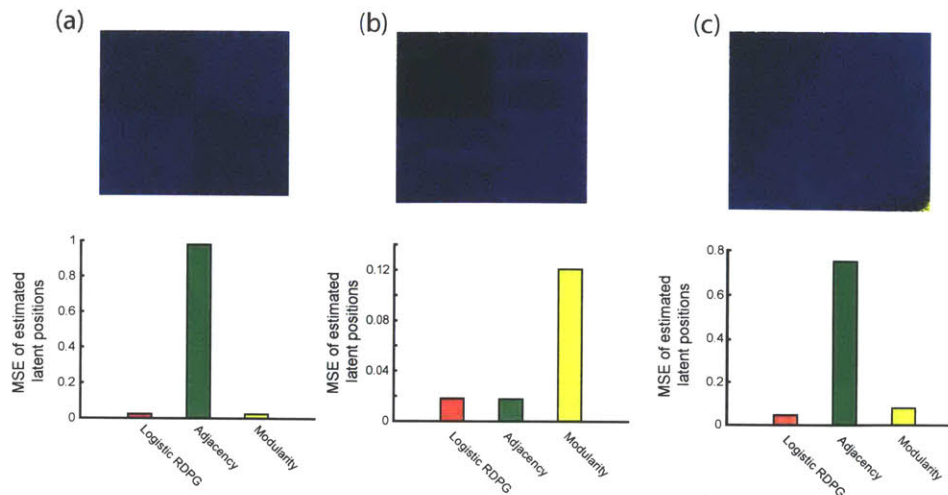


Figure 7-1: Normalized mean squared error (one minus squared correlation) of inferred latent positions for two SBMs (a-b) and a non-SBM logistic RDPG (c). The top eigenvectors of the adjacency matrix A and the modularity matrix M do not characterize the community structure in panel (a) and in panel (b), respectively. In panel (c), the top eigenvector of A does not recover the latent structure. In contrast, our method successfully recovers the underlying latent position vectors in all cases.

7.3 Performance Evaluation over Synthetic Networks

In this section, we compare the performance of our proposed method (Algorithm 5) with existing methods. First, we assess the performance of our algorithm against existing methods in inference of latent position vectors of two standard SBMs depicted in Figure 7-1. The network demonstrated in panel (a) has two dense clusters. In this case, the first eigenvector of the modularity matrix M leads to a good estimation of the latent position vector while the first eigenvector of the adjacency matrix A fails to characterize this vector. This is because the first eigenvector of the adjacency matrix correlates with node degrees. The modularity transformation regresses out the degree component and recovers the community structure. However, the top eigenvector of the modularity matrix fails to identify the underlying latent position vector when there is a single dense cluster in the network, and the community structure *is* correlated with node degrees (Figure 7-1-b). This discrepancy highlights the sensitivity of existing heuristic inference methods in different network models (the Modularity method has not previously been considered a latent-position inference method, but

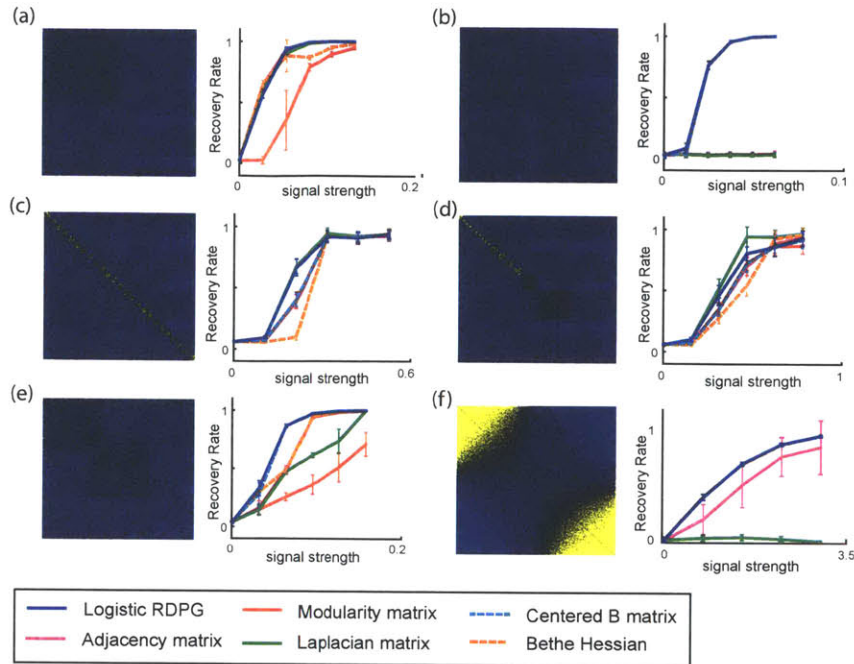


Figure 7-2: Performance comparison of our method (logistic RDPG) against spectral clustering methods in different clustering setups. Panels (a)-(e) illustrate networks that can be characterized by SBM, while panel (f) illustrates a non-SBM network model. The scale of the x axis is different in panel (f) than the rest of the panels. Our proposed method performs consistently well, while other methods exhibit sensitive and inconsistent performance in different network clustering setups.

we believe that its appropriate to do so). In contrast, our simple normalization allows the underlying latent position vectors to be accurately recovered in both cases. We also verified in panel (c) that our method successfully recovers latent positions for a non-SBM logistic RDPG. In this setup, the adjacency matrix’s first eigenvector again correlates with node degrees, and the modularity normalization causes an improvement. We found it remarkable that such a simple normalization (mean centering) enabled such significant improvements; using more sophisticated normalizations such as the Normalized Laplacian and the Bethe Hessian, no improvements over were observed (data not shown).

Second, we assessed the ability of our method to detect communities generated from the SBM. We compared against the following existing spectral network clustering methods:

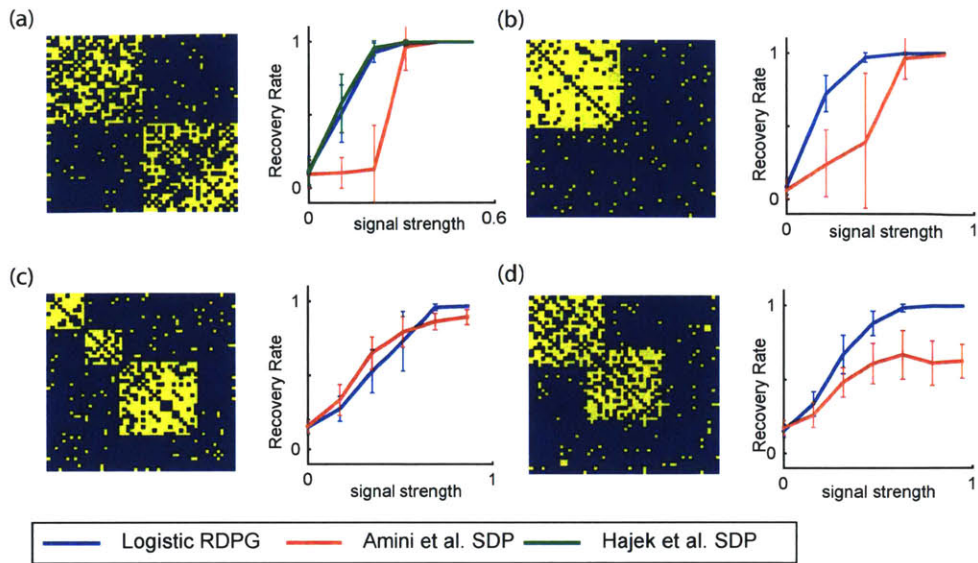


Figure 7-3: Performance comparison of our method (logistic RDPG) against semidefinite programming-based clustering methods.

- Modularity (Newman, 2006). We take the first d eigenvectors of the modularity matrix $M := A - d^T d / 2|\mathcal{E}|$, where d is the vector of node degrees and $|\mathcal{E}|$ is the number of edges in the network. We then perform k -means clustering on these eigenvectors.
- Normalized Laplacian (Chung, 1997). We take second- through $(d + 1)$ st- last eigenvectors of $L_{sym} := D^{-1/2}(D - A)D^{-1/2}$, where D is the diagonal matrix of degrees. We then perform k -means clustering on these eigenvectors.
- Bethe Hessian (Saade et al., 2014). We take the second- through $(d + 1)$ st- last eigenvectors of

$$H(r) := (r^2 - 1)\mathbf{1}_{n \times n} - rA + D,$$

where r^2 is the density of the graph as defined in [198].

- Unnormalized spectral clustering (Sussman et al., 2012). We take the first d eigenvectors of the adjacency matrix A , and perform k -means clustering on these eigenvectors.

- Spectral clustering on the mean-centered matrix B . We take the first d eigenvectors of the matrix B and perform k -means on them, without a scaling step.

Note that in our evaluation we include spectral clustering on the mean-centered adjacency matrix B without subsequent eigenvalue scaling of Algorithm 5 to demonstrate that the scaling step computed by logistic regression is essential to the performance of the proposed algorithm. When $d = 1$, the methods are equivalent. We also compare the performance of our method against two SDP-based approaches, the method proposed by Hajek et al. (2015) and the SDP-1 method proposed by Amini et al. (2014). For all methods we assume that the number of clusters k is given.

In our scoring metric, we distinguish between clusters and communities: For instance, in Figure 7-2-e, there are two clusters and four communities, comprised of nodes belonging only to cluster one, nodes belonging only to cluster two, nodes belonging to both clusters, and nodes belonging to neither. The score that we use is a normalized Jaccard index, defined as:

$$\frac{\max_{\sigma \in S_k} \sum_{l=1}^k \frac{|C_l \cap \hat{C}_{\sigma(l)}|}{|C_l|} - 1}{k - 1} \quad (7.31)$$

where C_l is the l -th community, \hat{C}_l is the l -th estimated community, and S_k is the group of permutations of k elements. Note that one advantage of using this scoring metric is that it weighs differently-sized clusters equally (it does not place higher weights on larger communities.).

Figure 7-2 presents a comparison between our proposed method and existing spectral methods in a wide range of clustering setups. Our proposed method performs consistently well, while other methods exhibit sensitive and inconsistent performance in different network clustering setups. For instance, in the case of two large clusters (b) the second-to-last eigenvector of the Normalized Laplacian fails to correlate with the community structure; in the case of having one dense cluster (a), the Modularity normalization performs poorly; when there are many small clusters (c), the performance of the Bethe Hessian method is poor. In each case, the proposed method performs at least as well as the best alternate method, except in the case of several

different-sized clusters (d), when the normalized Laplacian performs marginally better. In the case of overlapping clusters (e), our method performs significantly better than all competing methods. Spectral clustering on B without the scaling step also performs well in this setup; however, its performance is worse in panels (c-d) when d is larger, highlighting the importance of our logistic regression step.

The values of k and d for the different simulations were: $k = 2, d = 1$; $k = 2, d = 1$; $k = 25, d = 24$; $k = 17, d = 16$; $k = 4, d = 2$; $k = 2, d = 1$ for panels (a)-(f), respectively. The values of d are chosen based on the number of dimensions that would be informative to the community structure, if one knew the true latent positions. All networks have 1000 nodes, with background density 0.05.

While spectral methods are most prominent network clustering methods owing to their accuracy and efficiency, other approaches have been proposed, notably including SDP-based methods, which solve relaxations of the maximum likelihood problem over the SBM. We compare the performance of our method with the SDP-based methods proposed by Hajek et al. (2015) and Amini et al. (2014) (Figure 7-3). In the symmetric SBM, meaning the SBM with two equally-dense, equally-large communities, we find that our method performs almost equally well as the method of Hajek et al. (2015), which is a simple semidefinite relaxation of the likelihood in that particular case. Our method also performs better than the method of Amini et al., which solves a more complicated relaxation of the SBM maximum-likelihood problem in the more general case (Figure 7-3).

7.4 Proofs

Proof of Lemma 7.14. If not, then the optimal solution to Optimization (7.24) would be a better solution to Optimization (7.5) than X^* . \square

Proof of Lemma 7.15. The origin is in the feasible set for both optimizations. For each optimization, the objective function value satisfies $Tr(rCX) = rTr(CX)$. Thus, the optimum is either at the origin (if there is no positive solution) or at the boundary of the feasible set. If the optimum is at the origin, we have $s^* = \tilde{s} = 0$. If not, let X

be any solution to $\frac{1}{n^2} \sum_{i,j} f(X_{ij}) = h$. Let $r = \|X\|_F$, and let $r' = \sqrt{h/a_2}$. Claim: fixing γ , $r/r' \rightarrow 1$ uniformly as $h^* \rightarrow 0$.

Define

$$F_X(a) := \frac{1}{n^2} \sum_{i,j} f(aX_{ij}/\|X\|_F)$$

for $a > 0$. In addition, since r' has been defined such that the quadratic term of $F_X(r')$ is $a_2 \sum_{i,j} (\frac{r'}{r} X_{ij})^2 = h^*$, we have

$$F_X(r') = h^* + \frac{1}{n^2} O\left(\sum_{i,j} \left(\frac{r'}{r} X_{ij}\right)^3\right). \quad (7.32)$$

Moreover, the Taylor series for $f(\cdot)$ converges in a neighborhood of zero. Because of the constraint

$$\max_{i,j} X_{ij}^* = \max_i X_{ii}^* \leq \frac{\gamma}{n} \sum_i X_{ii}^*,$$

we can choose δ such that every entry X_{ij} falls within this neighborhood. This constraint also implies

$$\frac{1}{n^2} \sum X_{ij}^3 \leq \frac{1}{n^2} \sum |X_{ij}|^3 = O\left(\left(\frac{1}{n^2} \|X\|_F\right)^3\right).$$

Substituting this into (7.32), we have

$$F_X(r') = h^* + \frac{1}{n^2} O(r'^3). \quad (7.33)$$

Therefore, we have

$$\frac{|F_X(r) - F_X(r')|}{F_X(r')} = \frac{|h^* + O(r'^3) - h^*|}{h^*} = O(r'). \quad (7.34)$$

Note that $f(\cdot)$ is convex function with $f'(x) > 0$ for all $x > 0$ and $f'(x) < 0$ for all $x < 0$. Thus F_X is increasing, convex, and zero-valued at the origin: for any $a, b > 0$,

$$\frac{|a - b|}{b} < \frac{|F_X(a) - F_X(b)|}{F_X(b)}. \quad (7.35)$$

Thus $\frac{|r-r'|}{r'} = O(r')$ and $\frac{r}{r'} = 1 + O(r')$.

Let r_s be the norm of the $\arg \max$ to optimization (7.26); because the objective function is linear we have that $\tilde{s} \geq \frac{r'}{r_s} s^*$. Let r_t be the distance to the intersection of the boundary of the feasible set with the ray from the origin through the $\arg \max$ to optimization (7.27); then $s^* \geq \frac{r_t}{r'} \tilde{s}$. We have shown that both ratios tend uniformly to one. This completes the proof. \square

Proof of Lemma 7.16. First, suppose we have prior knowledge of eigenvalues of X^* . Denote its nonzero eigenvalues by $\lambda_1^*, \dots, \lambda_d^*$. Then we would be able to recover the optimal solution to Optimization (7.27) by solving the following optimization

$$\begin{aligned} \max_X \quad & Tr(BX) \\ & \lambda_i = \lambda_i^* \quad 1 \leq i \leq d \\ & rank(X) = d \end{aligned} \tag{7.36}$$

Note that the Frobenius norm of a matrix is determined by eigenvalues of the matrix as follows:

$$\|X\|_F^2 = Tr(XX^T) = Tr(X^2) = \sum \lambda_i^2. \tag{7.37}$$

Thus we can drop the Frobenius norm constraint in (7.36). Let X be an $n \times n$ psd matrix, whose non-null eigenvectors are the columns of a matrix $E \in \mathbb{R}^{n \times d}$, and whose respective eigenvalues are $\lambda_1, \dots, \lambda_d$. Let $V := E \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_d})$, so that $X = VV^T$. Rewrite the objective function as

$$Tr(BX) = Tr(V^T B V) = \sum_{i=1}^d \lambda_i e_i^T B e_i.$$

Therefore $\tilde{X} = EE^T$ and $X^* = VV^T$. \square

Proof of Lemma 7.17. The upper-triangular entries of A are independent Bernoulli random variables conditional on X and μ , with a logistic link function. The coeffi-

cients should be nonnegative, as X is constrained to be positive semidefinite. \square

Proof of Theorem 7.18. By Lemma 7.14, we have that the solution to optimization (7.24) is equal to the log-likelihood, up to addition of a constant. By Lemma 7.15, we have that for a fixed γ , as $h^* \rightarrow 0$, the quotient s^*/\tilde{s} converges uniformly to one, where s^* is the solution to (7.24) and \tilde{s} is the solution to optimization (7.27). The convergence is uniform over the choice of B that is needed for Theorem 7.18. Because s^* and \tilde{s} do not diverge to $\pm\infty$, this also implies that $s^* - \tilde{s}$, and therefore the log-likelihood ratio, converges uniformly to zero. By Lemma 7.16, the non-null eigenvectors of the arg max of optimization (7.27) are equivalent (up to rotation) to the first eigenvectors of B . Finally, by Lemma 7.17, the eigenvalues that maximize the likelihood can be recovered using a logistic regression step. By Lemma 7.15, the theorem would hold if we recovered the eigenvalues solving the approximate optimization (7.27). By finding the eigenvalues that exactly maximize the likelihood, we achieve a likelihood value at least as large. \square

Chapter 8

Network Biclustering

The term *biclustering* has been used to describe several distinct problems variants. In this chapter, we consider the problem of biclustering as a bipartite analogue of clustering: Given an $N \times M$ matrix, a bicluster is a subset of rows that are heavily connected to a subset of columns. In this framework, biclustering methods are data mining techniques allowing simultaneous clustering of the rows and columns of a matrix. We suppose there are two possible distributions for edge weights in the bipartite graph: a within-cluster distribution and a background distribution. Unlike in the traditional clustering problem, in our setup, biclusters may overlap, and a node may not belong to any cluster. We emphasize the distinction between biclustering and the bipartite analog of graph partitioning, which might be called *bipartitioning*.

Biclustering has several noteworthy applications. It has been used to find modules of coregulated genes using microarray gene expression data [290] and to predict tumor phenotypes from their genotypes [291]. It has been used for document classification, clustering both documents and related words simultaneously [292]. In all of these applications, biclusters are expected to overlap with each other, and these overlaps themselves are often of interest (e.g., if one wishes to explore the relationships between document topics).

The biclustering problem is NP-hard (see Proposition 8.2). However, owing to its practical importance, several heuristic methods using local search strategies have been developed. A popular approach is to search for one bicluster at a time by iteratively assigning rows to a bicluster based on the columns, and vice versa. Two algorithms based on this approach are ISA [293] and LAS [294]. Another approach is an exhaustive search for complete bicliques used by Bimax [295]. This approach fragments large noisy clusters into small complete ones. SAMBA [296] uses a heuristic combinatorial search for locally optimal biclusters, motivated by an exhaustive search algorithm that is exponential in the maximum degree of the nodes. FABIA [297] uses a variational approach to fit a model with fuzzy biclusters. For more details about existent biclustering algorithms, and performance comparisons, see references [295], [298] and [299]. Most of these methods have two shortcomings: first, they apply a local optimality criterion to each bicluster individually. Because a collection of locally

optimal biclusters might not be globally optimal, these local methods struggle to resolve overlapping clusters, which arise frequently in many applications. Second, the lack of a well-defined global objective function precludes an analytical characterization of their expected results.

Global optimization methods have been developed for problems closely related to biclustering, including clustering. Unlike most biclustering problem formulations, these are mostly partitioning problems: each node is assigned to one cluster or category. Major recent progress has been made in the development of spectral clustering methods (see references [300] and [282]) and message-passing algorithms (see [265], [6] and [301]). In particular, Affinity Propagation [6] maximizes the sum of similarities to one central exemplar instead of overall cluster density. Reference [302] uses variational expectation-maximization to fit the latent block model, which is a binary model in which each row or column is assigned to a row or column cluster, and the probability of an edge is dictated by the respective cluster memberships. Row and column clusters that are not paired to form biclusters.

In this chapter, we propose a message-passing algorithm that searches for a globally optimal collection of possibly overlapping biclusters. Our method maximizes a likelihood function using an approximation that separates a cluster-size penalty term into a row-count penalty and a column-count penalty. This decoupling enables the messages of the max-sum algorithm to be computed efficiently, effectively breaking an intractable optimization into a pair of tractable ones that can be solved in nearly linear time. When the underlying model parameters are unknown, they can be learned using an expectation-maximization approach.

Our approach has several advantages over existing biclustering algorithms: the objective function of our biclustering method has the flexibility to handle diverse statistical models; the max-sum algorithm is a more robust optimization strategy than commonly used iterative approaches; and in particular, our global optimization technique excels at resolving overlapping biclusters. In simulations, our method outperforms two of the best existing biclustering algorithms, ISA and LAS, when the planted clusters overlap. Applied to three gene expression datasets, our method

found biclusters of high quality in terms of cluster size and density.

8.1 Methods

8.1.1 Problem statement

Let $G = (V, W, E)$ be a weighted bipartite graph, with vertices $V = (1, \dots, N)$ and $W = (1, \dots, M)$, connected by edges with non-negative weights: $E : V \times W \rightarrow [0, \infty)$. Let $V_1, \dots, V_K \subset V$ and $W_1, \dots, W_K \subset W$. Let $(V_k, W_k) = \{(i, j) : i \in V_k, j \in W_k\}$ be a bicluster: Graph edge weights e_{ij} are drawn independently from either a *within-cluster distribution* or a *background distribution* depending on whether, for some k , $i \in V_k$ and $j \in W_k$. In this chapter, we assume that the within-cluster and background distributions are homogenous. However, our formulation can be extended to a general case in which the distributions are row- or column-dependent.

Let c_{ij}^k be the indicator for $i \in V_k$ and $j \in W_k$. Let $c_{ij} \triangleq \min(1, \sum_k c_{ij}^k)$ and let $\mathbf{c} \triangleq (c_{ij}^k)$.

Definition 8.1 (Biclustering Problem). Let $G = (V, W, E)$ be a bipartite graph with biclusters $(V_1, W_1), \dots, (V_K, W_K)$, within-cluster distribution f_1 and background distribution f_0 . The problem is to find the maximum likelihood cluster assignments (up to reordering):

$$\hat{\mathbf{c}} = \arg \max_{\mathbf{c}} \sum_{(i,j)} c_{ij} \log \frac{f_1(e_{ij})}{f_0(e_{ij})}, \quad (8.1)$$

$$c_{ij}^k = c_{rs}^k = 1 \Rightarrow c_{is}^k = c_{rj}^k = 1, \quad \forall i, r \in V, \forall j, s \in W.$$

Figure 8-1 demonstrates the problem qualitatively for an unweighted bipartite graph. In general, the combinatorial nature of a biclustering problem makes it computationally challenging.

Proposition 8.2. The clique problem can be reduced to the maximum likelihood problem of Definition (8.1). Thus, the biclustering problem is NP-hard.

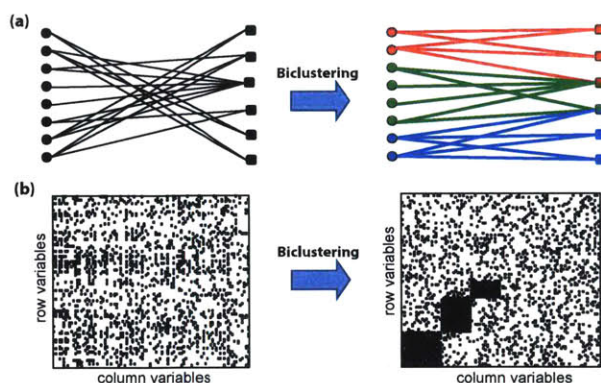


Figure 8-1: Biclustering is the analogue of clustering on a bipartite graph. (a) Biclustering allows nodes to be reordered in a manner that reveals modular structures in the bipartite graph. (b) The rows and columns of an adjacency matrix are similarly biclustered and reordered.

Proof. Proof is provided in Section 8.6.1. □

8.1.2 BCMP objective function

In this section, we introduce the global objective function considered in the proposed biclustering algorithm called *Biclustering using Message Passing* (BCMP). This objective function approximates the likelihood function of Definition 8.1. Let $l_{ij} = \log \frac{f_1(e_{ij})}{f_0(e_{ij})}$ be the log-likelihood ratio score of tuple (i, j) . Thus, the likelihood function of Definition 8.1 can be written as $\sum c_{ij} l_{ij}$. If there were no consistency constraints in the Optimization (8.1), an optimal maximum likelihood biclustering solution would be to set $c_{ij} = 1$ for all tuples with positive l_{ij} . Our key idea is to enforce the consistency constraints by introducing a cluster-size penalty function and shifting the log-likelihood ratios l_{ij} to recoup this penalty. Let N_k and M_k be the number of rows and columns, respectively, assigned to cluster k . We have,

$$\begin{aligned}
 \sum_{(i,j)} c_{ij} l_{ij} &\stackrel{(a)}{\approx} \sum_{(i,j)} c_{ij} \max(0, l_{ij} + \delta) - \delta \sum_{(i,j)} c_{ij} \\
 &\stackrel{(b)}{=} \sum_{(i,j)} c_{ij} \max(0, l_{ij} + \delta) + \delta \sum_{(i,j)} \max(0, -1 + \sum_k c_{ij}^k) - \delta \sum_k N_k M_k \\
 &\stackrel{(c)}{\approx} \sum_{(i,j)} c_{ij} \max(0, l_{ij} + \delta) + \delta \sum_{(i,j)} \max(0, -1 + \sum_k c_{ij}^k) - \frac{\delta}{2} \sum_k r_k N_k^2 + r_k^{-1} M_k^2.
 \end{aligned} \tag{8.2}$$

The approximation (a) holds when δ is large enough that thresholding l_{ij} at $-\delta$ has little effect on the resulting objective function. In equation (b), we have expressed the second term of (a) in terms of a cluster size penalty $-\delta N_k M_k$, and we have added back a term corresponding to the overlap between clusters. Because a cluster-size penalty function of the form $N_k M_k$ leads to an intractable optimization in the max-sum framework, we approximate it using a decoupling approximation (c) where r_k is a cluster shape parameter:

$$2N_k M_k \approx r_k N_k^2 + r_k^{-1} M_k^2, \tag{8.3}$$

when $r_k \approx M_k/N_k$. The cluster-shape parameter can be iteratively tuned to fit the estimated biclusters.

Following equation (8.2), the BCMP objective function can be separated into three terms as follows:

$$F(\mathbf{c}) = \sum_{i,j} \tau_{ij} + \sum_k \eta_k + \sum_k \mu_k, \tag{8.4}$$

$$\begin{cases} \tau_{ij} = \ell_{ij} \min(1, \sum_k c_{ij}^k) + \delta \max(0, \sum_k c_{ij}^k - 1) & \forall (i, j) \in V \times W, \\ \eta_k = -\frac{\delta}{2} r_k N_k^2 & \forall 1 \leq k \leq K, \\ \mu_k = -\frac{\delta}{2} r_k^{-1} M_k^2 & \forall 1 \leq k \leq K \end{cases} \quad (8.5)$$

Here τ_{ij} , the tuple function, encourages heavier edges of the bipartite graph to be clustered. Its second term compensates for the fact that when biclusters overlap, the cluster-size penalty functions double-count the overlapping regions. $\ell_{ij} \triangleq \max(0, l_{ij} - \delta)$ is the shifted log-likelihood ratio for observed edge weight e_{ij} . η_k and μ_k penalize the number of rows and columns of cluster k , N_k and M_k , respectively. Note that by introducing a penalty for each nonempty cluster, the number of clusters can be learned, and finding weak, spurious clusters can be avoided (see Section 8.2.3).

Now, we analyze BCMP over the following model for a binary or unweighted bipartite graph:

Definition 8.3. The binary biclustering model is a generative model for $N \times M$ bipartite graph (V, W, E) with K biclusters distributed by uniform sampling with replacement, allowing for overlapping clusters. Within a bicluster, edges are drawn independently with probability p , and outside of a bicluster, they are drawn independently with probability $q < p$.

In the following, we assume that p, q , and K are given. We discuss the case that the model parameters are unknown in Section 8.1.6. The following proposition shows that optimizing the BCMP objective function solves the problem of Definition 8.1 in the case of the binary model:

Proposition 8.4. Let (e_{ij}) be a matrix generated by the binary model described in Definition 8.3. Suppose p, q and K are given. Suppose the maximum likelihood assignment of edges to biclusters, $\arg \max(P(\text{data}|\mathbf{c}))$, is unique up to reordering. Let $r_k = M'_k/N'_k$ be the cluster shape ratio for the k -th maximum likelihood cluster. Then, by using these values of r_k , setting $\ell_{ij} = e_{ij}$, for all (i, j) , with cluster size

penalty

$$\frac{\delta}{2} = -\frac{\log\left(\frac{1-p}{1-q}\right)}{2\log\left(\frac{p(1-q)}{q(1-p)}\right)}, \quad (8.6)$$

we have,

$$\arg \max_{\mathbf{c}}(P(\text{data}|\mathbf{c})) = \arg \max_{\mathbf{c}}(F(\mathbf{c})). \quad (8.7)$$

Proof. The proof follows the derivation of equation (8.2). It is presented in Section 8.6.2. \square

Remark 8.5. In the special case when $q = 1 - p \in (0, 1/2)$, according to equation (8.6), we have $\frac{\delta}{2} = 1/4$. This is suggested as a reasonable initial value to choose when the true values of p and q are unknown; see Section 8.1.6 for a discussion of learning the model parameters.

The assumption that $r_k = N'_k/M'_k$ may seem rather strong. However, it is essential as it justifies the decoupling equation (8.3) that enables a linear-time algorithm. In practice, if the initial choice of r_k is close enough to the actual ratio that a cluster is detected corresponding to the real cluster, r_k can be tuned to find the true value by iteratively updating it to fit the estimated bicluster. This iterative strategy works well in our simulations. For more details about automatically tuning the parameter r_k , see Section 8.2.1.

In a more general statistical setting, log-likelihood ratios l_{ij} may be unbounded below, and the first step (a) of derivation (8.2) is an approximation; setting δ arbitrarily large will eventually lead to instability in the message updates.

8.1.3 The max-sum algorithm

In the following, we explain how message passing is used to optimize the BCMP objective function in approximately linear time $\mathcal{O}(K(N + \log M)(M + \log N))$. First, we briefly explain the max-sum algorithm that we use in our optimization. Consider

the following optimization:

$$\max_{X_1, \dots, X_n} \sum_{i=1}^m f_i(x_1, \dots, x_n) \quad (8.8)$$

where each variable X_i has alphabet \mathcal{X}_i . We assign a *function node* to each function f_i and a *variable node* to each variable. Messages $m_{X_i \rightarrow f_j}(x)$ and $m_{f_j \rightarrow X_i}(x)$ are vectors with length $|\mathcal{X}_i|$ defined as follows:

$$\begin{cases} m_{f_j \rightarrow X_i}(x) = \max_{X|X_i=x} (f_j(X_1, \dots, X_n) + \sum_{k \neq i} m_{X_k \rightarrow f_j}(X_k)), \\ m_{X_i \rightarrow f_j}(x) = \sum_{k \neq j} m_{f_k \rightarrow X_i}(x). \end{cases} \quad (8.9)$$

When messages converge after several iterations, an optimal value of variable X_i is then computed as,

$$x_i^* = \arg \max_{x \in \mathcal{X}_i} \sum_j m_{f_j \rightarrow X_i}(x). \quad (8.10)$$

If variables are binary (i.e., $|\mathcal{X}_i|=2$), scalar messages can be passed among nodes defined as follows,

$$\begin{cases} m_{f_j \rightarrow X_i} \triangleq m_{f_j \rightarrow X_i}(1) - m_{f_j \rightarrow X_i}(0), \\ m_{X_i \rightarrow f_j} \triangleq m_{X_i \rightarrow f_j}(1) - m_{X_i \rightarrow f_j}(0). \end{cases} \quad (8.11)$$

When the graphical model is a tree, the max-sum algorithm reduces to the Viterbi algorithm, which is exact. On graphical models with loops, it often obtains an approximately optimal solution in practice. It is the zero-temperature version of Belief Propagation (BP), also known as the sum-product algorithm or the cavity method. BP finds marginal probabilities in a graphical model given the joint probability distribution. To find a marginal probability, all possible configurations of the other variables must be considered, by summing instead of taking a maximum over the possible configurations. The fixed points of BP correspond to local minima of the Bethe free energy [303].

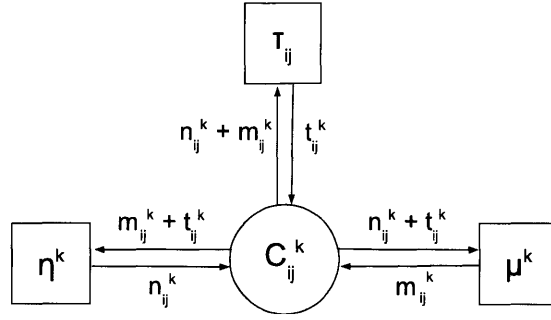


Figure 8-2: A graphical model of the BCMP objective function.

The max-sum algorithm was notably used in reference [6] for a clustering algorithm, Affinity Propagation, which inspired the proposed method.

8.1.4 Message updates

After solving the optimization of equation (8.15) in the main text, there are a few details involved in computing the actual message values. Let n be the value of N_1 that maximizes (8.15), and let y be the solution to equation (8.15).

(i) For $n_{12}^1(0)$, because the message sum is over $(i, j) \neq (1, 2)$, $m_{12}^1 + t_{12}^1$ is subtracted from s_i . This is only significant if row 1 was included in the arg max to (8.15) (i.e., if $s_1 \geq s_{(N+1-n)}$); in that case, it might no longer be optimal to include row i , in which case it may or may not be replaced by the $(n + 1)$ st largest row sum. If $l_1 = 0$, the message is

$$n_{12}^1(0) = y - \min(\max(0, m_{12}^1 + t_{12}^1), s_1 - s_{(N-n)}).$$

s_1 must be similarly adjusted for $n_{12}^1(1)$; however, the arg max is not affected:

$$n_{12}^1(1) = y - \max(0, m_{12}^1 + t_{12}^1).$$

(ii) For the message $n_{12}^1(1)$, row i must be included whether or not it is optimal in (8.15), thus it is possible that one fewer row to be included. If $s_1 < s_{(N+1-n)}$, then (i) is not relevant, and

$$n_{12}^1(1) = y + s_1 - \min\left(\frac{\delta}{2}(2n + 1), s_{(N+1-n)}\right)$$

(This also assumes that $l_{12} \geq 0$).

(iii) For calculating $n_{12}^1(0)$, if the initial penalty l_1 is used (see Section 8.2.3), it might be optimal to exclude all rows and set every $c_{12}^1 = 0$. If y' is the message value for $l_1 = 0$, then accounting for l_1 , the true value is $\max(0, y' - l_1)$. For $n_{12}^1(1)$, the $\arg \max$ is unaffected, and the message value is $y' - l_1$.

8.1.5 Biclustering Using Message Passing

In this section, we use the max-sum algorithm to optimize the objective function of equation (8.4). For a review of the max-sum message update rules, see Section 8.1.3. There are NM function nodes for the functions τ_{ij} , K function nodes for the functions η_k , and K function nodes for the functions μ_k . There are NMK binary variables, each attached to three function nodes: c_{ij}^k is attached to τ_{ij} , η_k , and μ_k (Figure 8-2). The incoming messages from these function nodes are named t_{ij}^k , n_{ij}^k , and m_{ij}^k , respectively. In the following, we describe messages for $c_{ij}^k = c_{12}^1$; other messages can be computed similarly.

First, we compute t_{12}^1 :

$$\begin{aligned}
 t_{12}^1(x) &\stackrel{(a)}{=} \max_{c_{12}^2, \dots, c_{12}^K} \left[\tau_{12}(x, c_{12}^2, \dots, c_{12}^K) + \sum_{k \neq 1} m_{12}^k(c_{12}^k) + n_{12}^k(c_{12}^k) \right] & (8.12) \\
 &\stackrel{(b)}{=} \max_{c_{12}^2, \dots, c_{12}^K} \left[\ell_{12} \min(1, \sum_k c_{12}^k) + \delta \max(0, \sum_k c_{12}^k - 1) + \sum_{k \neq 1} c_{12}^k (m_{12}^k + n_{12}^k) \right] + d_1
 \end{aligned}$$

where $d_1 = \sum_{k \neq 1} m_{12}^k(0) + n_{12}^k(0)$ is a constant. Equality (a) comes from the definition of messages according to equation (8.9). Equality (b) uses the definition of τ_{12} of equation (8.5) and the definition of the scalar message of equation (8.11). We can

further simplify t_{12} as follows:

$$\begin{cases} t_{12}^1(1) - d_1 \stackrel{(c)}{=} \ell_{12} + \sum_{k \neq 1} \max(0, \delta + m_{12}^k + n_{12}^k), \\ t_{12}^1(0) - d_1 \stackrel{(d)}{=} \ell_{12} - \delta + \sum_{k \neq 1} \max(0, \delta + m_{12}^k + n_{12}^k), & \text{if } \exists k, n_{12}^k + m_{12}^k + \delta > 0, \\ t_{12}^1(0) - d_1 \stackrel{(e)}{=} \max(0, \ell_{12} + \max_{k \neq 1} (m_{12}^k + n_{12}^k)), & \text{otherwise .} \end{cases} \quad (8.13)$$

If $c_{12}^1 = 1$, we have $\min(1, \sum_k c_{12}^k) = 1$, and $\max(0, \sum_k c_{12}^k - 1) = \sum_{k \neq 1} c_{12}^k$. These lead to equality (c). A similar argument can be made if $c_{12}^1 = 0$ but there exists a k such that $n_{12}^k + m_{12}^k + \delta > 0$. This leads to equality (d). If $c_{12}^1 = 0$ and there is no k such that $n_{12}^k + m_{12}^k + \delta > 0$, we compare the increase obtained by letting $c_{12}^k = 1$ (i.e., ℓ_{12}) with the penalty (i.e., $m_{12}^k + n_{12}^k$), for the best k . This leads to equality (e).

Remark 8.6. Computation of $t_{ij}^1, \dots, t_{ij}^k$ using equality (d) costs $\mathcal{O}(K)$, and not $\mathcal{O}(K^2)$, as the summation need only be computed once.

Messages m_{12}^1 and n_{12}^1 are computed as follows:

$$\begin{cases} m_{12}^1(x) = \max_{\mathbf{c}^1 | c_{12}^1 = x} [\mu_1(\mathbf{c}^1) + \sum_{(i,j) \neq (1,2)} t_{ij}^1(c_{ij}^1) + n_{ij}^1(c_{ij}^1)], \\ n_{12}^1(x) = \max_{\mathbf{c}^1 | c_{12}^1 = x} [\eta_1(\mathbf{c}^1) + \sum_{(i,j) \neq (1,2)} t_{ij}^1(c_{ij}^1) + m_{ij}^1(c_{ij}^1)], \end{cases} \quad (8.14)$$

where $\mathbf{c}^1 = \{c_{ij}^1 : i \in V, j \in W\}$. To compute n_{12}^1 in constant time, we perform a preliminary optimization, ignoring the effect of edge (1, 2):

$$\arg \max_{\mathbf{c}^1} -\frac{\delta}{2} N_1^2 + \sum_{(i,j)} t_{ij}^1(c_{ij}^1) + m_{ij}^1(c_{ij}^1). \quad (8.15)$$

Let $s_i = \sum_{j=1}^M \max(0, m_{ij}^1 + t_{ij}^1)$ be the sum of positive incoming messages of row i . The function η_1 penalizes the number of rows containing some nonzero c_{ij}^1 : if any message along that row is included, there is no additional penalty for including every positive message along that row. Thus, optimization (8.15) is computed by deciding which rows to include. This can be done efficiently through sorting: we sort row sums $s_{(1)}, \dots, s_{(N)}$ at a cost of $\mathcal{O}(N \log N)$. Then we proceed from largest to smallest,

including row $(N + 1 - i)$ if the marginal penalty $\frac{\delta}{2}(i^2 - (i - 1)^2) = \frac{\delta}{2}(2i - 1)$ is less than $s_{(N+1-i)}$. After solving optimization (8.15), the messages $n_{12}^1, \dots, n_{N2}^1$ can be computed in linear time, as we explain in Section 8.1.4.

Remark 8.7. Computation of n_{ij}^k through sorting costs $\mathcal{O}(N \log N)$.

Proposition 8.8 (Computational Complexity of BCMP). The computational complexity of BCMP over a bipartite graph with N rows, M columns, and K clusters is $\mathcal{O}(K(N + \log M)(M + \log N))$.

Proof. For each iteration, there are NM messages t_{ij} to be computed at cost $\mathcal{O}(K)$ each. Before computing (n_{ij}^k) , there are K sorting steps at a cost of $\mathcal{O}(M \log M)$, after which each message may be computed in constant time. Likewise, there are K sorting steps at a cost of $\mathcal{O}(N \log N)$ each before computing (m_{ij}^k) . \square

We provide an empirical runtime example of the algorithm in Figure 8-5.

8.1.6 Learning Model Parameters using Expectation-Maximization

In the BCMP objective function described in Section 8.1.2, the parameters of the generative model were used to compute the log-likelihood ratios (l_{ij}) . In practice, however, these parameters may be unknown. Expectation-Maximization (EM) can be used to estimate these parameters. The use of EM in this setting is slightly unorthodox, as we estimate the hidden labels (cluster assignments) in the M step instead of the E step. However, the distinction between parameters and labels is not intrinsic in the definition of EM [304] and the true ML solution is still guaranteed to be a fixed point of the iterative process. Note that it is possible that the EM iterative procedure leads to a locally optimal solution and therefore it is recommended to use several random re-initializations for the method.

The EM algorithm has three steps:

- **Initialization:** We choose initial values for the underlying model parameters θ and compute the log-likelihood ratios (l_{ij}) based on these values, denoting by F_0 the initial objective function.

- **M step:** We run BCMP to maximize the objective $F_i(\mathbf{c})$. We denote the estimated cluster assignments by $\hat{\mathbf{c}}_i$.
- **E step:** We compute the expected-log-likelihood function as follows:

$$F_{i+1}(\mathbf{c}) = E_{\theta}[\log P((e_{ij})|\theta)|\mathbf{c} = \hat{\mathbf{c}}_i] = \sum_{(i,j)} E_{\theta}[\log P(e_{ij}|\theta)|\mathbf{c} = \hat{\mathbf{c}}_i]. \quad (8.16)$$

Conveniently, the expected-likelihood function takes the same form as the original likelihood function, with an input matrix of *expected* log-likelihood ratios. These can be computed efficiently if conjugate priors are available for the parameters. Therefore, BCMP can be used to maximize F_{i+1} . The algorithm terminates upon failure to improve the estimated likelihood $F_i(\hat{\mathbf{c}}_i)$.

Below we describe how expectation-maximization (EM) was applied to the binary and Gaussian models used in our simulations. In the case of the binary model, we use uniform Beta distributions as conjugate priors for p and q , and in the case of the Gaussian model, we use inverse-gamma-normal distributions as the priors for the variances and means. Even when convenient priors are not available, EM is still tractable as long as one can sample from the posterior distributions.

In the case of a binary model with parameters p and q , the prior on p is $Unif(0, 1)$. If n tuples were placed in a cluster, and m out of those n had edges, the posterior distribution for p is $Beta(1 + m, 1 + n - m)$, and the posterior for q is computed similarly. The posteriors are independent of each other because the edge weights are assumed to be drawn independently (given the cluster assignments). There are only two integrals to be computed:

$$l_{ij} = \begin{cases} \int_{[0,1]^2} \log \frac{p}{q} dP(p)dP(q) & e_{ij} = 1 \\ \int_{[0,1]^2} \log \frac{1-p}{1-q} dP(p)dP(q) & e_{ij} = 0 \end{cases} \quad (8.17)$$

Here we use $P(p)$ and $P(q)$ to denote the posterior CDF of these parameters. These integrals can be computed numerically by sampling from the posterior distri-

butions.

In the case of a Gaussian model with two different means and variances, we use a normal-inverse gamma distribution on the mean and variance. Assume that $e_{ij}|c_{ij} = a$ is normally distributed with mean μ_a and variance σ_a^2 . Let the prior distribution both within clusters and outside of them be $(\mu_a, \sigma_a^2) \sim NIG(\mu, \lambda, \alpha, \beta)$, that is,

$$\begin{aligned}\sigma_a^2 &\sim \Gamma(\alpha, \beta), \\ \mu_a|\sigma_a^2 &\sim N(\mu, \sigma_a^2/\lambda).\end{aligned}\tag{8.18}$$

Let $\hat{\mathbf{c}} = (\hat{c}_{ij})$ be the current estimate of the cluster memberships. Let $A \triangleq \{(i, j) : \hat{c}_{ij} = 1\}$ and $B \triangleq \{(i, j) : \hat{c}_{ij} = 0\}$; let $\hat{\mu}_A, \hat{\sigma}_A^2, \hat{\mu}_B, \hat{\sigma}_B^2$ be the sample means and variances of the edge weights over the estimates \hat{A}, \hat{B} . Then the posterior distributions of the model parameters are

$$\begin{aligned}(\mu_1, \sigma_1^2)|\hat{A} &\sim NIG\left(\frac{\mu\lambda^{-1} + |A|\hat{\mu}_A}{\lambda^{-1} + |\hat{A}|}, \lambda^{-1} + |\hat{A}|, \alpha + \frac{|\hat{A}|}{2}, \beta + \frac{1}{2}\hat{\sigma}_A^2\right) \\ (\mu_0, \sigma_0^2)|\hat{B} &\sim NIG\left(\frac{\mu\lambda^{-1} + |B|\hat{\mu}_B}{\lambda^{-1} + |\hat{B}|}, \lambda^{-1} + |\hat{B}|, \alpha + \frac{|\hat{B}|}{2}, \beta + \frac{1}{2}\hat{\sigma}_B^2\right)\end{aligned}\tag{8.19}$$

Let ϕ be the standard Normal density function. The expression for the new input matrix is:

$$l_{ij} = \int_{\mathbb{R}^2} \log \sigma_1^{-1} \phi\left(\frac{e_{ij} - \mu_1}{\sigma_1}\right) dP(\mu_1, \sigma_1) - \int_{\mathbb{R}^2} \log \sigma_0^{-1} \phi\left(\frac{e_{ij} - \mu_0}{\sigma_0}\right) dP(\mu_0, \sigma_0)\tag{8.20}$$

This computation can be performed efficiently by drawing samples from the two distributions and reusing them for each of the NM integrations. Typically, in our simulations, approximately three iterations of EM were needed for convergence.

8.2 Automatic Parameter Tuning

8.2.1 Automatically tuning r_k

Automatically tuning r_k to fit the observed biclusters allows BCMP to find differently-shaped clusters. After each message update, it updates $r_{new} = \sqrt{M_k/N_k} \cdot \sqrt{r_{old}}$.

This strategy works well in practice. It can err by finding two biclusters with shapes closer to the original choice of r_k instead of a single, differently shaped cluster; to catch such errors, it is helpful to visualize the biclustering results or to check that no two clusters have mostly the same rows or columns.

8.2.2 Automatically tuning δ

The offset parameter δ is tuned by maximizing the unthresholded likelihood function. BCMP is run for several values of δ , and for each solution, the likelihood is computed. The solution with the greatest likelihood is kept.

8.2.3 Automatically detecting the number of clusters

When the true number of clusters is unknown, a penalty for each nonempty cluster can be used to find fewer than K clusters. A term $-l_1 \max(0, N_k - 1)$ is added to the η_k function, and $-l_1 \max(0, M_k - 1)$ is added to the μ_k function. This penalty can be interpreted as a significance threshold, as clusters with barely-positive scores might be found even in a dataset with no true clusters; or, it may be interpreted as a prior on the true number of clusters, which is $Geom(\exp(-2l_1))$ as $K \rightarrow \infty$.

It does not change the computational complexity to use these penalties. When computing $n_{ij}^k(1)$, there is no additional computation as the penalty l_1 is incurred automatically. When computing $n_{ij}^k(0)$, if y is the original message before accounting for l_1 , the new message is $\max(0, y - l_1)$.

In order to set this parameter, one strategy is: randomly permute the entries of the input matrix (l_{ij}) ; run BCMP on the scrambled data matrix; and choose $l_1 = \frac{1}{2} \max_k \sum_{i,j} c_{ij}^k l_{ij}$.

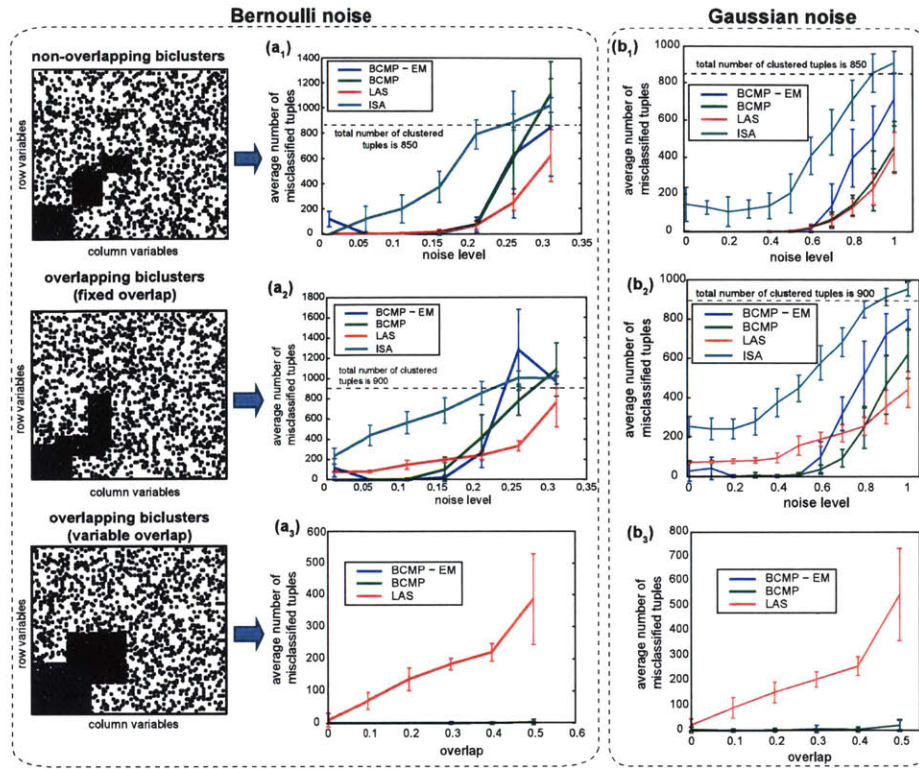


Figure 8-3: Performance comparison of the proposed method (BCMP) with ISA and LAS, for Bernoulli and Gaussian models, and for overlapping and non-overlapping biclusters. On the y axis is the total number of misclassified row-column pairs. Either the noise level or the amount of overlap is on the x axis.

8.3 Evaluation Results

We compared the performance of our biclustering algorithm with two methods, ISA and LAS, in simulations and in real gene expression datasets (Section 8.4). ISA was chosen because it performed well in comparison studies [295] [298], and LAS was chosen because it outperformed ISA in preliminary simulations. Both ISA and LAS search for biclusters using iterative refinement. ISA assigns rows iteratively to clusters fractionally in proportion to the sum of their entries over columns. It repeats the same for column-cluster assignments, and this process is iterated until convergence. LAS uses a similar greedy iterative search without fractional memberships, and it masks already-detected clusters by mean subtraction.

In our simulations, we generate simulated bipartite graphs of size 100x100. We

planted (possibly overlapping) biclusters as full blocks with two noise models:

- *Bernoulli noise*: we drew edges according to the binary model of Definition 8.3 with varying noise level $q = 1 - p$.
- *Gaussian noise*: we drew edge weights within and outside of biclusters from normal distributions $N(1, \sigma^2)$ and $N(0, \sigma^2)$, respectively, for different values of σ .

For each of these cases, we ran simulations on three setups (see Figure 8-3):

- *Non-overlapping clusters*: three non-overlapping biclusters were planted in a 100×100 matrix with sizes 20×20 , 15×20 , and 15×10 . We varied the noise level.
- *Overlapping clusters with fixed overlap*: Three overlapping biclusters with fixed overlaps were planted in a 100×100 matrix with sizes 20×20 , 20×10 , and 10×30 . We varied the noise level.
- *Overlapping clusters with variable overlap*: we planted two 30×30 biclusters in a 100×100 matrix with variable amount of overlap between them, where the amount of overlap is defined as the fraction of rows and columns shared between the two clusters. We used Bernoulli noise level $q = 1 - p = 0.15$, and Gaussian noise level $\sigma = 0.7$.

The methods used have some parameters to set. Here are the parameters that we used to run each method:

- *BCMP method with underlying parameters given*: We computed the input matrix of shifted log-likelihood ratios following the discussion in Section 8.1.2. The number of biclusters K was given. We initialized the cluster-shape parameters r_k at 1 and updated them as discussed in Section 8.2.1. In the case of Bernoulli noise, following Proposition 8.4 and Remark 8.5, we set $\ell_{ij} = e_{ij}$ and $\frac{\delta}{2} = 1/4$. In the case of Gaussian noise, we chose a threshold δ to maximize the unthresholded likelihood (see Section 8.2.2).

- *BCMP - EM method*: Instead of taking the underlying model parameters as given, we estimated them using the procedure described in Section 8.1.6. We used identical, uninformative priors on the parameters of the within-cluster and null distributions.
- *ISA method*: We used the same threshold ranges for both rows and columns, attempting to find best-performing threshold values for each noise level. These values were mostly around 1.5 for both noise types and for all three dataset types. We found positive biclusters, and used 20 reinitializations. Out of these 20 runs, we selected the best-performing run.
- *LAS method*: There were no parameters to set. Since K was given, we selected the first K biclusters discovered by LAS, which marginally increased its performance.

Evaluation results of both noise models and non-overlapping and overlapping biclusters are shown in Figure 8-3. In the non-overlapping case, BCMP and LAS performed similarly well, better than ISA. Both of these methods made few or no errors up until noise levels $q = 0.2$ and $\sigma = .6$ in Bernoulli and Gaussian cases, respectively. When the parameters had to be estimated using EM, BCMP performed worse for higher levels of Gaussian noise but well otherwise. ISA outperformed BCMP and LAS at very high levels of Bernoulli noise; at such a high noise level, however, the results of all three algorithms are comparable to a random guess.

In the presence of overlap between biclusters, BCMP outperformed both ISA and LAS except at very high noise levels. Whereas LAS and ISA struggled to resolve these clusters even in the absence of noise, BCMP made few or no errors up until noise levels $q = 0.2$ and $\sigma = .6$ in Bernoulli and Gaussian cases, respectively. Notably, the overlapping clusters were more asymmetrical, demonstrating the robustness of the strategy of iteratively tuning r_k in our method. In simulations with variable overlaps between biclusters, for both noise models, BCMP outperformed LAS significantly, while the results for the ISA method were very poor (data not shown). These results demonstrate that BCMP excels at inferring overlapping biclusters.

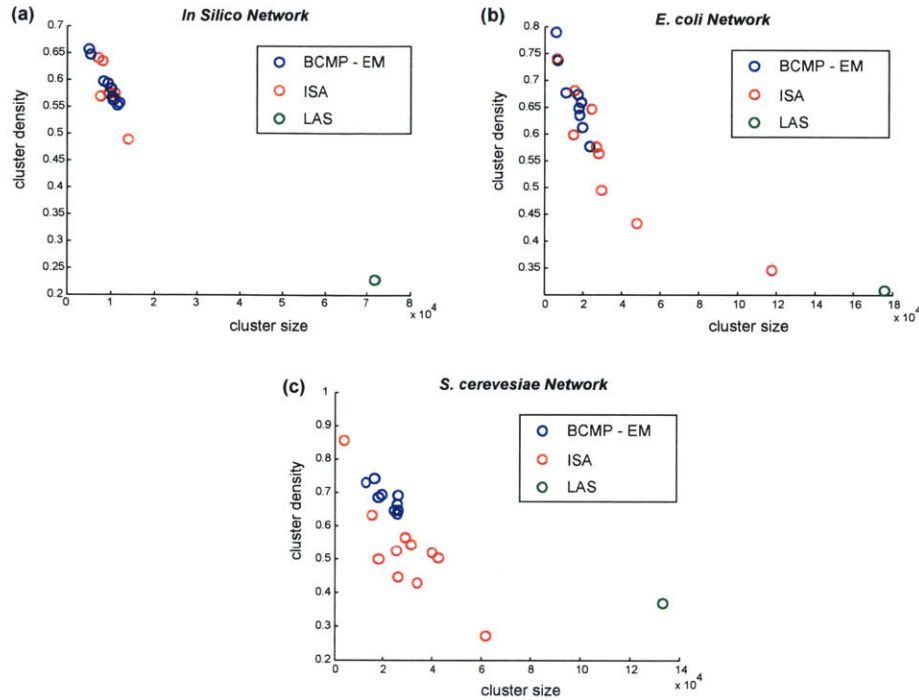


Figure 8-4: Total size and average density of biclusters reported by BCMP, ISA and LAS for three gene expression datasets. a) *In Silico*. b) *E. coli*. c) *S. cerevisiae*.

8.4 Gene Expression Data

We applied BCMP, ISA and LAS to three DREAM5 gene expression datasets: *In Silico*, *E. coli*, and *S. cerevisiae* [3]. We binarized these datasets, placing ones where the gene expression level was at least two standard deviations from its mean and zeroes elsewhere. We evaluated the three biclustering algorithms in terms of the total size and average density of the reported clusters; when reported clusters overlapped, these regions were not double-counted. We used BCMP to fit the stochastic block model, using different initial parameter settings for EM (see Section 8.1.6) and $K = 10$ (this is the number of clusters reported by LAS by default). ISA was run with different threshold-parameter settings, keeping the first ten clusters reported. LAS had no parameters to set. We found that the density of the clusters reported by LAS was diluted by a few large, low-density clusters. ISA and BCMP had similar results for the *In Silico* and *E. coli* datasets, and BCMP reported denser clusters on the *S. cerevisiae* datasets (Figure 8-4).

8.5 Document Classification Example

Biclustering has been used for document clustering [292]. Here we illustrate of how BCMP might be applied to such a problem.

Let $D = d_1, \dots, d_N$ be a collection of documents containing the words $W = w_1, \dots, w_M$, and let $(D_1, W_1), \dots, (D_K, W_K)$ be a hidden set of topics (each word and document can belong to any number of topics). Document d_i contains n_{ij} instances of word w_j out of $n_i = \sum_j n_{ij}$ total words. Assume that n_i is large enough to use a Poisson model: $n_{ij} \sim \text{Poisson}(\lambda_{ij})$ where $\lambda_{ij} = n_i r_j$ if document d_i concerns a topic relating to word w_j and $\lambda_{ij} = n_i s_j$ otherwise. Assume that K is given.

The problem is to recover the topics. It can be solved by estimating the parameters r_{ij}, s_{ij} and biclustering the bipartite graph $(D, W, (n_{ij}))$.

We can obtain initial estimates for r_j and s_j using k-means on n_{1j}, \dots, n_{Nj} . This is effectively a clustering problem along one axis with two clusters:

$$\max_{r_j, s_j, c_{1j}, \dots, c_{Nj}} P(n_{1j}, \dots, n_{Nj} | r_j, s_j, c_{1j}, \dots, c_{Nj})$$

where c_{ij} is the indicator variable for document i concerning some topic relating to word j . Words with no evidence for $r_j \neq s_j$ at some significance threshold can be discarded. After performing biclustering, these estimates can be improved using the newly estimated cluster memberships.

After estimating r_j and s_j , the likelihood ratios are

$$\frac{p_{ij}}{q_{ij}} = \frac{P(n_{ij} | \lambda_{ij} = n_i r_j)}{P(n_{ij} | \lambda_{ij} = n_i s_j)} = \frac{(n_i r_j)^{n_{ij}} e^{-n_i r_j} / n_{ij}!}{(n_i s_j)^{n_{ij}} e^{-n_i s_j} / n_{ij}!}$$

and the input matrix for BCMP is

$$\ell_{ij} = \max(0, \log \frac{p_{ij}}{q_{ij}} + \delta) = \max(0, n_{ij} \log \left(\frac{r_j}{s_j} \right) + n_i (s_j - r_j) + \delta). \quad (8.21)$$

EM can be used to iteratively update the input matrix and improve the parameter estimates as discussed in Section 8.1.6.

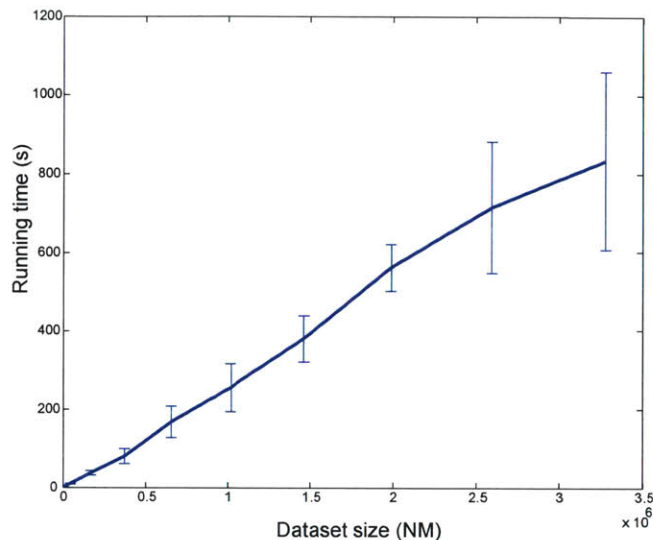


Figure 8-5: Running time for BCMP. In Proposition 8.8, we showed that the complexity of a round of message updates for BCMP is nearly linear, i.e., $\mathcal{O}(K(N + \log M)(M + \log N))$. This figure illustrates the running time of BCMP on a personal computer, with $N = M$.

8.6 Proofs

8.6.1 Proof of Proposition 1

In the case $K = 1$ the problem (8.1) is equivalent to the problem of finding a maximum weight biclique on the graph $G = (U, V, (l_{ij}))$ with edge weights $l_{ij} = \log \frac{f_1(e_{ij})}{f_0(e_{ij})}$. This problem contains as a special case the problem of finding a maximum size complete clique on a graph (V, E) : let $U = V$ and $l_{ii} = 1$. For $i \neq j$, let $l_{ij} = l_{ji} = 0$ if $(i, j) \in E$ and $-\infty$ otherwise. Suppose the maximal clique on (V, E) is W . Then the biclique (W, W) has weight $|W|$. For every other biclique (U', V') , if its weight is not $-\infty$, then $U' \cap V'$ must correspond to a complete clique on (V, E) ; its weight is $|U' \cap V'| \leq |W|$, thus W is a solution to (8.1) on G , and the problem (8.1) is NP-hard.

8.6.2 Proof of Proposition 2

This proof is essentially a pedantic version of the derivation of equation (8.2) in the main text. Given a putative set of biclusters $(V_1, W_1), \dots, (V_K, W_k)$, let $\mathbf{c} = (c_{ij})$ be the $N \times M$ matrix with ones where i, j are in the same bicluster and zeroes elsewhere. If the data is $E = (e_{ij})$, with ones where edges are observed, then the probability of the data given these biclusters is just the product of the probabilities of the tuples,

$$P(E|I) = \prod_{i \leq N, j \leq M} t_{ij}$$

where,

$$t_{ij} = \begin{cases} p & e_{ij} = c_{ij} = 1 \\ q & e_{ij} = 1, c_{ij} = 0 \\ 1 - p & e_{ij} = 0, c_{ij} = 1 \\ 1 - q & e_{ij} = c_{ij} = 0 \end{cases} \quad (8.22)$$

The maximum likelihood set of biclusters maximize this expression. Equivalently, we maximize the sum, over all tuples (i, j) placed in a bicluster, of the log-likelihood ratios, $\log(p/q)$ if $e_{ij} = 1$ and $\log((1-p)/(1-q))$ otherwise.

Let the input matrix be $(\ell_{ij}) = (e_{ij})$, and let the cluster-size penalty coefficient be

$$\frac{\delta}{2} = -\frac{\log(\frac{1-p}{1-q})}{2 \log(\frac{p(1-q)}{q(1-p)})}. \quad (8.23)$$

Let c_{ij} be the indicator for row-column pair ij being in cluster k for some k . First,

$$\begin{aligned} F(\mathbf{c}) &= \sum_{ij} c_{ij}(\ell_{ij} + (\sum_k c_{ij}^k - 1) \cdot \delta) - \frac{\delta}{2} \sum_k r_k N_k^2 + r_k^{-1} M_k^2 \\ &= \sum c_{ij}(\ell_{ij} + (\sum_k c_{ij}^k - 1) \cdot \delta) - \delta [\sum_k N_k M_k + \frac{1}{2}(r_k N_k - r_k^{-1} M_k)^2]. \end{aligned}$$

Now, if each $r_k = M_k/N_k$, then the terms $\frac{\delta}{2}(rN_k - r^{-1}M_k)^2$ drop out; otherwise, there is an additional penalty. The additional penalty cannot hurt, since it is zero when $\mathbf{c} = \mathbf{c}'$ and nonpositive otherwise. Notice that if a tuple is assigned to more than one cluster, the term $(\sum_k c_{ij}^k - 1)\frac{\delta}{2}$ cancels the fact that it is counted in two cluster-size penalty functions. By dropping the cluster shape term we get

$$\begin{aligned} F(\mathbf{c}) &\geq \sum c_{ij}(\ell_{ij} + (\sum_k c_{ij}^k - 1) \cdot \delta) - \delta \sum_k N_k M_k \\ &= \sum_{ij} c_{ij}(\ell_{ij} + \delta) \end{aligned} \tag{8.24}$$

$$\begin{aligned} &\propto |\{(i, j) : \text{tuple } ij \text{ is placed in a bicluster and } e_{ij} = 1\}| \cdot \log\left(\frac{p}{q}\right) \\ &\quad + |\{(i, j) : \text{tuple } ij \text{ is placed in a bicluster and } e_{ij} = 0\}| \cdot \log\left(\frac{1-p}{1-q}\right) \\ &= \log(P(\text{data}|\mathbf{c})) + \text{const} \end{aligned} \tag{8.25}$$

with equality when each $r_k = N_k/M_k$. Equality (8.24) follows from the fact that $\sum_k N_k M_k$ is the number of clustered tuples plus the amount of overlap between the clusters, and this overlap term cancels the term $\sum_{ij} c_{ij}(\sum_k c_{ij}^k - 1)$. Step (8.25) follows from the choice of ℓ_{ij} . This completes the proof of 2.3.

Chapter 9

Compressive Sensing over Networks

Data compression has been a research topic of many scholars in past years. However, recently, the field of compressive sensing, originated in [305], [306] and [307], looked at the compression problem from another point of view. In this chapter, we try to make a bridge between compressive sensing which may be viewed as a signal processing technique and both source coding and channel coding. By using this connection, we propose some non-trivial practical coding schemes that provide a trade-off between the rate and the decoding complexity.

Compressive sensing has provided a low complexity approximation to the signal reconstruction. Information theoretic has been mostly concerned with accuracy of the signal reconstruction under rate constraints. Here, we seek to provide new connections which use compressive sensing for traditional information theory problems such as Slepian-Wolf compression and channel coding in order to provide mechanisms to explicitly trade-off between the decoding complexity and the rate.

Some previous work investigated the connection between compressive sensing and information theory. For example, reference [308] studied the minimum number of noisy measurements required to recover a sparse signal by using Shannon information theory bounds. Reference [309] investigated the contained information in noisy measurements by viewing the measurement system as an information theoretic channel and using the rate distortion function. Also, reference [310] studied the trade-offs between the number of measurements, the signal sparsity level, and the measurement noise level for exact support recovery of sparse signals by using an analogy between support recovery and communication over the Gaussian multiple access channel.

In this chapter, we want to investigate applications of compressive sensing on source coding and channel coding. First, in Section 9.2, we consider distributed compression of correlated sources. In this section, by using compressive sensing, we propose a compression scheme for a Slepian-Wolf problem with a family of correlated sources. The proposed decoder is a modularized decoder, providing a trade-off between the compression rate and the decoding complexity. We call this scheme *Sparse Distributed Compression*. Then, We use this compression scheme for a general multicast network with correlated sources. Here, we first decode some of the sources by

a network decoding technique and then, we use the compressive sensing decoder to obtain the whole sources.

Then, in Section 9.3, we investigate applications of compressive sensing on channel coding. We propose a coding scheme that combines compressive sensing and random channel coding for a high-SNR point-to-point Gaussian channel. We call this scheme *Sparse Channel Coding*. We propose a modularized decoder, providing a trade-off between the capacity loss and the decoding complexity (i.e., the higher the capacity loss, the lower the decoding complexity.). The idea is to add intentionally some correlation to transmitted signals in order to decrease the decoding complexity. At the receiver side, first, we use a compressive sensing decoder on a noisy signal to obtain a noisy estimate of the original signal and then, we apply a traditional channel coding decoder to find the original signal.

9.1 Technical Background and prior work

In this section, we review some prior results, in both compressive sensing and distributed compression, which will be used in the rest of the chapter.

9.1.1 Compressive Sensing Background

Let $\mathbf{X} \in \mathbb{R}^n$ be a signal vector. We say this signal is k -sparse if k of its coordinates are non-zero and the rest are zero. Define $\alpha = \frac{k}{n}$ as the sparsity ratio of the signal. Suppose Φ is a $n \times n$ measurement matrix (hence, $\mathbf{Y} = \Phi\mathbf{X}$ is the measurement vector.). If Φ is non-singular, by having \mathbf{Y} and Φ , one can recover \mathbf{X} . However, [306] and [307] showed that, one can recover the original sparse signal by having far fewer measurements than n . Suppose for a given n and k , m satisfies the following:

$$m \geq \rho k \log(n/k) \tag{9.1}$$

where ρ is a constant. Let \mathcal{S} be the power set of $\{1, 2, \dots, n\}$, and \mathcal{S}_m represent a member of \mathcal{S} with cardinality m (e.g., $\mathcal{S}_m = \{1, 2, \dots, m\}$). A $m \times n$ matrix $\Phi_{\mathcal{S}_m}$

is formed of rows of Φ whose indices belong to \mathcal{S}_m . $\Phi_{\mathcal{S}_m}$ represents an incomplete measurement matrix.

Reference [311] showed that, if m satisfies (9.1) and $\Phi_{\mathcal{S}_m}$ satisfies the *Restricted Isometry Property* (which will be explained later), having these m measurements is sufficient to recover the original sparse vector. In fact, it has been shown in [311] that, \mathbf{X} is the solution of the following linear programming:

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_{\mathcal{L}_1} \\ \text{subject to} \quad & \mathbf{Y} = \Phi_{\mathcal{S}_m} \mathbf{X}, \end{aligned} \tag{9.2}$$

where $\|\cdot\|_{\mathcal{L}_1}$ represents the \mathcal{L}_1 norm of a vector. Let us refer to the complexity of this optimization by $CX_{cs}^{noiseless}(m, n)$. We say $\Phi_{\mathcal{S}_m}$ satisfies the restricted isometry property (RIP) iff there exists $0 < \delta_k < 1$ such that, for any vector $\mathbf{V} \in \mathbb{R}^n$, we have:

$$(1 - \delta_k) \|\mathbf{V}\|_{\mathcal{L}_2}^2 \leq \|\Phi_{\mathcal{S}_m} \mathbf{V}\|_{\mathcal{L}_2}^2 \leq (1 + \delta_k) \|\mathbf{V}\|_{\mathcal{L}_2}^2. \tag{9.3}$$

Reference [311] showed if,

$$\delta_k + \delta_{2k} + \delta_{3k} < 1, \tag{9.4}$$

then, the optimization (9.2) can recover the original k -sparse signal. This condition leads to (9.1). Intuitively, if $\Phi_{\mathcal{S}_m}$ satisfies the RIP condition, it preserves the Euclidean length of k sparse signals. Hence, these sparse signals cannot be in the null-space of $\Phi_{\mathcal{S}_m}$. Moreover, any k subset of columns of $\Phi_{\mathcal{S}_m}$ are nearly orthogonal if $\Phi_{\mathcal{S}_m}$ satisfies RIP. Most of matrices which satisfy RIP are random matrices. Reference [312] showed a connection between compressive sensing, n -widths, and the Johnson-Lindenstrauss Lemma. Through this connection, [313] provided an elegant way to find δ_k by using the concentration theorem for random matrices and proposed some data-base friendly matrices (binary matrices) satisfying RIP.

Now, suppose our measurements are noisy (i.e., $\mathbf{Y} = \Phi_{\mathcal{S}_m} \mathbf{X} + \mathbf{Z}$, where $\|\mathbf{Z}\|_{\mathcal{L}_2} =$

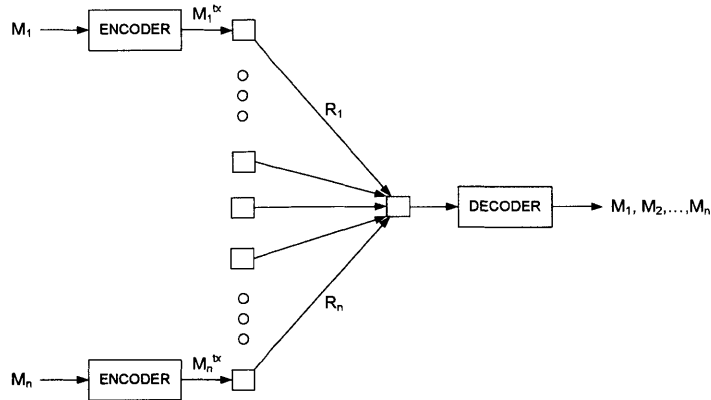


Figure 9-1: A depth one multicast tree with n sources.

ϵ). Then, reference [305] showed that if Φ_{S_m} satisfies RIP (9.3) and \mathbf{X} is sufficiently sparse, the following optimization recovers the original signal up to the noise level:

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_{\mathcal{L}_1} \\ \text{subject to} \quad & \|\mathbf{Y} - \Phi_{S_m} \mathbf{X}\|_{\mathcal{L}_2} \leq \epsilon. \end{aligned} \tag{9.5}$$

Hence, if \mathbf{X}^* is a solution of (9.5), then $\|\mathbf{X} - \mathbf{X}^*\|_{\mathcal{L}_2} \leq \beta\epsilon$, where β is a constant. To have this, the original vector \mathbf{X} should be sufficiently sparse. Specifically, [305] showed that, if

$$\delta_{3k} + 3\delta_{4k} < 2, \tag{9.6}$$

then, the results hold. We refer to the complexity of this optimization by $CX_{cs}^{noisy}(m, n)$.

9.1.2 Distributed Compression Background

In this section, we briefly review some prior results in distributed compression. In Section 9.2, we will use these results along with compressive sensing results, making a bridge these two subjects.

Consider the network depicted in Figure 9-1. This is a depth one tree network with n sources. Suppose each source i , has messages referred by M_i , wishes to send to

the receiver. To do this, each source i sends its message with a rate R_i (i.e., say \mathbf{m}_i has m messages of source i . This source maps \mathbf{m}_i to $\{1, 2, \dots, 2^{mR_i}\}$). Let $R = \sum_{i=1}^n R_i$. According to Slepian-Wolf compression [314], one needs to have,

$$R \geq H(M_1, M_2, \dots, M_n). \quad (9.7)$$

By using linear network coding, reference [315] extended Slepian-Wolf compression for a general multicast problem. Note that, the proposed decoders in [314] and [315] (a minimum entropy or a maximum probability decoder) have high complexity [316]. For a depth one tree network, references [317], [318], [319], [320], [321] and [322] provide some low complexity decoding techniques for Slepian-Wolf compression. However, none of these methods can explicitly provide a trade-off between the rate and the decoding complexity. For a general multicast network, it has been shown in [323] that, there is no separation between Slepian-Wolf compression and network coding. For this network, a low complexity decoder for an optimal compression scheme has not been proposed yet.

9.2 Compressive Sensing and Distributed Source Coding

In this part, we propose a distributed compression scheme for correlated sources providing a trade-off between the decoding complexity and the compression rate. References [314] and [315] propose optimal distributed compression schemes for correlated sources, for depth-one trees and general networks, respectively. However, their proposed compression schemes need to use a minimum entropy or a maximum probability decoder, which has high complexity [316]. On the other hand, compressive sensing provides a low complexity decoder, by using a convex optimization, to recover a sparse vector from incomplete and contaminated observations (e.g., [305], [306], [307] and [311]). Our aim is to use compressive sensing over networks to design compression schemes for correlated sources.

Note that, our proposed techniques can be applied on different network topologies, with different assumptions and constraints. Here, to illustrate the main ideas, first we consider a noiseless tree network with depth one and with correlated sources. Then, in Section 9.2.2, we extend our techniques to a general multicast network.

9.2.1 A One-stage Tree with n -Correlated Sources

Consider the network shown in Figure 9-1, which is a one-stage noiseless multicast tree network with n -sources. First, consider a case where sources are independent. Hence, the problem reduces to a classical source coding problem. Say each source is transmitting with a rate R_i . Denote $R = \sum_{i=1}^n R_i$. It is well-known that the following sum-rate for this network is the minimum required sum-rate:

$$R \geq \sum_{i=1}^n H(M_i), \quad (9.8)$$

where M_i is the message random variable of source i . Let us refer to the complexity of its decoder at the receiver as $CX_{indep}(n)$, since sources are independent.

Now, we consider a case where sources are correlated. We formulate the correlation of sources as follows:

Definition 9.1. Suppose $\mu_{i,t}$ represents a realization of the i^{th} source message random variable at time t (i.e., a realization of M_i at time t). Hence, the vector $\overline{\mu}_t = \{\mu_{1,t}, \dots, \mu_{n,t}\}$ is the sources' message vector at time t . We drop the subscript t when no confusion arises. Suppose sources are correlated so that there exists a transform matrix under which the sources' message vector is k -sparse (i.e., there exists a $n \times n$ transform matrix Φ and a k -sparse vector $\overline{\mu}'_t$ such that $\overline{\mu}_t = \Phi \overline{\mu}'_t$). Let \mathcal{S} be the power set of $\{1, 2, \dots, n\}$, and \mathcal{S}_m denote members of \mathcal{S} whose cardinality are m . Suppose $\Phi_{\mathcal{S}_m}$, defined as in Section 9.1.1, satisfies RIP (9.3), where m satisfies (9.1) for a given n and k . We have, $\overline{\mu_{t,\mathcal{S}_m}} = \Phi_{\mathcal{S}_m} \overline{\mu}'_t$, where $\overline{\mu_{t,\mathcal{S}_m}}$ is components of $\overline{\mu}_t$ whose indices belong to \mathcal{S}_m . We refer to these sources as *k -sparsely correlated sources*.

If we know zero locations of $\overline{\mu}'_t$ and if the transform matrix is non-singular, by

having each subset containing k sources, the whole vector $\bar{\mu}_t$ and $\bar{\mu}_t'$ can be computed. Hence, by the data processing inequality, we have,

$$H(M_{i_1}, \dots, M_{i_k}) + nH_b(\alpha) = H(M_1, \dots, M_n), \quad (9.9)$$

for and $1 \leq i_1 < i_2 < \dots < i_k \leq n$ where M_{i_j} is the message of source i_j . Also, by [306], having m sources, where m satisfies (9.1) for a given k and n , allows us to recover the whole vector $\bar{\mu}_t'$ and therefore, $\bar{\mu}_t$. Hence,

$$H(M_{i_1}, \dots, M_{i_m}) = H(M_1, \dots, M_n). \quad (9.10)$$

Example 9.2. For an example of k -sparsely correlated sources, suppose each coordinates of $\bar{\mu}_t'$ is zero with probability $1 - \alpha$, otherwise it is uniformly distributed over $\{1, 2, \dots, 2^R\}$. Suppose the entries of Φ are independent realizations of Bernoulli random variables:

$$\Phi_{i,j} = \begin{cases} \frac{1}{\sqrt{m}} & \text{with prob. } \frac{1}{2} \\ -\frac{1}{\sqrt{m}} & \text{with prob. } \frac{1}{2} \end{cases}. \quad (9.11)$$

Thus, as shown in [313], Φ_{S_m} satisfies RIP (9.3). Also, we have,

$$\begin{aligned} H(M_1, \dots, M_n) &= nH_b(\alpha) + k(R + 1) \\ &\approx nH_b(\alpha) + kR. \end{aligned} \quad (9.12)$$

The entropy of each source random variable for large n can be computed approximately as follows:

$$H(M_i) \approx R + \frac{1}{2} \log(k). \quad (9.13)$$

Note that, the individual entropies of k -sparsely correlated sources are roughly the same. Now, we want to use compressive sensing to have a distributed source coding

Table 9.1: Comparison of Distributed Compression Methods of Correlated Sources

Compression Methods	Minimum Min-Cut Rate	Decoding Complexity
SDCIC	$\sum_{j=1}^m H(M_{i_j})$	$CX_{indep}(m) + CX_{cs}^{noiseless}(m, n)$
Slepian-Wolf (SW)	$H(M_1, \dots, M_n)$	$CX_{sw}(n)$
Combination of SDC and SW	$H(M_1, \dots, M_n)$	$CX_{sw}(m) + CX_{cs}^{noiseless}(m, n)$
Naive Correlation Ignorance	$\sum_{i=1}^n H(M_i)$	$CX_{indep}(n)$

providing a trade-off between the compression rate and the decoding complexity.

We shall show that, if sources are k -sparsely correlated, if we have m of them at the receiver at each time t , by using a linear programming, all sources can be recovered at that time. Hence, instead of sending n correlated sources, one needs to transmit m of them where $m \ll n$. We assume that the entropies of sources are the same. Hence, each source transmits with probability $\gamma = \frac{m}{n}$. Therefore, by the law of large numbers, for large enough n , we have m transmitting sources. For the case of non-equal source entropies, one can a priori choose m sources with the lowest entropy as the transmitting sources.

In Table 9.1, we compare four different schemes:

- The first scheme which Theorem 9.3 is about, is called *Sparse Distributed Compression* with independent coding (SDCIC). The idea is to send m sources from n of them, assuming they are independent (i.e., their correlation is not used in coding). In this case, at the receiver, first we decode these m transmitted messages and then, we use a compressive sensing decoder to recover the whole sources. The decoding complexity is $CX_{indep}(m) + CX_{cs}^{noiseless}(m, n)$. The required min-cut rate between each receiver and sources is $\sum_{j=1}^m H(M_{i_j})$.
- In the second scheme referred in Table 9.1 as the *Slepian-Wolf* coding, one performs an optimal source coding on correlated sources. The required min-cut rate between sources and the receiver is $H(M_1, \dots, M_n)$ which is less than $\sum_{j=1}^m H(M_{i_j})$. However, at the receiver, we need to have a Slepian-Wolf decoder with complexity $CX_{sw}(n)$.
- One can have a combination of sparse distributed compression and Slepian-Wolf

compression. To do this, instead of n sources, we transmit m of them by using a Slepian-Wolf coding. The required min-cut rate would be $H(M_{i_1}, \dots, M_{i_m})$ which by (9.10), $H(M_1, \dots, M_n) = H(M_{i_1}, \dots, M_{i_m}) \leq \sum_{j=1}^m H(M_{i_j})$. At the receiver, first we use a Slepian-Wolf decoder with complexity $CX_{sw}(m)$ to decode the m transmitted sources' information. Then, we use a compressive sensing decoder to recover the whole sources.

- The fourth method is a naive way of transmitting n correlated sources so that we have an easy decoder at the receiver. We call this scheme a *Naive Correlation Ignorance* method. In this method, we simply ignore the correlation in the coding scheme. In this naive scheme, the required min-cut rate is $\sum_{i=1}^n H(M_i)$ with a decoding complexity equal to $CX_{indep}(n)$. Note that, in sparse distributed compression with independent coding, the required min-cut rate is much less than this scheme ($\sum_{j=1}^m H(M_{i_j}) \ll \sum_{i=1}^n H(M_i)$).

The following theorem is about sparse distributed compression with independent coding:

Theorem 9.3. *For a depth one tree network with n -sources, k -sparsely correlated, the required min-cut rate between the receiver and sources in sparse distributed compression with independent coding (SDCIC) is,*

$$R \geq \sum_{j=1}^m H(M_{i_j}), \quad (9.14)$$

where m is the smallest number satisfies (9.1) for a given n and k and i_j is the index of j^{th} active source. Also, the decoding complexity of SDCIC method is $CX_{indep}(m) + CX_{cs}^{\text{noiseless}}(m, n)$.

Proof. For a given n (number of sources) and k (the sparsity factor), we choose m to satisfy (9.1). Now, suppose each source is transmitting its block with probability γ where $\gamma = \frac{m}{n}$. Hence, by the law of large numbers, the probability that we have m transmitting sources' blocks goes to one for sufficiently large n . Without loss of

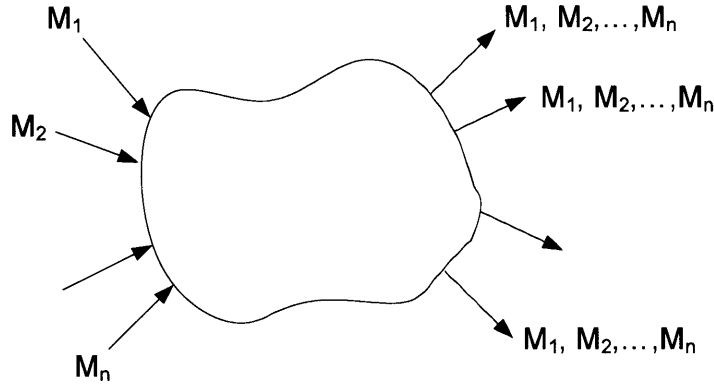


Figure 9-2: A general multicast network with n sources.

generality, say sources $1, \dots, m$ are transmitting (i.e., active sources) and other sources are idle. If $R_i \geq H(M_i)$ for active sources and zero for idle sources, active sources can transmit their messages to the receiver and the required min-cut rate between the receiver and sources would be $R \geq \sum_{j=1}^m H(M_{i_j})$. Let \mathcal{S}_m indicate active sources' indices. First, at the receiver, we decode these m active sources. The complexity of this decoding part is $CX_{indep}(m)$ because we did not use their correlation in the coding scheme. Hence, we have $\overline{\mu_{t,\mathcal{S}_m}}$, active sources' messages at time t at the receiver (in this case, $\overline{\mu_{t,\mathcal{S}_m}} = (\mu_{1,t}, \dots, \mu_{m,t})$). Then, by having $\overline{\mu_{t,\mathcal{S}_m}}$ and using the sparsity of $\overline{\mu}_t'$, the following optimization can recover the whole sources:

$$\begin{aligned} \min \quad & \|\overline{\mu}_t'\|_{\mathcal{L}_1} \\ \text{subject to} \quad & \overline{\mu_{t,\mathcal{S}_m}} = \Phi_{\mathcal{S}_m} \overline{\mu}_t'. \end{aligned} \tag{9.15}$$

The overall decoding complexity is $CX_{indep}(m) + CX_{cs}^{noiseless}(m, n)$. □

9.2.2 A General Multicast Network with Correlated Sources

In this section, we extend results of Section 9.2.1 to a general noiseless multicast network. Note that, for a depth one tree network, references [317], [318], [319], [320],

[321] and [322] provide some low complexity decoding techniques for Slepian-Wolf compression. However, none of these methods can explicitly provide a trade-off between the rate and the decoding complexity. For a general multicast network, it has been shown in [323] that, there is no separation between Slepian-Wolf compression and network coding. For this network, a low complexity decoder for an optimal compression scheme has not been proposed yet. There are few practical approaches and they rely on small topologies [324]. In this section, by using compressive sensing, we propose a coding scheme which provides a trade-off between the compression rate and the decoding complexity.

Consider a multicast network shown in Figure 9-2 which has n sources and T receivers. If sources are independent, reference [325] showed that the minimum required min-cut rate between each receiver and sources is as (9.8). This rate can be achieved by using linear network coding over the network. Reference [315] showed that random linear network coding can perform arbitrarily closely to this rate bound. Say $CX_{indep}^{nc}(n)$ is the complexity of its decoder.

Now, consider a case when sources are correlated. Reference [315] extended Slepian-Wolf compression result [314] to a multicast network by using network coding. The proposed decoder is a minimum entropy or a maximum probability decoder whose complexity is high [316]. We refer to its complexity by $CX_{sw}^{nc}(n)$.

In this section, we consider sources to be k -sparsely correlated. We want to use compressive sensing along with network coding to propose a distributed compression scheme providing a trade-off between the compression rate and the decoding complexity.

In the following, we compare four distributed compression methods for a general multicast network:

- In the first scheme, we use sparse distributed compression with independent coding. With high probability, we have m active sources. Then, we perform network coding on these active sources over the network, assuming they are independent. At the receiver, first we decode these m active sources and then, by a compressive sensing decoder, we recover the whole sources. The decod-

ing complexity is $CX_{indep}^{nc}(m) + CX_{cs}^{noiseless}(m, n)$. The required min-cut rate between each receiver and sources is $\sum_{j=1}^m H(M_{i_j})$.

- In the second scheme, we use an extended version of Slepian-Wolf compression for a multicast network [8]. Reference [8] considers a vector linear network code that operates on blocks of bits. The required min-cut rate between sources and the receiver is $H(M_1, \dots, M_n)$. At the receiver, one needs to have a minimum entropy or a maximum probability decoder with complexity $CX_{sw}^{nc}(n)$.
- A combination of sparse distributed compression and Slepian-Wolf compression can provide a trade-off between the compression rate and the decoding complexity. For example, instead of n sources, one can transmit m of them by using a Slepian-Wolf coding for multicast networks. The required min-cut rate would be $H(M_{i_1}, \dots, M_{i_m})$. At the receiver, first we use a Slepian-Wolf decoder with complexity $CX_{sw}^{nc}(m)$ to decode m active sources. Then, we use a compressive sensing decoder to recover the whole sources.
- In a naive correlation ignorance method, we simply ignore the correlation in the coding scheme. In this naive scheme, the required min-cut rate is $\sum_{i=1}^n H(M_i)$ with a decoding complexity equal to $CX_{indep}^{nc}(n)$.

The following theorem is about sparse distributed compression with independent coding for a multicast network:

Theorem 9.4. *For a general multicast network with n -sources, k -sparsely correlated, the required min-cut rate between each receiver and sources in sparse distributed compression with independent coding is,*

$$R \geq \sum_{j=1}^m H(M_{i_j}), \quad (9.16)$$

where m is the smallest number satisfies (9.1) for a given n and k . Also, the decoding complexity is $CX_{indep}^{nc}(m) + CX_{cs}^{noiseless}(m, n)$.

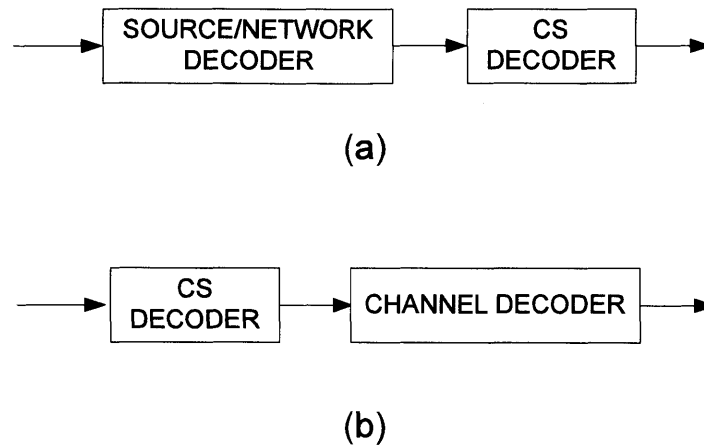


Figure 9-3: A modularized decoder for a) compressive sensing with source coding, (b) compressive sensing with channel coding.

Proof. The proof is similar to the one of Theorem 9.3. The only difference is that, here, one needs to perform linear network coding on active sources over the network without using the correlation among them. At the receiver, first, these m active sources' blocks are decoded (the decoding complexity of this part is $CX_{indep}^{nc}(m)$). Then, a compressive sensing decoder (9.15) is used to recover the whole sources. \square

Remark: In this section, we explained how compressive sensing can be used with distributed source coding. At the receiver side, first, we decode the active sources by a network decoding technique and then, we use a compressive sensing decoder to recover the whole sources. This high level modularized decoding scheme of sparse distributed compression is depicted in Figure 9-3-a. On the other hand, when we use compressive sensing for channel coding (which will be explained in detail in Section 9.3), we switch these two decoding modules (as shown in Figure 9-3-b). In other words, first, we use a compressive sensing decoder on a noisy signal to obtain a noisy estimate of the original signal and then, we use a channel decoder to find the original message from this noisy estimate.

9.3 Compressive Sensing and Channel Coding

In this section, we make a bridge between compressive sensing, which is more a signal processing technique, and channel coding. We consider a high-SNR Gaussian point-to-point channel depicted in Figure 9-4. We propose a coding scheme that combines compressive sensing and random channel coding with a modularized decoder to provide a trade-off between the capacity loss and the decoding complexity (i.e., the higher the capacity loss, the lower the decoding complexity.). We call this scheme *Sparse Channel Coding*. We add intentionally some correlation to transmitted signals to decrease the decoding complexity. This is an example of how compressive sensing can be used with channel coding and could be extended to other types of channels.

Consider a point to point channel with additive Gaussian noise with noise power N and transmission power constraint P , shown in Figure 9-4. Suppose we are in a high SNR regime (i.e., $\frac{P}{N} \gg 1$). The capacity of this channel is,

$$C \approx \frac{1}{2} \log \left(\frac{P}{N} \right). \quad (9.17)$$

An encoding and a decoding scheme for this channel can be found in [316]. As explained in [316], at the receiver, one needs to use a m -dimensional maximum likelihood decoder where m , the code length, is arbitrarily large. We refer to the complexity of such a decoder as $CX_{ml}(m)$. Our aim is to design a coding scheme to have a trade-off between the rate and the decoding complexity.

Suppose m is arbitrarily large. Choose n and k to satisfy (9.6).

Definition 9.5. A $(2^{mR}, m)$ sparse channel code for a point-to-point Gaussian channel with power constraint P satisfies the following:

- Encoding process: message i from the set $\{1, 2, \dots, 2^{mR}\}$ is assigned to a m -length vector \mathbf{W}_i , satisfying the power constraint; that is, for each codeword, we have,

$$\frac{1}{m} \sum_{j=1}^m W_{ij}^2 \leq P, \quad (9.18)$$

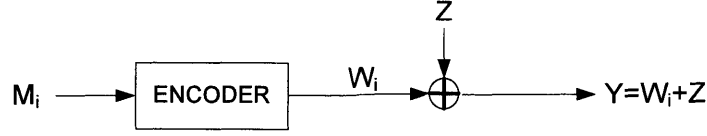


Figure 9-4: A Point-to-Point Gaussian Channel

where W_{ij} is the j^{th} coordinate of codeword W_i .

- Decoding process: the receiver receives Y , a noisy version of the transmitted codeword W_i (i.e., $Y = W_i + Z$, where Z is a Gaussian noise vector with i.i.d. coordinates with power N). A decoding function g maps Y to the set $\{1, 2, \dots, 2^{mR}\}$. The decoding complexity is $CX_{ml}(k) + CX_{cs}^{noisy}(m, n)$.

A rate R is achievable if $P_e^m \rightarrow 0$ as $m \rightarrow \infty$, where,

$$P_e^m = \frac{1}{2^{mR}} \sum_{i=1}^{2^{mR}} Pr(g(\mathbf{Y}) \neq i | M = i), \quad (9.19)$$

and M represents the transmitted message.

Theorem 9.6. *For a high-SNR Gaussian point-to-point channel described in Definition 9.5, the following rate is achievable, while the decoding complexity is $CX_{ml}(k) + CX_{cs}^{noisy}(m, n)$:*

$$R = \frac{k}{m}C + \frac{n}{m}H_b(\alpha) + \frac{k}{2m} \log\left(\frac{1}{\beta(1 + \delta_k)}\right), \quad (9.20)$$

where C is the channel capacity (9.17), α , β and δ_k are parameters defined in Section 9.1.1, and $H_b(\cdot)$ is the binary entropy function.

Before expressing the proof, let us make some remarks:

- if we have (9.1) (i.e., $m/n = \rho\alpha \log(1/\alpha)$), the achievable rate can be approximated as follows:

$$R \approx \frac{1}{\log(1/\alpha)}C + \frac{1}{\alpha \log(1/\alpha)}H_b(\alpha) \quad (9.21)$$

where $\alpha = k/n$ is the sparsity ratio. The first term of (9.21) shows that the capacity loss factor is a log function of the sparsity ratio. The term $\frac{1}{\alpha \log(1/\alpha)} H_b(\alpha)$ is the rate gain that we obtain by using compressive sensing in our scheme. Note that, since the SNR is sufficiently high, the overall rate would be less than the capacity.

- The complexity term $CX_{ml}(k)$ is an exponential function, while $CX_{cs}^{noisy}(m, n)$ is a polynomial.

Proof. We use random coding and compressive sensing arguments to propose a concatenated channel coding scheme satisfying (9.20). First, we explain the encoding part:

- Each message i is mapped randomly to a k -sparse vector \mathbf{X}_i with length n so that its non-zero coordinates are i.i.d. drawn from $\mathcal{N}(0, \frac{m}{k(1+\delta_k)}P)$. n is chosen to satisfy (9.6), for a given k and m .
- To obtain \mathbf{W}_i , we multiply \mathbf{X}_i by a $m \times n$ matrix (Φ_{S_m}) satisfying the RIP condition (9.3) (i.e., $\mathbf{W}_i = \Phi_{S_m} \mathbf{X}_i$). We send \mathbf{W}_i through the channel.

In fact, it is a concatenated channel coding scheme [326]. The outer layer of this coding is the sparsity pattern of \mathbf{X}_i . The inner layer of this coding is based on random coding around each sparsity pattern. To satisfy transmission power constraint P , we generate each non-zero coordinate of \mathbf{X}_i by $\mathcal{N}(0, \frac{m}{k(1+\delta_k)}P)$. Note that, the encoding matrix Φ_{S_m} is known to both encoding and decoding sides.

At the receiver, we have a noisy version of the transmitted signal. Suppose message i has been transmitted. Hence, $\mathbf{Y} = \mathbf{W}_i + \mathbf{Z} = \Phi_{S_m} \mathbf{X}_i + \mathbf{Z}$. The decoding part is as follows:

- Since \mathbf{X}_i is k -sparse, first, we use a compressive sensing decoder to find $\tilde{\mathbf{X}}_i$ such that $\|\mathbf{X}_i - \tilde{\mathbf{X}}_i\|_{\mathcal{L}_2}^2 \leq \beta mN$, as follows, where β is the parameter of (9.5):

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_{\mathcal{L}_1} & (9.22) \\ \text{subject to} \quad & \frac{1}{m} \|\mathbf{Y} - \Phi_{S_m} \mathbf{X}_i\|_{\mathcal{L}_2}^2 \leq N. \end{aligned}$$

We assume that the complexity of this convex optimization is smaller than the one of a maximum likelihood decoder.

- Since we are in a high-SNR regime, by having $\tilde{\mathbf{X}}_i$, we can find the sparsity pattern of \mathbf{X}_i . It gives us the outer layer code of message i . The higher the outer layer code rate, the higher the required SNR. We shall develop this argument with more detail later (9.26).
- Having the sparsity pattern of \mathbf{X}_i , we use a maximum likelihood decoder in a k -dimensional space (i.e., non-zero coordinates of \mathbf{X}_i) to find non-zero coordinates of \mathbf{X}_i (the inner layer code). The complexity of this decoder is denoted by $CX_{ml}(k)$.

Before presenting the error probability analysis, let us discuss different terms of (9.20). Since in our coding scheme, for each message i , we send a k -sparse signal by sending a m -length vector, we have the fraction $\frac{k}{m}$ before the capacity term C . The capacity term C in (9.20) comes from the inner layer coding scheme. The term $\frac{n}{m} H_b(\alpha)$ comes from the outer layer coding scheme. Note that, since α is a small number, $H_b(\alpha)$ is close to one. Also, the ratio $\frac{n}{m}$ depends on the outer layer code rate. Here, we assumed the SNR is sufficiently high so that we can use all possible outer layer codes. The term $\frac{k}{2m} \log\left(\frac{1}{\beta(1+\delta_k)}\right)$ is because of the power constraint and the RIP condition. Note that, if one performs a time-sharing based channel coding, a rate $\frac{k}{m}C$ is achievable with a decoding complexity $CX_{ml}(k)$. By using the compressive sensing, we obtain additional rate terms because of the outer layer coding with approximately the same complexity.

We proceed by the error probability analysis. Define Ξ_i as the sparsity pattern of \mathbf{X}_i (i.e., $\Xi_{ij} = 0$ if $X_{ij} = 0$, otherwise, $\Xi_{ij} = 1$). Also, define $\tilde{\Xi}_i$ as the decoded

sparsity pattern of message i . Without loss of generality, assume that message 1 was transmitted. Thus, $\mathbf{Y} = \Phi_{S_m} \mathbf{X}_1 + \mathbf{Z}$. Define the following events:

$$\begin{aligned}
 E_0 &= \left\{ \frac{1}{m} \sum_{j=1}^m W_{1j}^2 > P \right\} \\
 E_i &= \left\{ \frac{1}{n} \|\mathbf{X}_i - \tilde{\mathbf{X}}_1\|_{\mathcal{L}_2}^2 \leq \frac{\beta m}{n} N \right\} \\
 E_{p_1} &= \{ \tilde{\Xi}_1 \neq \Xi_1 \}.
 \end{aligned} \tag{9.23}$$

Hence, an error occurs when E_0 occurs (i.e., the power constraint is violated), or E_{p_1} occurs (i.e., the outer layer code is wrongly decoded), or E_1^c occurs (i.e., the underlying sparse signal of the transmitted message \mathbf{X}_1 and its decoded noisy version $\tilde{\mathbf{X}}_1$ are in a distance greater than the noise level), or one of E_i occurs while $i \neq 1$. Say M is the transmitted message (in this case $M = 1$) and \hat{M} is the decoded message. Let \mathcal{E} denote the event $M \neq \hat{M}$. Hence,

$$\begin{aligned}
 Pr(\mathcal{E} | M = 1) &= Pr(\mathcal{E}) \\
 &= Pr(E_0 \cup E_{p_1} \cup E_1^c \cup \bigcup_{j=2}^{2^{mR}} E_j) \\
 &\leq Pr(E_0) + Pr(E_{p_1}) + Pr(E_1^c) \\
 &\quad + \sum_{j=2}^{2^{mR}} Pr(E_j),
 \end{aligned} \tag{9.24}$$

by union bounds of probabilities. We bound these probabilities term by term as follows:

- By the law of large numbers, we have $Pr(\|\mathbf{X}_1\|_{\mathcal{L}_2}^2 > \frac{m}{1+\delta_k} P + \epsilon) \rightarrow 0$ as $n \rightarrow \infty$.
By using the RIP condition (9.3) and with probability one, we have,

$$\|\Phi_{S_m} \mathbf{X}_1\|_{\mathcal{L}_2}^2 \leq (1 + \delta_k) \|\mathbf{X}_1\|_{\mathcal{L}_2}^2 \leq mP. \tag{9.25}$$

Hence, $Pr(E_0) \rightarrow 0$ as $n \rightarrow \infty$.

- After using the compressive sensing decoder mentioned in (9.22), we have $\tilde{\mathbf{X}}_1$ such that $\frac{1}{n}\|\mathbf{X}_1 - \tilde{\mathbf{X}}_1\|_{\mathcal{L}_2}^2 \leq \frac{\beta m}{n}N$. Suppose this error is uniformly distributed over different coordinates. We use a threshold comparison to determine the sparsity pattern $\tilde{\Xi}_1$. The probability of error in determining $\tilde{\Xi}_1$ determines how high the SNR should be. Intuitively, if we use all possible outer layer codes (all possible sparsity patterns), we should not make any error in determining the sparsity pattern of each coordinate j (i.e., $\tilde{\Xi}_{1j}$). Hence, we need sufficiently high SNR for a given n . On the other hand, if we decrease the outer layer code rate, the required SNR would be lower than the case before. Here, to illustrate how the required SNR can be computed, we assume that we use all possible sparsity patterns. Say $E_{p_{j1}}$ is the event of making an error in determining whether the j^{th} coordinate is zero or not. We say $\tilde{\Xi}_{1j} = 0$ if $X_{ij} < \tau$. Otherwise, $\tilde{\Xi}_{1j} = 1$. Hence, for a given threshold τ , we have,

$$\begin{aligned} Pr(E_{p_{j1}}|\Xi_{1j} = 0) &= \phi\left(\frac{\tau}{\sqrt{\frac{\beta m}{n}N}}\right) \\ Pr(E_{p_{j1}}|\Xi_{1j} = 1) &= 1 - 2\phi\left(\frac{\tau}{\sqrt{P + \frac{\beta m}{n}N}}\right) \end{aligned}$$

where $\phi(\cdot)$ is the cumulative distribution function of a normal random variable.

Hence,

$$\begin{aligned} Pr(E_{p_1}) &= 1 - Pr(E_{p_1}^c) & (9.26) \\ &= 1 - \left(1 - \phi\left(\frac{\tau}{\sqrt{\frac{\beta m}{n}N}}\right)\right)^{n(1-\alpha)} \\ &\quad \times \left(2\phi\left(\frac{\tau}{\sqrt{P + \frac{\beta m}{n}N}}\right)\right)^{n\alpha} \end{aligned}$$

Note that, one can choose τ and $\frac{P}{N}$ large enough to have $Pr(E_{p_1})$ arbitrarily small.

- By [306], $Pr(E_1^c)$ is zero.
- At the last step, we need to bound $Pr(E_j)$ for $j \neq 1$. We have,

$$\begin{aligned}
 Pr(E_j) &= Pr(\Xi_j = \Xi_1)Pr(E_j|\Xi_j = \Xi_1) \\
 &+ Pr(\Xi_j \neq \Xi_1)Pr(E_j|\Xi_j \neq \Xi_1) \\
 &= \alpha^k(1 - \alpha)^{n-k}Pr(E_j|\Xi_j = \Xi_1).
 \end{aligned} \tag{9.27}$$

Note that, $Pr(E_j|\Xi_j \neq \Xi_1)$ goes to zero in a high SNR regime. By bounding the noise power of non-zero coordinates by βmN (the whole noise power) and using random coding argument in a k dimensional space, we have,

$$Pr(E_j|\Xi_j = \Xi_1) \leq \frac{1}{\left(1 + \frac{1}{\beta(1+\delta_k)} \frac{P}{N}\right)}, \tag{9.28}$$

where \leq indicates the right hand side is less or equal than the left one in an exponential rate. Therefore, by (9.27) and (9.28), and doing some manipulations, we have,

$$P(E_j) \leq 2^{-n(H_b(\alpha) + \frac{\alpha}{2} \log(\frac{1}{\beta(1+\delta_k)} \frac{P}{N}))}, \tag{9.29}$$

for $j = 2, \dots, 2^{mR}$.

Thus, in a sufficiently high SNR regime, by (9.24) and (9.29), we have,

$$\begin{aligned}
 P_e^m = Pr(\mathcal{E}) &= Pr(\mathcal{E}|M = 1) \\
 &\leq 2^{mR} 2^{-n(H_b(\alpha) + \frac{\alpha}{2} \log(\frac{1}{\beta(1+\delta_k)} \frac{P}{N}))} \\
 &\leq 2^{m(R - \frac{n}{m} H_b(\alpha) - \frac{k}{m} C - \frac{k}{2m} \log(\frac{1}{\beta(1+\delta_k)} \frac{P}{N}))}.
 \end{aligned} \tag{9.30}$$

For any given $\epsilon > 0$, we can choose n large enough to have $R \leq \frac{k}{m}C + \frac{n}{m}H_b(\alpha) + \frac{k}{2m} \log\left(\frac{1}{\beta(1+\delta_k)}\right) - \epsilon$ achievable. \square

Chapter 10

Joint Source-Channel-Network

Coding using Compressive Sensing

The power budget and communication rates are two main limitations in wireless network applications. Information theory focuses much on minimizing transmission rates, but the complexity of coding processes. For example, in point-to-point channels, equivalence theory allows us to separate the problem of channel coding from that of joint network/source coding [316]. However, any approach that separates source coding from channel coding may suffer from potential inefficiencies since we reduce redundancy for source coding but later reintroduce it for channel coding. Moreover, even if we allow separate channel coding from source/network coding (which is necessary because of lack of separation [323]), the known schemes have generally high complexity inherited from the general Slepian-Wolf problem [314].

Compressive sensing ([305], [306], [307]) can be used to perform channel coding ([327]) as well as source coding, even for distributed sources ([327]). In both of these cases, the performance of the compressive sensing approach can be shown to be close to that of the information theoretic optimum. The fact that, a single approach may successfully be applied for channel coding and source coding renders it an attractive candidate for use in a joint problem. Moreover, the fact that, linear analog network coding is asymptotically optimal in the high SNR regime [328], together with the fact that, compressive sensing-based channel and source techniques are also linear, provides a compelling reason for considering a compressive sensing-based joint source-channel-network coding.

Here, we propose a joint source-channel-network coding scheme, based on compressive sensing principles, for wireless networks with AWGN channels (that may include multiple access and broadcast). We assume that, sources exhibit temporal and spatial dependencies. We consider a multicast wireless network with multiple sources and receivers (see Figure 10-1). Our goal is to provide a reconstruction of all sources within an allowed distortion level at all receivers in a low-complexity manner. A key idea of our proposed scheme is to reshape inherent redundancies among source samples (temporal or spatial) in a joint source-channel-network coding scheme. Our scheme does not rely on the distribution knowledge of sources. Also, it provides an attractive relationship between the number of uses of the network and the reconstruc-

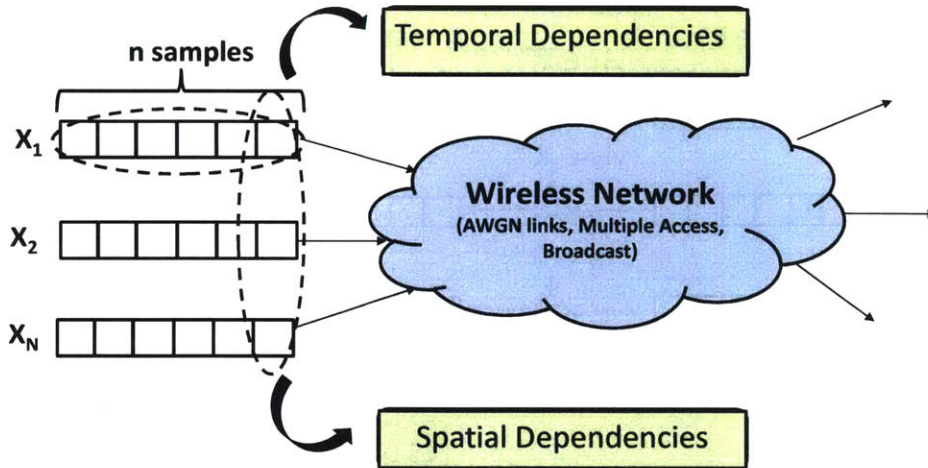


Figure 10-1: The considered problem setup where sources have temporal and spatial dependencies, and receivers desire to reconstruct sources within an allowed distortion level.

tion error.

A schematic view of our proposed framework is illustrated in Figure 10-2. In this framework, we perform joint source-channel coding at each source by randomly projecting source values to a lower dimensional space. We consider sources that spatially satisfy the restricted eigenvalue (RE) condition and may use sparse networks as well as more general sources, which use the randomness of the network to allow mapping to lower dimensional spaces. Our scheme relies on using random linear network coding in the real field. The receiver uses compressive sensing decoders (spatial and temporal ones) to reconstruct source signals.

Our main insight is that, the main characteristics required for operating compressive sensing are preserved under very general linear operations, which may themselves be arise from cascade of linear transformations in a network. We can therefore cascade compressive sensing techniques as well as linear operations tied to coding in the interior of the network. This preservation of the main characteristics of the sources under linear transformations lies at the core of both compressive sensing and analog network coding. This work constitutes a step in bringing together these two techniques, which have natural similarities.

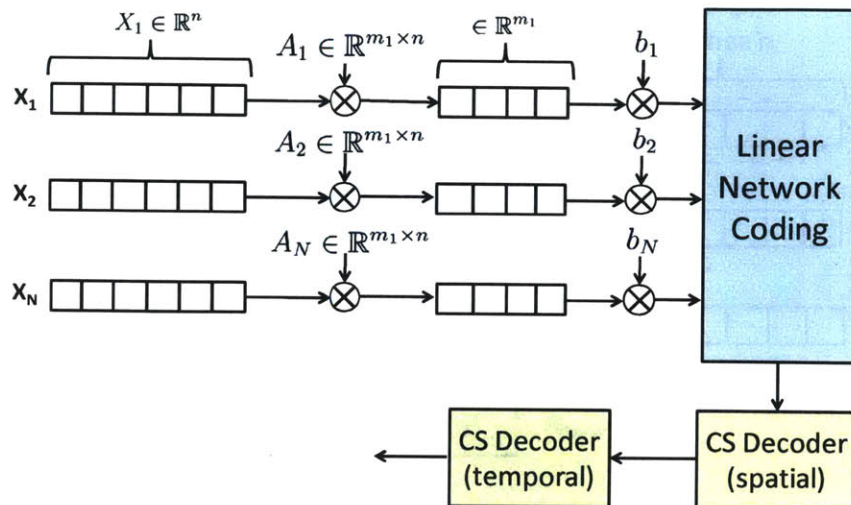


Figure 10-2: A schematic view of our proposed joint source-channel-network coding scheme. Each source projects its high dimensional signal to a lower dimensional space by using a random projection matrix followed by a pre-coding on-off module. Intermediate nodes perform linear network coding. Each receiver uses compressive sensing decoders (LASSO) to reconstruct source signals.

This chapter is organized as follows. In Section 10.1, we review some prior results in compressive sensing and random projections. We propose our scheme in Section 10.2 and characterize its performance. Proofs are presented in Section 10.3.

10.1 Compressive Sensing Background

In this section, we review some prior results in compressive sensing. Let $\mathbf{Y} \in \mathbb{R}^N$ be an unknown signal vector. We say this signal is k -sparse if k of its coordinates are non-zero and the rest are zero. Note that, the assumption of having $N - k$ components of \mathbf{Y} to be exactly zero may seem unrealistic for practical cases. This assumption is called a hard sparsity assumption. Weakly sparse models are also considered in the compressive sensing literature [329]. Roughly, a vector \mathbf{Y} is weakly sparse if it can be approximated by a sparse vector. In this chapter, we consider hard sparsity models although all discussions can be extended to weakly sparse signals.

Let $G \in \mathbb{R}^{m \times N}$ be a measurement matrix where $m \ll N$. The observation vector

\mathbf{Z} is defined by the following linear model:

$$\mathbf{Z} = \mathbf{G}\mathbf{Y} + \mathbf{W},$$

where $\mathbf{W} \in \mathbb{R}^m$ is a noise vector. We assume $w_i \sim \mathcal{N}(0, \sigma^2)$, where $\mathcal{N}(0, \sigma^2)$ represents a Gaussian random variable with mean zero and variance σ^2 .

By obtaining observations \mathbf{Z} and knowing measurement matrix \mathbf{G} , the goal is to reconstruct $\tilde{\mathbf{Y}}$ such that

$$\|\mathbf{Y} - \tilde{\mathbf{Y}}\|_{l_2}^2 \leq D,$$

where D is a distortion allowance threshold. If the signal \mathbf{Y} is a sparse signal (i.e., $k \ll N$) and the measurement matrix satisfies the *Restricted Eigenvalue (RE) Condition* [330] (which will be explained later) or more restrictively, *Restricted Isometry Property* ([311]), then having m measurements where $m \ll N$ is sufficient to recover the original sparse vector. In fact, $\tilde{\mathbf{Y}}$ is a solution for the following convex optimization called LASSO ([331]):

$$\min_{\mathbf{Y}} \frac{1}{2m} \|\mathbf{Z} - \mathbf{G}\mathbf{Y}\|_{l_2}^2 + \xi \|\mathbf{Y}\|_{l_1}, \quad (10.1)$$

where $\|\cdot\|_{l_p}$ represents the l_p norm of a vector and $\xi > 0$ is a regularization parameter. If $\tilde{\mathbf{Y}}$ is a solution of this optimization, then, [330] shows that,

$$\|\mathbf{Y} - \tilde{\mathbf{Y}}\|_{l_2}^2 \leq \frac{\delta}{\gamma^2} \sigma^2 \frac{k \log(N)}{m} \quad (10.2)$$

with high probability, where δ is a constant and γ is a parameter related to the restricted eigenvalue condition of the matrix \mathbf{G} .

Now, we define the restricted eigenvalue condition (RE) for the measurement matrix $\mathbf{G} \in \mathbb{R}^{m \times N}$ ([330]). Suppose \mathcal{S} is the support of a sparse vector \mathbf{Y} (i.e., \mathcal{S} is the set of non-zero indices of \mathbf{Y} where $|\mathcal{S}| = k$). \mathcal{S}^c is the complement set of \mathcal{S} . We

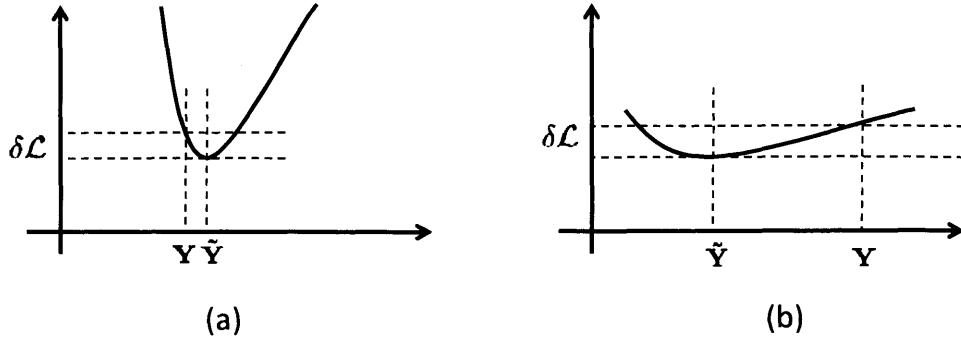


Figure 10-3: To control the reconstruction error, the loss function should have high curvature around the optimal solution of LASSO (equation 10.1) along the directions of $\mathcal{C}(\mathcal{S}; \alpha)$.

define the subset

$$\mathcal{C}(\mathcal{S}; \alpha) = \{\mathbf{Y} \in \mathbb{R}^N : \|\mathbf{Y}_{S^c}\|_{l_1} \leq \alpha \|\mathbf{Y}_S\|_{l_1}\}.$$

This corresponds to the cone of vectors where l_1 norm on their support dominates the l_1 norm off the support. For $\alpha = 1$, this cone is the same as the one used in restricted null-space property ([332]).

Definition 10.1. The matrix \mathbf{G} satisfies the restricted eigenvalue (RE) condition over \mathcal{S} with parameter (γ, α) if

$$\frac{\frac{1}{m} \|\mathbf{G}\mathbf{Y}\|_{l_2}^2}{\|\mathbf{Y}\|_{l_2}^2} \geq \gamma > 0 \quad (10.3)$$

for all $\mathbf{Y} \in \mathcal{C}(\mathcal{S}; \alpha) - \{0\}$.

The proof of equation (10.2) can be found in reference ([330]). By abuse of notation, we say a square matrix satisfies the restricted eigenvalue condition with parameter (γ, α) if all of its sub-matrices with $|\mathcal{S}|$ rows satisfies the RE condition with parameter (γ, α) . Also, roughly we say a matrix satisfy the RE condition, if γ is large enough. In the following, we develop some intuitions about the compressive sensing decoder of equation (10.1) and the RE condition.

The restricted eigenvalue condition bounds away eigenvalues of the matrix $\frac{1}{m} \mathbf{G}^T \mathbf{G}$

from zero in the direction of the considered cone. Note that, to be able to control error of LASSO in equation (10.1), around the optimal point, the loss function should have high curvature to translate small l_2 loss to small error in \mathbf{Y} (see Figure 10-3). For this loss function, Hessian (curvature) can be calculated as $\frac{1}{m}\mathbf{G}^T\mathbf{G}$. Therefore, we want to bound eigenvalues of $\frac{1}{m}\mathbf{G}^T\mathbf{G}$ away from zero to have high curvature of the loss function. However, since \mathbf{G} is a fat matrix (i.e., $m \ll N$), some of its eigenvalues are zero. That is why we require to have eigenvalues away from zero for a subset $\mathcal{C}(\mathcal{S}; \alpha)$ of vectors. This subset is well-aligned with the directions of the Hessian in which loss is observed ([330]).

It is shown in different references that random matrices satisfy the RE condition. Reference [312] shows a connection between compressive sensing, n -widths, and the Johnson-Lindenstrauss Lemma. Through this connection, [313] proposes some database friendly matrices (binary matrices) satisfying the RE condition.

Our goal in this chapter is to use random projection principles to obtain a power-rate efficient sensing-communication scheme. New results and bounds in compressive sensing can potentially improve results of this chapter. One can find more details on compressive sensing results in [305], [306], [307], etc.

In the next section, we present main results of this chapter.

10.2 Main Results

In this section, we present main results of this chapter of the thesis. First, we present a Lemma demonstrating how the RE condition changes when matrices are cascaded (multiplied to each other). Then, after introducing the problem setup and some notations, we present our proposed joint source-channel-network coding scheme. Finally, we characterize the performance of this proposed scheme. Proofs are presented in the next section.

10.2.1 Background Results

In this section, we present a Lemma which is central to implementing the cascading linear transformations which form the basis of our approach.

Lemma 10.2 (Cascading Lemma). Suppose $\mathbf{G} \in \mathbb{R}^{m \times N}$ satisfies the RE condition with parameter (γ, α) . Say $\mathbf{C}_1 \in \mathbb{R}^{m \times m}$ and $\mathbf{C}_2 \in \mathbb{R}^{N \times N}$. Suppose minimum absolute eigenvalues of \mathbf{C}_1 and \mathbf{C}_2 are λ_1 and λ_2 , respectively. Also, say for $\mathbf{Y} \in \mathcal{C}(\mathcal{S}; \alpha)$, $\mathbf{Y}\mathbf{C}_2 \in \mathcal{C}(\mathcal{S}; \alpha)$. Then, $\mathbf{C}_1\mathbf{G}$ and $\mathbf{G}\mathbf{C}_2$ satisfy the RE condition with parameters $(\gamma\lambda_1^2, \alpha)$ and $(\gamma\lambda_2^2, \alpha)$, respectively.

The proof of this Lemma is presented in Section 10.3. Intuitively, this lemma says, if \mathbf{G} , which satisfies the RE condition, is cascaded with other matrices whose eigenvalues are bounded away from zero, the overall matrix would satisfy the RE condition.

10.2.2 The problem Setup and Notations

We consider a wireless network with N sources and l receivers. Each source has n correlated samples to transmit to all receivers (see Figure 10-1). Also, we assume that, there are temporal and spatial correlations among source samples. Vector $\mathbf{X}_i \in \mathbb{R}^n$ represents samples of source i and $\mathbf{X} \in \mathbb{R}^{N \times n}$ is a matrix of all source samples. x_i^t represents a sample of source i at time t .

We have the following assumptions on source samples capturing their temporal and spatial redundancies:

Assumption 10.3. Sources have temporal and spatial dependencies:

1) Each source has temporal dependency: there exists a matrix $\Phi \in \mathbb{R}^{n \times n}$ such that $\mathbf{X}_i = \Phi\eta_i$, where η_i is k_1 sparse and $k_1 \ll n$.

2) Source samples have spatial dependency: suppose $\mathbf{Y} \in \mathbb{R}^N$ is a vector containing $y_j = \mathbf{A}_j^T \mathbf{X}_j$ for some vector $\mathbf{A}_j \in \mathbb{R}^n$, for $1 \leq j \leq N$. There exists a matrix $\Psi \in \mathbb{R}^{N \times N}$ such that, $\mathbf{Y} = \Psi\mu$, where μ is k_2 sparse, and $k_2 \ll N$.

Note that, matrices Φ and Ψ are not required to be known at sources, but at receivers only. This is an important issue in hardware implementations of sensors ([333]). We will show that, our proposed scheme is generic.

Each receiver desires to obtain all sources within an allowed distortion level D . In other words, for all $1 \leq i \leq N$ and all receivers, we should have

$$\frac{1}{n} \|\mathbf{X}_i - \tilde{\mathbf{X}}_i\|_{l_2}^2 \leq D, \quad (10.4)$$

with high probability, where $\tilde{\mathbf{X}}_i$ is the reconstructed signal of source i .

Links of the network are AWGN channels. Nodes in the network are allowed to perform linear network coding in the real field (see analog network coding references [334], [328], etc.). We assume that, the network can deliver m linear combinations of source samples to each receiver at each use of the network. In other words, at each use of the network, a receiver obtains $\mathbf{Z} \in \mathbb{R}^m$ where, $\mathbf{Z} = \mathbf{G}\mathbf{Y} + \mathbf{W}$. $\mathbf{G} \in \mathbb{R}^{m \times N}$ is the network coding matrix, $\mathbf{Y} \in \mathbb{R}^N$ is a vector that sources transmit at that time (i.e., source i transmits y_i). \mathbf{W} is the noise vector where we assume, $w_i \sim N(0, \sigma^2)$, for all $1 \leq i \leq m$. The quantity m can be viewed as a real field version of the min-cut rate of the network defined in the finite field ([325], [8]). We assume that, the number of source samples (i.e., $n \times N$) is much larger than m . We define \mathcal{C}_{use} as the number of network uses to deliver all source samples to receivers so that, the reconstruction distortion requirement (equation 10.4) holds for each receiver. We characterize \mathcal{C}_{use} for our proposed sensing-communication framework and compare it with a naive information theoretic approach.

10.2.3 Joint Source-Channel-Network Coding Scheme

In this section, we present our proposed joint source-channel-network coding scheme based on compressive sensing principles. We also characterize the rate-distortion performance of our proposed scheme in Theorem 10.6.

Our proposed sensing-communication scheme can be performed in the following steps, as illustrated in Figure 10-2. We will explain each step in more detail.

Algorithm 10.4. Our proposed joint source-channel-network coding can be performed in following steps.

- **Step 1 (Temporal Pre-Coding):** Each source projects its n dimensional signal to a lower dimensional space by using random projection matrices.
- **Step 2 (Spatial Pre-Coding):** Each source i transmits with certain probability b_i , at each time.
- **Step 3 (Network Coding):** Nodes in the network perform analog random linear network coding.
- **Step 4 (Decoding):** Each receiver uses compressive sensing decoders (LASSO) to reconstruct source signals.

These steps are shown in a schematic way in Figure 10-2. In the following, we explain each step in more details:

Step 1 (Temporal Pre-Coding): Each source i projects its sample vector \mathbf{X}_i to a lower dimensional space, by using a matrix multiplication. That is, $\mathbf{Y}_i = \mathbf{A}_i \mathbf{X}_i$, where $\mathbf{X}_i \in \mathbb{R}^n$ and $\mathbf{A}_i \in \mathbb{R}^{m_1 \times n}$. Note that, $m_1 \ll n$, and m_1 is chosen according to Theorem 10.6 to control the reconstruction error. Matrix \mathbf{A}_i is a random matrix that satisfies the RE condition with parameter (γ_1, α) .

Step 2 (Spatial Pre-Coding): At time t , each source transmits y_i^t (the t -th component of \mathbf{Y}_i) with probability b_i . In the next step, we will show how by choosing appropriate b_i , the transmission power can be decreased further.

Step 3 (Network Coding): Nodes in the network perform random linear network coding over the real field. At each use of the network, it can deliver m linear combinations of transmitted samples to receivers. Therefore, m_2/m network transmissions are required to gather m_2 linear combinations of y_i^t at each receiver, with high probability. Note that, $m_2 \ll N$ is chosen according to Theorem 10.6 to control the reconstruction error. By abuse of notation, we drop superscript t . Hence, a receiver obtains:

$$\mathbf{Z} = \mathbf{GBY} + \mathbf{W} = \mathbf{GB}\Psi\boldsymbol{\mu} + \mathbf{W}, \quad (10.5)$$

where $\mathbf{Y} \in \mathbb{R}^N$ contains y_i for $1 \leq i \leq N$, $\mathbf{B} \in \mathbb{R}^{N \times N}$ is a diagonal matrix whose i -th diagonal element is one with probability b_i , and zero otherwise. $\mathbf{G} \in \mathbb{R}^{m_2 \times N}$ is the network coding matrix for that receiver. Note that, by using Assumption 10.3, $\mathbf{Y} = \Psi\mu$, where μ is k_2 sparse.

To be able to use compressive sensing principles and decoders, the matrix $\mathbf{GB}\Psi$ should satisfy the RE condition explained in 10.1. We consider two cases: (1) Ψ satisfies the RE condition, (2) Ψ does not satisfy the RE condition.

In case (1), since Ψ satisfies the RE condition, by using Lemma 10.2, $\mathbf{GB}\Psi$ would satisfy the RE condition. Since the rank of $\mathbf{GB}\Psi$ needs to be m_2 with high probability, we choose $b_i = 1$ with probability m_2/N , and $b_i = 0$ with probability $1 - m_2/N$. In other words, only a fraction of sources are transmitting. Also, by performing random linear network coding over the real field, rank of the matrix \mathbf{G} would be m_2 , with high probability ([315]). Hence, for a general network structure, receivers can use compressive sensing decoders to reconstruct source signals within an allowed distortion level.

In case (2) where Ψ does not satisfy the RE condition, we choose $b_i = 1$ for all i (i.e., all sources are transmitting). Hence, \mathbf{B} is the identity matrix. In this case, if \mathbf{G} satisfies the RE condition, by using Lemma 10.2, $\mathbf{G}\Psi$ would satisfy the RE condition, and therefore, receivers can use compressive sensing decoders to reconstruct source signals.

Intuitively, to have \mathbf{G} satisfy the RE condition, the network should be dense to randomly mix source samples. For networks that are not sufficiently dense, such as tree networks, \mathbf{G} never satisfies the RE condition. The following example presents a family of networks for which \mathbf{G} can be designed to satisfy the RE condition.

Example 10.5. Consider the network topology with N sources depicted in Figure 10-4. Sources are connected randomly to m intermediate nodes. Expected in-degree of each of these intermediate nodes is greater than or equal to $N/3$ (i.e., high in-degree nodes). These intermediate nodes are connected to another network. The transform matrix of this network (i.e., $\mathbf{G} \in \mathbb{R}^{m \times N}$) can be decoupled into two matrices: the one for the first stage of the network $\mathbf{G}_1 \in \mathbb{R}^{m \times N}$ and, the one for the rest, $\mathbf{G}_2 \in \mathbb{R}^{m \times m}$.

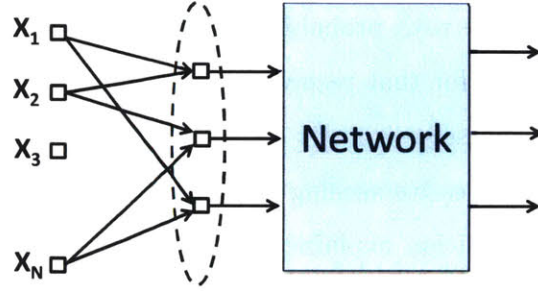


Figure 10-4: An example of a network topology where the network coding matrix can be designed to satisfy the RE condition. Hence, the network can transform sources with spatial dependencies which do not satisfy the RE condition, to a proper lower dimensional space.

That is, $\mathbf{G} = \mathbf{G}_2\mathbf{G}_1$. If sources choose network coding coefficients uniformly from $\{-1, 1\}$, as shown in reference [313], the matrix \mathbf{G}_1 satisfies the RE condition. By performing random linear network coding, the rank of \mathbf{G}_2 is m with high probability. Therefore, by using Lemma 10.2, \mathbf{G} satisfies the RE condition.

Step 4 (Decoding): Each receiver first uses a compressive sensing decoder as in equation (10.1) to reconstruct \mathbf{Y}^t for each $1 \leq t \leq m_1$ as follows:

$$\min_{\mu} \frac{1}{2m_2} \|\mathbf{Z} - \mathbf{G}\mathbf{B}\Psi\mu\|_{l_2}^2 + \xi \|\mu\|_{l_1}. \quad (10.6)$$

Therefore, for each source i , $\tilde{\mathbf{Y}}_i$ is obtained where $\tilde{\mathbf{Y}}_i = \mathbf{A}_i\Phi\mathbf{X}_i^s + \mathbf{U}_i$. \mathbf{U}_i is the error occurred in the reconstruction of \mathbf{Y}^t . Then, another compressive sensing decoder is used to reconstruct $\tilde{\mathbf{X}}_i$:

$$\min_{\eta_i} \frac{1}{2m_1} \|\tilde{\mathbf{Y}}_i - \mathbf{A}_i\Phi\eta_i\|_{l_2}^2 + \xi \|\eta_i\|_{l_1}. \quad (10.7)$$

The number of network uses of this proposed sensing-communication scheme is equal to $\mathcal{C}_{\text{use}} = \frac{m_1 m_2}{m}$.

$\mathbf{A}_i\Phi$ satisfies the RE condition with parameters $(\gamma_1 \lambda_1^2, \alpha)$ when \mathbf{A}_i satisfies the RE

condition with parameters (γ_1, α) and λ_1 is the minimum absolute eigenvalue of the matrix Φ (according to Lemma 10.2). The matrix $\mathbf{GB}\Psi$ satisfies the RE condition with parameters $(\gamma_2\lambda_2^2, \alpha)$ (i.e., in case 1, matrix $\mathbf{B}\Psi_S$ satisfies the RE condition with parameter (γ_2, α) and λ_2 is the minimum absolute eigenvalue of \mathbf{G}^S where \mathbf{G}^S contains columns of \mathbf{G} corresponding to active sources. In case 2, \mathbf{G} satisfies the RE condition with parameters (γ_2, α) and λ_1 is the minimum absolute eigenvalue of the matrix Ψ). Let λ_3 and λ_4 be maximum eigenvalues of matrices Ψ and Φ , respectively. By using these parameters, define the following constant:

$$c = \delta \frac{(\lambda_3\lambda_4)^2}{(\gamma_1\gamma_2)^2(\lambda_1\lambda_2)^4}, \quad (10.8)$$

where δ is a constant. The following Theorem illustrates how to choose m_1 and m_2 in the proposed joint source-channel-network coding scheme to reconstruct each source signal within an allowed distortion level:

Theorem 10.6. *A joint source-channel-network coding scheme of Algorithm 10.4 satisfies the distortion requirement of equation (10.4) with high probability if*

$$C_{use} = c \frac{k_1 k_2 \log(n) \log(N) \sigma^2}{m} \frac{1}{D}, \quad (10.9)$$

where $k_1 \ll n$ and $k_2 \ll N$ represent temporal and spatial sparsity of sources, σ^2 is the noise variance, D is the allowed reconstruction distortion level, and c is a constant defined in equation (10.8).

The proof of this Theorem is presented in Section 10.3. This theorem provides a way of choosing parameters m_1 and m_2 in Algorithm 10.4, where each source signal can be reconstructed within an allowed distortion level (equation (10.4)).

Remarks:

- Our proposed scheme is efficient: Sources use naturally occurring redundancies to perform joint source-channel coding. Therefore, there is no need of performing source coding and channel coding separately which would cause inefficiencies. Also, in the case that Ψ satisfies the RE condition, only a fraction

of sources are required to transmit, which provides additional power saving in transmissions.

- Our proposed scheme is broadly structure independent: Matrices Φ and Ψ are not required to be known at sources, but at receivers. It is important in the hardware implementation of sensors for different applications where measured signals can have different correlation structures ([333]).
- Our proposed scheme has a continuous rate-distortion performance: as it can be seen from Theorem 10.6, if the noise variance σ^2 increases, by using the same sensing-communication scheme, the reconstruction distortion would increase in a continuous manner (i.e., there is no performance drop off). In fact, our proposed scheme can be viewed as a *noise modulating* sensing-communication scheme.
- Our proposed scheme has low complexity decoders. Unlike Slepian-Wolf ([314]) decoders of information theoretic coding schemes, decoders used in our scheme are convex optimizations and therefore have lower complexity.
- It is insightful to compare the rate-distortion performance of our scheme with a naive information theoretic scheme ignoring source correlations. Following the rate-distortion Theorem ([316]), to represent a Gaussian random variable $\mathcal{N}(0, \sigma^2)$ within a distortion level D , one needs $R(D) = \frac{1}{2} \log(\frac{\sigma^2}{D})$ bits. Suppose sources transmit their samples ignoring correlations. Since at one network use, the network can deliver m linear combinations of source samples, $\frac{nN}{m}$ network uses are required. Moreover, since the noise variance is σ^2 but the allowed distortion level is D , roughly speaking, one needs at least $R(D)$ extra factor of transmissions. Overall,

$$\frac{nN}{m} \log\left(\frac{\sigma^2}{D}\right)$$

network uses are required. Note that, $k_1 \ll n$ and $k_2 \ll N$. Comparing this with the one of our proposed scheme (Theorem 10.6) illustrates that, although our scheme is power efficient, it performs closely to an optimal scheme where

only an informative subspace of source samples are transmitted.

- Note that, projection matrices A_i should be known at receivers. Transmitting these matrices may cause a rate overhead in a practical sensing-communication scheme. In practical implementations, since this matrix is a random matrix, each source can use a pseudo random matrix generated (i.e., pseudo-random bit sequence (PRBS) generators) and transmit its starting seed, instead of transmitting the whole matrix. This would have a negligible effect in the sensing-communication performance. Analyzing our proposed scheme with pseudo random matrices are beyond the scope of this chapter.
- Note that, the result of Theorem 10.6 is not deterministic and holds with high probability. There is always a small chance that, reconstructed signals have higher distortion than the allowed threshold.
- As the last point, our proposed framework allows to have *distributed sensing-communication*. That is, suppose there is one signal that receivers desire to reconstruct it. Each sensor has a low power budget to take all measurements and perform codings by itself. However, if several sensors (sources) are used, instead, each can take few measurements from the signal. By using our proposed scheme, receivers can reconstruct the original signal by using these measurements. In fact, our framework makes a virtual *strong* sensor by using several *weak* sensors.

10.3 Proofs

In this section, we present proofs.

10.3.1 Proof of Lemma 10.2

We first show that, C_1G satisfies the RE condition with parameters $(\gamma\lambda_1^2, \alpha)$ where G satisfies the RE condition with parameters (γ, α) and λ_1 is the minimum absolute eigenvalue of C_1 . We have,

$$\begin{aligned}
 \frac{1}{m} \|\mathbf{C}_1 \mathbf{G} \mathbf{Y}\|_{l_2}^2 &= \frac{1}{m} \mathbf{Y}^T \mathbf{G}^T \mathbf{C}_1^T \mathbf{C}_1 \mathbf{G} \mathbf{Y} \\
 &= \frac{1}{m} (\mathbf{G} \mathbf{Y})^T (\mathbf{C}_1^T \mathbf{C}_1) (\mathbf{G} \mathbf{Y}) \\
 &\stackrel{(1)}{\geq} \frac{1}{m} \lambda_1^2 \|\mathbf{G} \mathbf{Y}\|_{l_2}^2 \\
 &\stackrel{(2)}{\geq} \gamma \lambda_1^2 \|\mathbf{Y}\|_{l_2}^2
 \end{aligned}$$

Inequality (1) comes from a linear algebra lemma and inequality (2) comes from the definition of the RE condition for the matrix \mathbf{G} .

In the following, we prove the second half of this Lemma. We assume that, for $\mathbf{Y} \in \mathcal{C}(\mathcal{S}; \alpha)$, $\mathbf{Y} \mathbf{C}_2 \in \mathcal{C}(\mathcal{S}; \alpha)$. Also, suppose λ_2 is the minimum absolute eigenvalue of the matrix \mathbf{C}_2 . Therefore,

$$\begin{aligned}
 \frac{1}{m} \|\mathbf{G} \mathbf{C}_2 \mathbf{Y}\|_{l_2}^2 &= \frac{1}{m} \mathbf{Y}^T \mathbf{C}_2^T \mathbf{G}^T \mathbf{G} \mathbf{C}_2 \mathbf{Y} \\
 &= \mathbf{Y}^T \left(\frac{1}{m} \mathbf{C}_2^T \mathbf{G}^T \mathbf{G} \mathbf{C}_2 \right) \mathbf{Y} \\
 &= (\mathbf{C}_2 \mathbf{Y})^T \left(\frac{1}{m} \mathbf{G}^T \mathbf{G} \right) (\mathbf{C}_2 \mathbf{Y}) \\
 &\stackrel{(3)}{\geq} \gamma \|\mathbf{C}_2 \mathbf{Y}\|_{l_2}^2 \\
 &\stackrel{(4)}{\geq} \gamma \lambda_2^2 \|\mathbf{Y}\|_{l_2}^2
 \end{aligned}$$

Inequality (3) comes from the definition of the RE condition for the matrix \mathbf{G} . Note that, in this part, we use the assumption that, for $\mathbf{Y} \in \mathcal{C}(\mathcal{S}; \alpha)$, $\mathbf{Y} \mathbf{C}_2 \in \mathcal{C}(\mathcal{S}; \alpha)$. Inequality (4) comes from a linear algebra lemma. This completes the proof.

10.3.2 Proof of Theorem 10.6

In this part, we prove Theorem 10.6. Suppose $\mathbf{A}_i \Phi$ satisfies the RE condition with parameters $(\gamma_1 \lambda_1^2, \alpha)$ (i.e., \mathbf{A}_i satisfies the RE condition with parameters (γ_1, α) and λ_1 is the minimum absolute eigenvalue of the matrix Φ . Then, use Lemma 10.2).

Also, say the matrix $\mathbf{G}\mathbf{B}\Psi$ satisfies the RE condition with parameters $(\gamma_2\lambda_2^2, \alpha)$ (i.e., in case 1, matrix $\mathbf{B}\Psi_S$ satisfies the RE condition with parameter (γ_2, α) and λ_2 is the minimum absolute eigenvalue of \mathbf{G}^S where \mathbf{G}^S contains columns of \mathbf{G} , corresponding to active sources. In case 2, \mathbf{G} satisfies the RE condition with parameters (γ_2, α) and λ_1 is the minimum absolute eigenvalue of the matrix Ψ). Hence, by using the LASSO decoder of equation (10.6) and the error bound of equation (10.2), we have

$$\|\mu - \tilde{\mu}\|_{l_2}^2 \leq \frac{\delta_1}{\gamma_2^2\lambda_2^4} \frac{k_2 \log(N)}{m_2} \sigma^2. \quad (10.10)$$

Hence,

$$\begin{aligned} \|y_i - \tilde{y}_i\|_{l_2}^2 &\stackrel{(1)}{\leq} \|\mathbf{Y} - \tilde{\mathbf{Y}}\|_{l_2}^2 \\ &\stackrel{(2)}{=} \|\Psi\mu - \Psi\tilde{\mu}\|_{l_2}^2 \\ &\stackrel{(3)}{\leq} \lambda_3^2 \|\mu - \tilde{\mu}\|_{l_2}^2 \\ &\stackrel{(4)}{\leq} \frac{\delta_1\lambda_3^2}{\gamma_2^2\lambda_2^4} \frac{k_2 \log(N)}{m_2} \sigma^2. \end{aligned} \quad (10.11)$$

Inequality (1) is true for all vectors. Equality (2) uses the spatial correlation matrix property where $\mathbf{Y} = \Psi\mu$. Inequality (3) uses a linear algebra lemma. Inequality (4) uses equation (10.10).

Say $\sigma_u^2 = \|y_i - \tilde{y}_i\|_{l_2}^2$. For the decoder of equation (10.7), we can write,

$$\begin{aligned} \|x_i^t - \tilde{x}_i^t\|_{l_2}^2 &\stackrel{(5)}{\leq} \|\mathbf{X}_i - \tilde{\mathbf{X}}_i\|_{l_2}^2 \\ &\stackrel{(6)}{=} \|\Phi\eta - \Phi\tilde{\eta}\|_{l_2}^2 \\ &\stackrel{(7)}{\leq} \lambda_4^2 \|\eta - \tilde{\eta}\|_{l_2}^2 \\ &\stackrel{(8)}{\leq} \frac{\delta_2\lambda_4^2}{\gamma_1^2\lambda_1^4} \frac{k_1 \log(n)}{m_1} \sigma_u^2 \\ &\stackrel{(9)}{\leq} c \frac{k_1 k_2 \log(n) \log(N)}{m_1 m_2} \sigma^2 \end{aligned}$$

where $c = \delta_1 \delta_2 \frac{(\lambda_3 \lambda_4)^2}{(\gamma_1 \gamma_2)^2 (\lambda_1 \lambda_2)^4}$.

Inequality (5) is true for all vectors. Equality (6) uses the temporal correlation matrix property where $\mathbf{X}_i = \Phi \eta$. Inequality (7) uses a linear algebra lemma. Inequality (8) uses equation (10.2). Inequality (9) uses equation (10.11). Having the distortion requirement $\frac{1}{n} \|\mathbf{X}_i - \tilde{\mathbf{X}}_i\|_{l_2}^2 \leq D$ and $C_{\text{use}} = \frac{m_1 m_2}{D}$ completes the proof of Theorem 10.6.

Chapter 11

Nonuniform Sampling of Stochastic Signals

Sampling is essential in any digital system that interfaces with the analog world. All else being equal, it is desirable to minimize the number of samples while maintaining an acceptable reconstruction distortion. In some applications, minimizing the number of samples can be translated into having power-efficient sampling, since the power consumption at an analog-to-digital converter (ADC) is approximately proportional to its sampling rate [335]. Also, having fewer samples can increase the efficiency of other processing of these measurements. For example, if these samples should be transmitted to another place via a communication channel, having fewer samples will improve power and bandwidth efficiencies.

One can view *sampling* as a *query* to obtain information from a signal or a function that can only be measured remotely. We get a sample of this function at an arbitrary time when we *query* for it over a communication medium. Here, a portion of the operational cost (e.g., power) is proportional to the number of samples that we acquire. Hence, again it is desirable to minimize the number of samples taken.

A uniform sampling at the Nyquist rate of the signal may cause some redundant samples, since the global signal bandwidth may not be a good measure of local variations of the signal. Although traditional nonuniform sampling schemes (e.g., [336, 337]) deal with this problem, they have certain limitations. Firstly, they are mostly designed to operate under specific conditions for restrictive signal models (e.g., [338–342]) and, secondly, sampling times (i.e., time stamps) must be stored/transmitted to be used in the reconstruction process. This may cause power/bandwidth inefficiencies in sampling/communication procedures.

In this chapter, we introduce a new framework for an adaptive nonuniform sampling scheme (see Figure 11-1). The key idea of this framework is that *time increments between samples are computed by using a function of previously taken samples*. Therefore, keeping sampling times (time stamps), except initialization times, is not necessary. The function by which sampling time intervals is computed is called the *sampling function*. The aim of this sampling framework is to have a balance between the reconstruction distortion and the average sampling rate. We refer to this sampling framework as *Time-stampless Adaptive Nonuniform Sampling* (TANS). The TANS

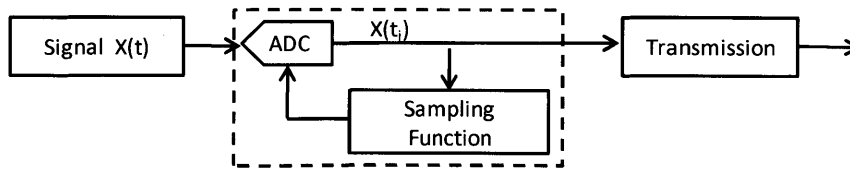


Figure 11-1: A schematic view of the TANS framework: sampling times are determined by a function of m most recently taken samples. Hence, it is not necessary to save sampling times (time stamps) for use in the reconstruction process.

concept can be applied on continuous- or discrete-time signals, and the design and analysis can be based on deterministic or stochastic models.

The TANS framework is described in general terms in Section 11.1. Section 11.2 then formalizes the problem setup for this chapter, where we focus on discrete-time signals and stochastic models. A greedy method is developed in Section 11.3, and a method based on dynamic programming (DP) is developed in Section 11.4. Simulation results are provided in Section 11.5, and some proofs are relegated to Section 11.6.

11.1 TANS Framework

In this section, we introduce the TANS framework and sketch design approaches that will be developed in more detail in later sections. Fix some nonnegative integer m and suppose the i th sample of signal $X(t)$ is taken at time t_i . We take the $(i+1)$ st sample after a time increment of

$$T_i \triangleq f(t_{i-m+1}, \dots, t_{i-1}, t_i, X(t_{i-m+1}), \dots, X(t_{i-1}), X(t_i)),$$

where f is called the sampling function. This makes the sampling rate adapt to local characteristics of the signal. Since the time increment is a function of the m most recently taken samples, we say the *order* of the sampling function f is m . The sampling is nonuniform except in the trivial cases when f is a constant-valued function (e.g., $m = 0$). Some initialization of the first m sampling times is necessary, but the effect of this initialization on the rate is amortized.

The sampling function is known at the reconstruction side. Assuming that the *state*

$$S_i \triangleq \{t_{i-m+1}, \dots, t_{i-1}, t_i, X(t_{i-m+1}), \dots, X(t_{i-1}), X(t_i)\}$$

is also known at the reconstruction side when reconstructing $X(t)$ on $[t_i, t_{i+1}]$, there is *no need for the sampling times (time stamps) to be transmitted*. These times can be computed by using the sampling function and previously taken samples:

$$t_{i+1} = t_i + f(S_i).$$

This type of synchronization in an adaptive system without explicit communication is often called backward adaptation [343]. In a practical setting involving both sampling and quantization, backward adaptivity requires using the quantized values to drive the adaptation [344–347]. Here, to maintain focus on sampling rate and adaptation of sampling increments, we do not explicitly include quantization effects. Note that, while the sampling time selection is causal, the reconstruction method can be causal or non-causal.

The aim of TANS is to balance between the average sampling rate and the reconstruction distortion. This objective is different from the one considered in change point analysis [348] or active learning [349]. There, the objective is to find points of the signal at which statistical behaviors of the signal change, by causal or non-causal sampling, respectively.

Suppose $\hat{X}(t)$ is the reconstructed signal computed by some reconstruction method. For the case of discrete time and a stochastic signal model, define $d(S_i, T_i)$ as the expected reconstruction distortion over samples from time $t_i + 1$ until time $t_{i+1} - 1$. That is,

$$d(S_i, T_i) = \mathbb{E}_{\mathcal{X}} \left[\sum_{t=t_i+1}^{t_i+T_i-1} D(X(t), \hat{X}(t)) \right],$$

where \mathcal{X} is the known probabilistic model of the signal $X(t)$ and $D(X(t), \hat{X}(t))$ represents the distortion at time t .¹ In this chapter, we use mean-squared error (MSE)

¹An analogous formulation for continuous time would replace the sum with an integral over

as our distortion metric. Note that, at times t_i and t_{i+1} , the reconstruction distortion is zero since exact sample values are known at these times. Note that, for a given state S_i , $d(S_i, T_i)$ is an increasing function with respect to T_i . On the other hand, the greater the next sampling step, the larger the rate benefit. Hence, a rate penalty can be defined as $a(S_i, T_i) = \rho/f(S_i) = \rho/T_i$, where ρ is a rate award parameter. We define the cost of each sampling state as the sum of the expected reconstruction distortion and the rate penalty, that is,

$$c(S_i, T_i) \triangleq d(S_i, T_i) + a(S_i, T_i). \quad (11.1)$$

The overall cost of the sampling process is the sum of different sampling state costs, that is, $\sum_i c(S_i, T_i)$.

Finding an appropriate sampling function for TANS depends on requirements such as average sampling rate, maximum distortion, etc. In this chapter, we investigate two general approaches to computing appropriate sampling functions for given sampling setups: *greedy methods* and *dynamic programming (DP) methods*.

In greedy methods, a sampling function at state S_i chooses the next sampling increment T_i to minimize the sampling state cost $c(S_i, T_i)$. As depicted in Figure 11-2, for any given state S_i , $d(S_i, T_i)$ is an increasing function of T_i , while $a(S_i, T_i)$ is a decreasing function. Therefore, there is a trade-off between the sampling rate and the expected reconstruction distortion. A greedy method balances this trade-off by choosing the state cost minimizer as the next sampling increment. In certain cases, greedy sampling schemes can perform closely to an optimal scheme.

Since $c(S_i, T_i)$ depends on the current sampling state S_i , a greedy sampling function does not take into account characteristics of the next sampling state. Intuitively, the larger the sampling increment T_i at the sampling state S_i , the lower the *quality* of the next sampling state. Hence, in general, greedy methods are not optimal sampling schemes considering the overall sampling cost $\sum_i c(S_i, T_i)$ as a comparison measure. We consider effects of the next sampling states' *quality* in DP methods.

$t \in [t_i, t_{i+1}]$. Without a stochastic model, a maximum error criterion could be used.

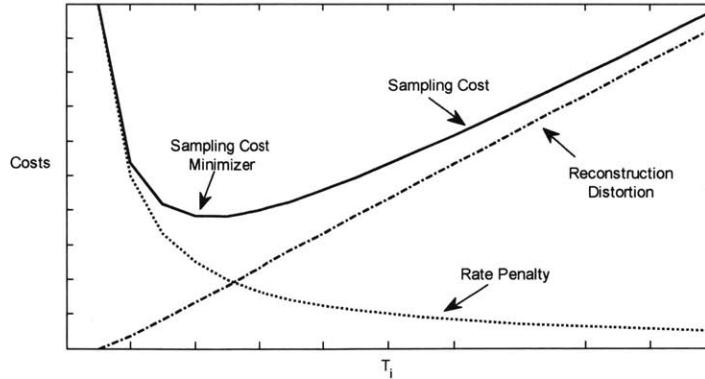


Figure 11-2: Demonstrating the behavior of different parts of sampling state cost $c(S_i, T_i) = d(S_i, T_i) + a(S_i, T_i)$, where $d(S_i, T_i)$ is the reconstruction distortion and $a(S_i, T_i)$ is the rate penalty function.

We show that, an exact Bellman-Ford equation (BFE) can be written and solved for some sampling setups. For those cases, the solution of BFE provides an optimal sampling function which minimizes the overall sampling cost. In cases where solving the BFE is not practically feasible (for instance because the number of possible sampling states is large or exact sampling states are not known), we propose sampling functions based on approximate dynamic programming (ADP) algorithms. Sampling functions derived by greedy methods can be used in ADP-based sampling methods. In fact, greedy methods can be viewed as DP-based methods with a unit time horizon.

In this chapter, we consider two examples of stochastic signals: stationary signals and Markovian signals. Unlike stationary signals, Markovian processes have sudden changes in their statistical properties based on an underlying hidden Markov chain. Note that, the general sampling framework can be applied on other signal models.

11.2 Problem Setup, Signal Models and Background Results

In this section, we first present the problem setup of TANS. Then, we introduce signal models considered in this chapter. At the end of this section, a *generalized linear prediction filter* is proposed, which linearly predicts the future samples of a

stationary process by using a set of nonuniform samples from its past.

11.2.1 Problem Setup

Consider a discrete-time signal $X(t)$. TANS with order m is used to take samples from this signal. Hence, the next sampling increment at time t_i is a function of m most recently taken samples at that time (i.e., $T_i = f(S_i)$). The set of samples taken using the sampling function $f(\cdot)$ is denoted by $\Gamma_{X(t)}^f$. A function $g(\cdot)$ is used to reconstruct the original signal from its samples, that is, $\hat{X}(t) = g(\Gamma_{X(t)}^f)$. The reconstruction error signal is $e(t) = X(t) - \hat{X}(t)$. The overall sampling cost is the sum of sampling state costs over different states, that is, $c_{\text{tot}}(f, g) = \sum_i c(S_i, T_i)$.

A system optimization problem under the TANS framework can be stated as follows:

Definition 11.1 (Optimal TANS sampling problem). For a class of signals $X(t)$, a given reconstruction function $g(\cdot)$, and an order m , a sampling function $f^*(\cdot)$ is desired to minimize the overall expected sampling cost:

$$f^* = \arg \min_f c_{\text{tot}}(f, g). \quad (11.2)$$

The resulting cost is denoted $c_{\text{tot}}^* \triangleq c_{\text{tot}}(f^*, g)$.

11.2.2 Signal Models

Here, we consider the following signal models:

- Case 1: an autoregressive signal with memory of one (i.e., AR(1)):

$$X(t+1) = \alpha X(t) + Z(t+1), \quad (11.3)$$

where $Z(t+1)$ is a Gaussian noise with zero mean. If the power of the signal is assumed to be one, the noise variance is $1 - \alpha^2$.

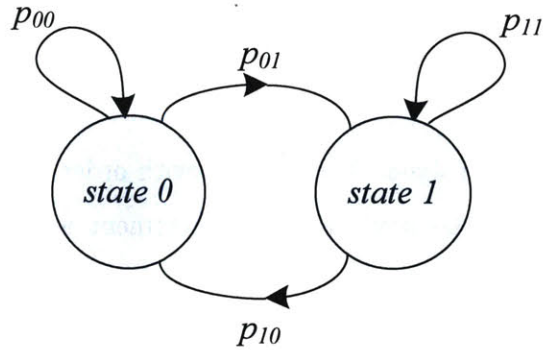


Figure 11-3: A hidden Markov chain considered in Markovian signal model of equation (11.4).

- Case 2: a Markovian signal:

$$X(t+1) = \alpha_{\theta_t} X(t) + Z_{\theta_t}(t+1), \quad (11.4)$$

where θ_t represents the state of a hidden Markov chain (MC) with state transition probabilities depicted in Figure 11-3. At time t , if the MC is at state 0, $\theta_t = 0$; otherwise, $\theta_t = 1$. Depending on the value of θ_t , the signal is generated by a first-order AR model with parameter α_{θ_t} and the noise variance $1 - \alpha_{\theta_t}^2$. Note that, in this model, unlike the previous case, the coefficient of the AR model has a sudden change in time depending on the state of the underlying hidden Markov chain.

11.2.3 Generalized Linear Prediction Filter

Suppose $X(t)$ is a stationary signal. Assume we have m samples of $X(t)$ at times $\{t_{i-m+1}, \dots, t_i\}$. Our aim is to predict linearly $X(t)$, where $t = t_i + T$ for some $T \geq 1$, by using these known samples and their times so that the expected mean square error is minimized (MMSE predictor).

Consider $\hat{X}(t)$ as a predicted value of $X(t)$ by using these m sample values and their times. The prediction error is $e(t) = X(t) - \hat{X}(t)$. Define $\tau_k \triangleq t - t_{i-k}$ for

$0 \leq k \leq m - 1$. We want to find optimal linear prediction weights $w_{\tau_0}, \dots, w_{\tau_{m-1}}$ so that the prediction error power is minimized:

$$\begin{aligned} \min_{w_{\tau_0}, \dots, w_{\tau_{m-1}}} \quad & E[|e(t)|^2 | S_i] \\ \text{subject to} \quad & \hat{X}(t) = \sum_{k=0}^{m-1} w_{\tau_k} X(t_{i-k}). \end{aligned} \quad (11.5)$$

A solution to this linear optimization is referred to as $w_{\tau_i}^*$, for $0 \leq i \leq m - 1$. Note that, unlike a regular linear prediction filter ([350]), a generalized linear prediction filter predicts $X(t)$ by using a set of nonuniform samples.

Since $X(t)$ is a real signal, the auto-correlation function of $X(t)$ can be written as

$$r(i) = E[X(t)X(t - i)]. \quad (11.6)$$

To simplify notations, we define the following matrices:

$$\begin{aligned} \boldsymbol{\tau} &= [\tau_0, \dots, \tau_{m-1}]^T \\ \mathbf{p} &= [r(-\tau_0), \dots, r(-\tau_{m-1})]^T \\ \mathbf{w}_{\boldsymbol{\tau}}^* &= [w_{\tau_0}^*, \dots, w_{\tau_{m-1}}^*]^T. \end{aligned} \quad (11.7)$$

Also, a $m \times m$ auto-correlation matrix \mathbf{R} is defined whose component in the i th row and j th column is $r(\tau_i - \tau_j)$.

The following theorem provides optimal weights for the generalized linear prediction filter:

Theorem 11.2.

$$\mathbf{p} = \mathbf{R}\mathbf{w}_{\boldsymbol{\tau}}^*.$$

Proof. See Section 11.6.1. □

For $X(t)$ with zero mean, the variance of the prediction error is defined as follows:

$$\sigma_{e^*}^2(S_i, T) = E[|e^*(t)|^2 | S_i].$$

The following theorem provides a way to compute the variance of the prediction error for a generalized linear prediction filter:

Theorem 11.3. $\sigma_{e^*}^2(S_i, T) = r(0) - \mathbf{p}^T \mathbf{w}_r^*$.

Proof. See Section 11.6.2. □

11.3 Greedy TANS

In this section, we investigate greedy sampling functions for various signal models. In Section 11.4, we evaluate performance of these schemes and compare them with some other schemes, including uniform sampling setups. In all of these sampling schemes, the reconstruction function is assumed to be a generalized linear prediction filter, introduced in Section 11.2.3. Note that, it is a causal reconstruction function.

In greedy methods, a sampling function is computed as follows:

$$\begin{aligned} T_i &= \arg \min_T c(S_i, T) & (11.8) \\ &\text{subject to } T \geq 1, \end{aligned}$$

where $f(S_i) = T_i$. The sampling function of a greedy method depends on the current sampling state and does not take into account characteristics of the next sampling states. Intuitively, the larger the sampling increment T_i at the sampling state S_i , the lower the *quality* of the next sampling state. Therefore, we have a trade-off between the sampling rate award of the current state and the sampling cost of the next state (Figure 11-2). Hence, greedy methods usually are not optimal solutions of a TANS sampling problem presented in Definition 11.1. However, these greedy solutions can be used to approximate optimal solutions, since the latter may have high computational complexity.

Here, we investigate the greedy sampling function of equation (11.8) for two signal classes described in Section 11.2.2. We analyze their sampling rate versus the expected reconstruction distortion in Theorems 11.4 and 11.8. Simulation results of proposed schemes are shown in Section 11.5.

11.3.1 Greedy TANS for Autoregressive Signals

In this section, we consider greedy sampling functions for an AR(1) signal model described in (11.3).

Suppose $f_{\text{greedy}}^*(S_i)$ is an optimal greedy sampling function. Also, suppose d_{tot}^* is the expected reconstruction distortion per sample corresponding to this sampling function. The following theorem introduces an optimal greedy sampling function and its expected reconstruction distortion per sample.

Theorem 11.4. *For an AR(1) signal with parameter α described in equation (11.3), over a large enough time interval $[0, T_{\text{tot}}]$, an optimal greedy sampling function is*

$$f_{\text{greedy}}^*(S_i) = T^* \tag{11.9}$$

where $T^* = \arg \min_T \sum_{j=1}^{T-1} (1 - \alpha^{2j}) + \rho/T$. This sampling function leads to the following average reconstruction distortion per sample:

$$d_{\text{tot}}^* = \frac{1}{T^*} \sum_{j=1}^{T^*-1} 1 - \alpha^{2j}. \tag{11.10}$$

Proof. See Section 11.6.3. □

Note that, for this signal model, an optimal greedy sampling function yields to a uniform sampling with the sampling rate $1/T^*$. In the following corollary, we present a formula to compute T^* .

Corollary 11.5. Suppose T_{root} is a solution of the following equation:

$$(1 - \alpha^{2T}) - \frac{\rho}{T(T+1)} = 0 \quad (11.11)$$

where $0 < \alpha < 1$ and $\rho > 0$. If $(1 - \alpha^2) < \frac{\rho}{2}$, then

- (i) T_{root} is unique; and
- (ii) $T^* = \lfloor T_{\text{root}} \rfloor$ or $T^* = \lfloor T_{\text{root}} \rfloor + 1$.

Proof. See Section 11.6.4. □

We shall illustrate this Corollary by simulation in Section 11.5. Also, arguments of Theorem 11.4 can be extended for a general stationary signal.

For an AR signal, since statistical properties of the signal do not vary in time, there is no rate adaption with respect to sample values. This is not the case for Markovian signals. We investigate the greedy TANS framework for Markovian signals in the next section; a rate adaptation with respect to sample values would be helpful to minimize the sampling cost and leads to a nonuniform sampling scheme.

11.3.2 Greedy TANS for Markovian Signals

Consider a Markovian signal described by equation (11.4), where θ_t represents the state of a hidden underlying Markov chain depicted in Figure 11-3. In this section, for simplicity we assume the MC is symmetric (i.e., $p_{01} = p_{10}$). However, all arguments can be extended for a general MC. We also assume that α_0 and α_1 are known. However, the state of the Markov chain (i.e., θ_t) is unknown and needed to be estimated by using the taken samples. We use a generalized linear prediction filter for the reconstruction. Note that, if α_0 and α_1 are also unknowns, one can learn these parameters at the beginning of the process by taking more samples. Then, our proposed scheme can be applied for the rest of the process.

Extending the previous notation, define θ_{S_i} as the state of the MC during the sampling state S_i . If during S_i the MC state stays at zero, $\theta_{S_i} = 0$. Similarly, if the

MC state stays at one, $\theta_{S_i} = 1$. Otherwise, if there is an MC transition within this sampling state, $\theta_{S_i} = 2$. We assume that θ_{S_i} is unknown and needs to be estimated by using the taken samples. The estimated value of θ_{S_i} is referred by $\hat{\theta}_{S_i}$. The error probability of this estimation is referred by $P_e(S_i) = Pr(\hat{\theta}_{S_i} \neq \theta_{S_i})$.

Algorithm 11.6. A greedy sampling function for the considered Markovian signal has the following steps:

- Step i,0: Compute $\hat{\theta}_{S_i}$ and $P_e(S_i)$.
- Step i,1: Compute $T_i = \arg \min_T c(S_i, T | \hat{\theta}_{S_i})$, where $c(S_i, T | \hat{\theta}_{S_i})$ is the sampling state cost given $\hat{\theta}_{S_i}$ (see (11.14) and (11.15)).
- Step i,2: Take a sample at time $t_i + T_i$.
- Step i,3: Compute S_{i+1} . Repeat.

For simplicity, we assume that the sampling increment T_i is small enough that the probability of having more than one MC transition is negligible. Formally, we have the following assumptions.

Assumption 11.7. The transition probabilities p_{01} and p_{10} are small and the sampling increments are bounded, so

$$0 \leq 1 - p_{00}^{T_i} \lambda 1/2, \tag{11.12}$$

$$0 \leq 1 - p_{11}^{T_i} \lambda 1/2$$

for all T_i .

The sampling state cost, $c(S_i, T)$, can be conditioned on the value of $\hat{\theta}_{S_i}$ as follows:

$$\begin{aligned}
 & c(S_i, T | \hat{\theta}_{S_i} = 0) \\
 & \stackrel{(a)}{=} d(S_i, T | \hat{\theta}_{S_i} = 0) + a(S_i, T | \hat{\theta}_{S_i} = 0) \\
 & \stackrel{(b)}{=} (1 - P_e(S_i)) \left\{ \underbrace{p_{00}^{T-1}}_{\text{I}} \underbrace{\left(\sum_{l=1}^{T-1} (1 - \alpha_0^{2l}) \right)}_{\text{II}} \right. \\
 & \quad \left. + \sum_{j=0}^{T-2} \underbrace{p_{00}^j p_{01}}_{\text{III}} \left(\underbrace{\sum_{l=1}^j (1 - \alpha_0^{2l})}_{\text{IV}} + \underbrace{(T-1-j)\sigma_{\max}^2}_{\text{V}} \right) \right\} \\
 & \quad + \underbrace{P_e(S_i)(T-1)\sigma_{\max}^2}_{\text{VI}} + \underbrace{\frac{\rho}{T}}_{\text{VII}}. \tag{11.13}
 \end{aligned}$$

Equality (a) follows from the definition of the sampling state cost in equation (11.1). To compute the expected distortion, we consider two cases: when the estimation is correct (with probability $1 - P_e(S_i)$) and when the estimation is incorrect. If the estimation is correct, with probability p_{00}^{T-1} , there is no MC transition over the interval of length $T - 1$ (part I), and the expected distortion is $\sum_{l=1}^{T-1} (1 - \alpha_0^{2l})$ (part II). A MC transition happens at time $t_i + j + 1$, with probability $p_{00}^j p_{01}$ (part III). The expected distortion in this case is $\sum_{l=1}^j (1 - \alpha_0^{2l}) + (T - 1 - j)\sigma_{\max}^2$ (parts IV and V), where σ_{\max}^2 is the maximum prediction error variance (in this example, $\sigma_{\max}^2 = 1$). If the estimation process fails (with probability $P_e(S_i)$), a maximum prediction error variance σ_{\max}^2 occurs over all $(T - 1)$ samples (part VI). Finally, part VII shows the rate award.

If the probability of having a Markov chain transition over the time interval of length T_i is negligible, for all T_i , parts (I-V) of equation (11.13) can be simplified further, by only considering the case when there is no MC transition over this interval. This can happen when $p_{00}^{T_i} \approx 1$, for all T_i . Therefore, under this assumption, equation (11.13) can be approximated as follows:

$$\begin{aligned}
 c(S_i, T | \hat{\theta}_{S_i} = 0) & \\
 \approx (1 - P_e(S_i)) \sum_{\ell=1}^{T-1} (1 - \alpha^{2\ell}) & \\
 + P_e(S_i)(T - 1)\sigma_{\max}^2 + \frac{\rho}{T}. & \quad (11.14)
 \end{aligned}$$

The sampling state cost function conditioned on $\hat{\theta}_{S_i} = 1$ (i.e., $c(S_i, T | \hat{\theta}_{S_i} = 1)$) can be written similarly.

Finally, for the case $\hat{\theta}_{S_i} = 2$, we assume that the prediction variance is the maximum prediction error variance σ_{\max}^2 :

$$c(S_i, T | \hat{\theta}_{S_i} = 2) = (T - 1)\sigma_{\max}^2 + \frac{\rho}{T}. \quad (11.15)$$

We analyze the performance of the proposed greedy sampling scheme in the following theorems. We derive upper and lower bounds for average sampling rate and expected reconstruction distortion per sample of the proposed sampling scheme. Simulation results for this sampling scheme are presented in Section 11.5.

Before presenting theorems, we introduce some notations. Suppose that, for all S_i , $P_e^{\text{low}} \leq P_e(S_i) \leq P_e^{\text{up}}$. By considering an upper bound on $P_e(S_i)$, we define

$$\begin{aligned}
 T_0^{\text{low}} = \arg \min_T \quad & (1 - P_e^{\text{up}}) \sum_{\ell=1}^{T-1} (1 - \alpha_0^{2\ell}) \\
 & + P_e^{\text{up}}(T - 1)\sigma_{\max}^2 + \frac{\rho}{T}. \quad (11.16)
 \end{aligned}$$

T_0^{up} is defined similarly by considering a lower bound on $P_e(S_i)$. Analogously, T_1^{up} and T_1^{low} can be defined.

Also, d_0^{up} , an upper bound on the expected reconstruction distortion per sample given $\hat{\theta}_{S_i} = 0$ is defined as follows:

$$d_0^{\text{up}} = \frac{1}{T_0^{\text{up}}} \left\{ (1 - P_e^{\text{up}}) \sum_{\ell=1}^{T_0^{\text{up}}-1} (1 - \alpha_0^{2\ell}) + P_e^{\text{up}}(T_0^{\text{up}} - 1)\sigma_{\max}^2 \right\}.$$

Quantities d_0^{low} , d_1^{up} and d_1^{low} are defined similarly.

The following theorem provides analytical upper and lower bounds on the average sampling rate and the expected reconstruction distortion of the greedy sampling scheme introduced in Algorithm 11.6.

Theorem 11.8. *Consider a Markovian signal defined in equation (11.4) over a large enough time interval $[0, T_{tot}]$. Under Assumption 11.7, an achievable rate-distortion pair (R, D) of the greedy sampling scheme of Algorithm 11.6 can be bounded as follows:*

$$\frac{1}{2T_0^{up}} + \frac{1}{2T_1^{up}} \leq R \leq \frac{1}{2T_0^{low}} + \frac{1}{2T_1^{low}} \quad (11.17)$$

$$\frac{d_0^{low}}{2} + \frac{d_1^{low}}{2} \leq D \leq \frac{d_0^{up}}{2} + \frac{d_1^{up}}{2}. \quad (11.18)$$

Proof. See Section 11.6.5. □

Similarly to Corollary 11.5, T_0^{low} , T_0^{up} , T_1^{low} and T_1^{up} can be calculated by finding roots of some equations. For example:

Corollary 11.9. Suppose T_{root} is a solution of the following equation:

$$(1 - P_e^{up})(1 - \alpha_0^{2T}) + P_e^{up}(T - 1)\sigma_{max}^2 - \frac{\rho}{T(T + 1)} = 0 \quad (11.19)$$

where $0 < \alpha < 1$ and $\rho > 0$. If $(1 - P_e^{up})(1 - \alpha_0^2) < \frac{\rho}{2}$, then,

- (i) T_{root} is unique; and
- (ii) $T_0^{low} = \lfloor T_{root} \rfloor$ or $T_0^{low} = \lfloor T_{root} \rfloor + 1$.

The proof of Corollary 11.9 is similar to the one of Corollary 11.5. Also, a similar corollary can be stated for T_0^{up} , T_1^{low} and T_1^{up} .

In Theorem 11.8, the performance of the proposed sampling scheme (i.e., its average sampling rate and the expected reconstruction distortion) is bounded. However, it is insightful to compare its performance to a genie-aided sampling scheme where the state of the underlying Markov chain is known. For a genie-aided scheme, $P_e^{low} = P_e^{up} = 0$, and therefore upper and lower bounds of Theorem 11.8 match:

Corollary 11.10. For a genie-aided sampling scheme, the following rate-distortion pair is achievable:

$$R = \frac{1}{2T_0^{\text{genie}}} + \frac{1}{2T_1^{\text{genie}}}, \quad (11.20)$$

$$D = \frac{d_0^{\text{genie}}}{2} + \frac{d_1^{\text{genie}}}{2}, \quad (11.21)$$

where

$$T_0^{\text{genie}} = \arg \min_T \sum_{\ell=1}^{T-1} (1 - \alpha_0^{2\ell}) + \frac{\rho}{T},$$

$$d_0^{\text{genie}} = \frac{1}{T_0^{\text{genie}}} \sum_{\ell=1}^{T_0^{\text{genie}}-1} (1 - \alpha_0^{2\ell}),$$

and T_1^{genie} and d_1^{genie} are defined similarly.

Corollary 11.10 can be derived by using $P_e^{\text{low}} = P_e^{\text{up}} = 0$ in Theorem 11.8.

In the first step of this sampling function, we need to estimate θ_{S_i} by using m most recently taken samples of the signal and compute its probability of error $P_e(\hat{\theta}_{S_i})$. In the following, we give an example of such a scheme.

Example 11.11. Suppose $m = 2$ and we use a maximum likelihood estimator. If $\theta_{S_i} = 0$, the probability distribution of $X(t_i)$ is $\mathcal{N}(\alpha_0^{T_i-1} X(t_{i-1}), 1 - \alpha_0^{2T_i-1})$, where $\mathcal{N}(\mu, \sigma^2)$ represents a Gaussian distribution with mean μ and variance σ^2 . The prior probability of this event is $p_{00}^{T_i-1}$. Similarly, if $\theta_{S_i} = 1$, $X(t_i)$ is distributed by $\mathcal{N}(\alpha_1^{T_i-1} X(t_{i-1}), 1 - \alpha_1^{2T_i-1})$. The prior probability of this event is $p_{11}^{T_i-1}$. Otherwise, the distribution of $X(t_i)$ is $\mathcal{N}(\alpha_0^j \alpha_1^{T_i-1-j} X(t_{i-1}), 1 - \alpha_0^{2j} \alpha_1^{2(T_i-1-j)})$ with the prior probability $\binom{T_i-1}{j} p_{00}^j p_{11}^{T_i-1-j}$, for $1 \leq j \leq T_i-1$. Each of these events corresponds to the case $\theta_{S_i} = 2$ (i.e., there is a transition within the sampling state.). Therefore, by having the observed value of $X(t_i)$ and using these distributions, a maximum likelihood estimator can estimate $\hat{\theta}_{S_i}$ and compute its error probability.

Remarks:

1. If the state space is not large, computations can be performed off-line and results can be used in the sampling function.

2. When the order of the sampling function is large and/or autocorrelation coefficients change continuously, a maximum likelihood estimator may not be practically interesting. In these cases, we can use previously taken samples within a window of size W from the last sample (i.e., all taken samples from time $t_i - W + 1$ to the time t_i) to update or estimate autocorrelation coefficients to use in the sampling function. The quality of this estimation process depends on the window size W , the variation rate of autocorrelation coefficients, and the technique used. Two possible methods for estimating autocorrelation coefficients are as follows:

- A gradient-based method. Suppose at the sampling state S_{i-1} , the set of estimated autocorrelation coefficients is $\{\hat{r}(j) : j \geq 1\}$. By taking a sample at time n_i , these coefficients are updated as follows:

$$\hat{r}(j) := \hat{r}(j) + \gamma(X(t_i)X(t_i - j) - \hat{r}(j)) \quad (11.22)$$

for all possible j 's, where $:=$ represents an update sign, and $\gamma > 0$ is a gradient step size. This gradient-based update method can be useful when W is not large.

- A window-based method. If the window size W is large and there are enough known samples within the window, an empirical value for each autocorrelation coefficient can be computed.

The estimated autocorrelation coefficients $\{\hat{r}(j) : j \geq 0\}$ can be used in the generalized linear prediction filter in order to design a sampling function similar to the one of Algorithm 11.6.

11.4 Dynamic Programming-Based TANS

In greedy TANS, sampling functions are derived based on minimizing the sampling cost at each sampling state. Therefore, in general, greedy methods may not provide optimal sampling functions with respect to the overall sampling cost.

We consider *quality* of next sampling states in dynamic programming-based TANS methods.

For the sampling state S_i , a cost-to-go function $J_f(S_i)$ is defined as follows:

$$\begin{aligned} J_f(S_i) &= c(S_i, T_i) + \beta c(S_{i+1}, T_{i+1}) + \cdots \\ &= \sum_{j \geq i} \beta^{j-i} c(S_j, T_j). \end{aligned} \tag{11.23}$$

In this setup, $0 < \beta < 1$ is called a *discount factor*. An interpretation of this factor is that, the cost of the next sampling state is less important for the current state policy by a factor of β . Another interpretation of this factor is that, the process may be ended at each sampling state with probability $1 - \beta$. An optimal cost-to-go function $J(S_i)$ of the sampling state S_i is defined as follows:

$$J(S_i) = \inf_f J_f(S_i). \tag{11.24}$$

A Bellman-Ford equation (BFE) can be written for this problem by using cost-to-go functions of different sampling states as follows:

$$J(S_i) = \inf_{T_i} (c(S_i, T_i) + \beta E[J(S_{i+1})]). \tag{11.25}$$

A solution of this Bellman-Ford equation (BFE) is an optimal solution for the sampling problem presented in Definition 11.1 when the reconstruction function is causal. We investigate this problem for various signal models and sampling setups in this section. For some cases where the number of sampling states is not large, an optimal sampling function can be derived. In other cases where finding this optimal solution is computationally difficult, we propose sampling functions based on approximate dy-

dynamic programming (ADP) algorithms. We define a *quality* function $q(S_i)$ for each sampling state S_i . A greedy solution is used to define this quality function. Then, a sampling function can be computed as follows:

$$T_i = \arg \inf_T c(S_i, T) + \beta E[q(S_{i+1})] \quad (11.26)$$

subject to $T \geq 1$.

In this setup, we consider quality effects of just one sampling state ahead. Also, note that the solution of the BFE is optimal if the reconstruction function $g(\cdot)$ is causal. It is a necessary assumption of the dynamic programming setup to have a separation of sampling state costs at different stages.

11.4.1 An Online Source Coding Scheme Based on TANS

In this section, we consider a sampling problem where an exact DP-based solution can be derived as a solution of the BFE equation (11.25). Consider a Markovian signal with an underlying hidden Markov chain depicted in Figure 11-3. For simplicity, suppose $p_{01} = \epsilon_0 \lambda 1$ and $p_{10} = \epsilon_1 \lambda 1$. Other regimes of transition probabilities can be analyzed similarly. Suppose $X(t)$ is a binary signal generated by this underlying hidden Markov chain so that, at state 0, $X(t) = 0$, and at state 1, $X(t) = 1$. For the reconstruction, we use a causal function that selects a most-probable binary sequence to fill missing places. Hamming distance is used as an error measure. We call this problem an *online source coding* problem, since it does not have a compression delay as in an unconstrained block source coding scheme. In this problem, TANS can provide some compression gain without having any delay. This problem can be extended to Markov chains with more states and different regimes of transition probabilities.

In this sampling problem, there are two different sampling states: If $X(t_i) = 0$, then, $S_i = 0$; otherwise, $S_i = 1$. Suppose $f(\cdot)$ is the sampling function (i.e., $T_i = f(S_i)$). Suppose $f(S_i = 0) = T_0$ and $f(S_i = 1) = T_1$. Hence, T_0 and T_1 represent sampling steps at different sampling states.

Since $\epsilon_0, \epsilon_1 \ll 1$, the reconstruction method would choose all-0 and all-1 sequences when $S_i = 0$ and $S_i = 1$, respectively. Suppose $S_i = 0$. Hence, the next sample is taken at time $t_i + T_0$. An error happens if there are one or more Markov chain transitions over the time interval $[t_i + 1, t_i + T_0 - 1]$. To simplify our analysis, we only consider first-order error terms (i.e., we assume at most one transition happens over a time interval $[t_i + 1, t_i + T_0 - 1]$). To have this simplifying assumption, we need to restrict the sampling increments such that $\max(T_0, T_1) \ll \min\{\frac{1}{\epsilon_0}, \frac{1}{\epsilon_1}\}$. Therefore, the probability of having more than one Markov chain transition over a time interval of a length T_0 or T_1 is negligible.

By considering Hamming distance as a distortion measure, the sampling state cost at $S_i = 0$ can be written as follows:

$$c(S_i, T_0) = \sum_{j=1}^{T_0-1} (1 - \epsilon_0)^{j-1} \epsilon_0 (T_0 - j) + \frac{\rho}{T_0}. \quad (11.27)$$

Note that, $(1 - \epsilon_0)^{j-1} \epsilon_0$ is the probability of not having a transition over the first $j - 1$ samples of the sampling interval and having a transition at the j th sample. Hence, the considered reconstruction method makes errors on samples from time $t_i + j$ to $t_i + T_0 - 1$, corresponding to a Hamming distance of $T_0 - j$. A similar argument can be made for the case of $S_i = 1$.

By considering these sampling state cost functions, the BFE can be written as follows:

$$\begin{aligned} J(S_i = 0) &= \min_{T_0, T_1} [c(S_i = 0, T_0) \\ &\quad + \beta(1 - \epsilon_0)^{T_0} J(S_i = 0) \\ &\quad + \beta(1 - (1 - \epsilon_0)^{T_0}) J(S_i = 1)], \\ J(S_i = 1) &= \min_{T_0, T_1} [c(S_i = 1, T_1) \\ &\quad + \beta(1 - \epsilon_1)^{T_1} J(S_i = 1) \\ &\quad + \beta(1 - (1 - \epsilon_1)^{T_1}) J(S_i = 0)]. \end{aligned}$$

Since these BFEs have only two variables, various numerical and analytical methods can be applied to find their solution (e.g., see [351]), which in turn corresponds to an optimal sampling scheme in TANS. Simulation results for this sampling scheme are given in Section 11.5.

11.4.2 Approximate Dynamic Programming Methods

Finding a solution of the BFE of a DP-based sampling function may not be practically feasible if the sampling state space is large or is unknown (or is known partially). In these cases, approximate dynamic programming (ADP) algorithms can provide suboptimal solutions with a reasonable computation complexity. In this section, we investigate an ADP-based sampling function of TANS introduced in equation (11.26). To use this approximate algorithm, for any sampling state S_i , a *quality* measure $q(S_i)$ is assigned. We use a greedy sampling solution to define this quality function.

Consider a Markovian signal $X(t)$ described in Section 11.2.2. A greedy sampling function for this signal is introduced in Algorithm 11.6. We refer to this greedy sampling function as $f^{\text{greedy}}(\cdot)$, where $T_i^{\text{greedy}} = f^{\text{greedy}}(S_i)$.

A *quality* function of each state is defined as

$$q(S_i) = \gamma T_i^{\text{greedy}} \tag{11.28}$$

where γ is a scaling parameter. Intuitively, the larger the greedy sampling step, the higher the quality of the sampling state. Therefore, an ADP-based sampling function of TANS can be derived by using the optimization setup of equation (11.26).

For computing the expected quality of the next sampling state, transition probabilities among different sampling states should be known. Suppose, at a sampling state S_i , the next sample is taken after a time interval T_i . The value of the sample $X(t_i + T_i)$ is a random variable with a mean $\hat{X}(t_i + T_i)$, which can be computed by a generalized linear prediction filter. The probability distribution of $X(t_i + T_i)$ determines the probability distribution of the next sampling state S_{i+1} . Therefore, the expected value of the quality function of the next state can be computed by using

this probability distribution. However, to simplify this sampling function further, one may approximate this expected quality by the quality of the most probable next state, which has a sample value of $\hat{X}(t_i + T_i)$ at time t_{i+1} . We call this state \hat{S}_{i+1} . Therefore, a more simplified sampling function based on an ADP algorithm can be written as

$$T_i = \arg \min_T c(S_i, T) + \beta q(\hat{S}_{i+1}).$$

For any given state S_i and $T_i = T$, $q(\hat{S}_{i+1})$ can be computed. Intuitively, the term $\beta q(\hat{S}_{i+1})$ is a correction term for the greedy solution considering the quality of the next sampling state. In this scheme, the quality of only one future sampling state is considered. However, one can extend this algorithm to consider the quality of more than one future sampling state.

Algorithm 11.12. An approximate dynamic programming-based sampling function for the Markovian signal of equation (11.4) can be summarized as follows:

- Step i,0: Compute $\hat{\theta}_{S_i}$, $P_e(S_i)$ and $q(\hat{S}_{i+1})$.
- Step i,1: Compute $T_i = \arg \min_T (c(S_i, T | \hat{\theta}_{S_i}) + \beta q(\hat{S}_{i+1}))$.
- Step i,2: Take a sample at time $t_i + T_i$.
- Step i,3: Compute S_{i+1} . Repeat.

Simulation results for this sampling procedure are presented in Section 11.5.

11.5 Simulation Results

In this section, we evaluate the performance of the proposed sampling schemes by simulations and compare their performance against uniform sampling. In uniform sampling, the sampling rate is always in the form of $R = 1/T_{\text{uni}}$, where T_{uni} is a positive integer. To be able to compare the performance of different methods with uniform sampling at different rates, we modify the uniform sampling setup to capture

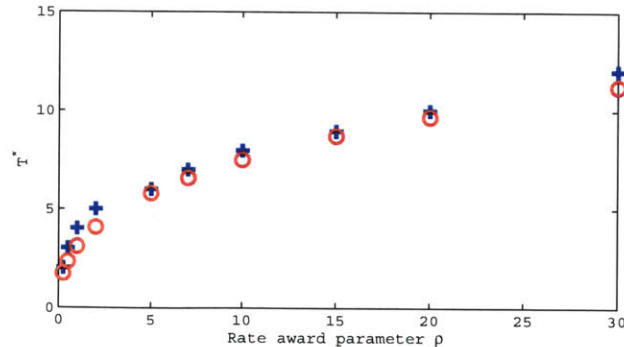


Figure 11-4: Illustration of Corollary 11.5 by simulation. The red curve is the solution of (11.11). The blue curve is optimal sampling rate T^* . Note that, their difference is always less than 1.

all possible sampling rates. To do this, for a given rate $R = 1/T_{\text{uni}}$ where T_{uni} is not an integer number, the i th sample is taken at time $t_i = \text{round}(T_{\text{uni}})$.

First, we consider an autoregressive signal model introduced in equation (11.3). In Theorem 11.4, we show that, an optimal sampling scheme for this signal is uniform. Also, Corollary 11.5 provides a straightforward way to compute the optimal sampling rate. Figure 11-4 illustrates this corollary for the case of $\alpha = 0.99$ (i.e., noise power is 0.02) and for different ρ values.

Now, we consider Markovian signals introduced in equation (11.4). Transition probabilities of the underlying MC are assumed to be 0.001. We demonstrate the performance of different sampling methods on a rate-distortion plots (Figures 11-5, 11-6 and 11-7). Rate refers to average sampling rate and distortion refers to average reconstruction distortion per sample. We use MSE as our distortion measure. For greedy TANS, we use generalized linear prediction (GLP) filter as the reconstruction method. Note that, this reconstruction is causal. For uniform sampling, we use three reconstruction methods: causal line-connecting (CLC), non-causal line-connecting (NCLC) and GLP filtering. Lower curves in these plots indicate better performance.

Figure 11-5 shows analytical rate-distortion curves of Theorem 11.8 for a Markovian signal of equation (11.4) for different estimation error probabilities. Here, the noise power in state 0 of the Markov chain is 0.05 (i.e., $\alpha_0 \approx 0.97$) and the noise power in state 1 is 0.5 (i.e., $\alpha_1 \approx 0.7$). As illustrated in this plot, the lower the error

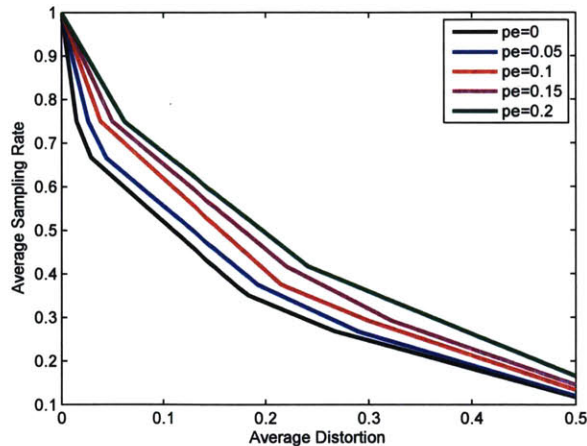


Figure 11-5: Rate-distortion curves of Theorem 11.8 for a Markovian signal of equation (11.4) with the signal parameters $\alpha_0 = 0.01$ and $\alpha_1 = 0.99$. P_e represents error probability of the estimation process. $P_e = 0$ curve corresponds to the so-called genie-aided scheme.

probability, the better the performance. Note that the case where $P_e = 0$ is referred to as a genie-aided scheme.

Figures 11-6 and 11-7 show rate-distortion curves achieved by simulations for various schemes. The signal model is Markovian as in equation (11.4) with parameters $\alpha_0 = 0.01$ and $\alpha_1 = 0.99$ in schemes considered in Figure 11-6 and, $\alpha_0 = 0.7$ and $\alpha_1 = 0.97$ in the ones of Figure 11-7. In greedy schemes, we use a maximum likelihood estimation block with $m = 10$ to estimate the state of the underlying Markov chain. Although, it is not fair to compare performance of greedy TANS with causal reconstruction with a uniform sampling scheme with non-causal reconstruction, in the case of Figure 11-6, greedy TANS outperforms uniform sampling schemes including the one with a non-causal reconstruction method. In the case of Figure 11-7, greedy TANS outperforms uniform schemes with causal reconstructions. In the low-distortion regime (distortion less than 0.08), it also outperforms uniform sampling with non-causal reconstruction.

As illustrated in these figures, genie-aided greedy TANS provides an analytical lower bound for greedy TANS in the rate-distortion plot. The proposed greedy TANS performs closely to this lower bound. Also, by choosing the estimation error probability 0.05 (estimated from simulations), an analytical upper bound for greedy TANS can

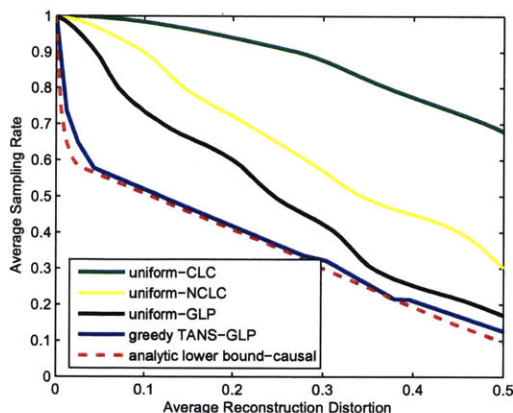


Figure 11-6: Average sampling rate versus average reconstruction distortion for a Markovian signal with parameters $\alpha_0 = 0.01$ and $\alpha_1 = 0.99$ for methods: (i) uniform sampling with causal line-connecting (CLC) reconstruction, (ii) uniform sampling with non-causal line-connecting (NCLC) reconstruction, (iii) uniform sampling with generalized linear prediction (GLP) filtering (iv) greedy TANS with generalized linear prediction (GLP) filtering, and (v) analytical lower bound for greedy TANS based on Theorem 11.8.

be achieved as proposed in Theorem 11.8. Moreover, by comparing two greedy TANS schemes of these figures, we notice that, the gain provided by the TANS framework is larger in the case that, noise variances differ more in two MC states.

Figure 11-8 shows the performance of a dynamic programming-based TANS scheme for an online source coding application explained in Section 11.4.1. Here, we assume $\epsilon_0 = 0.1$ and $\epsilon_1 = 0.01$. To solve the Bellman-Ford equation (11.28), a value-iteration method is used [351]. As shown in this plot, a DP-based TANS scheme outperforms uniform sampling.

Figure 11-9 illustrates performance of a TANS scheme based on approximate dynamic programming for a Markovian signal explained in Algorithm 11.12. Here, we assume that underlying Markov chain transition probabilities are 0.1 (i.e., $p_{01} = p_{10} = 0.1$). The signal parameters are assumed to be $\alpha_0 = 0.7$ and $\alpha_1 = 0.99$. As illustrated in this figure, a TANS scheme based on ADP outperforms the greedy one.

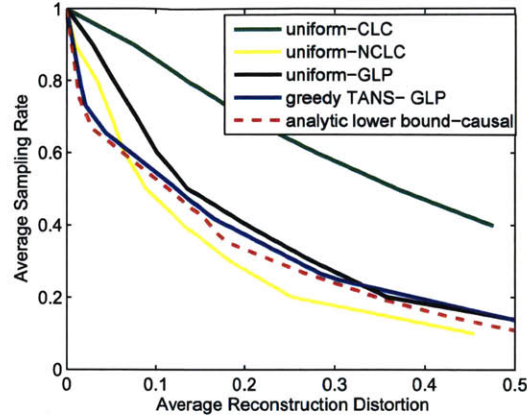


Figure 11-7: Average sampling rate versus average reconstruction distortion for a Markovian signal with parameters $\alpha_0 = 0.7$ and $\alpha_1 = 0.97$ for methods: (i) uniform sampling with causal line-connecting (CLC) reconstruction, (ii) uniform sampling with non-causal line-connecting (NCLC) reconstruction, (iii) uniform sampling with generalized linear prediction (GLP) filtering (iv) greedy TANS with generalized linear prediction (GLP) filtering, and (v) analytical lower bound for greedy TANS based on Theorem 11.8.

11.6 Proofs

11.6.1 Proof of Theorem 11.2

To find a solution of the optimization problem equation (11.5), we use similar techniques as for the regular linear prediction filter [350]. Note that, in an optimal scheme, the error term should be orthogonal to all known samples:

$$E[X(t_{i-k})e^*(t)] = 0 \quad (11.29)$$

for $k = 0, \dots, m-1$, where $e^*(t) = X(t) - \sum_{k=0}^{m-1} w_{\tau_k}^* X(t_{i-k})$.

By using equations (11.6) and (11.29), optimal weights $w_{\tau_i}^*$, for $0 \leq i \leq m-1$ should satisfy the following set of linear equations:

$$r(-\tau_k) = \sum_{i=0}^{m-1} w_{\tau_i}^* r(\tau_i - \tau_k) \quad (11.30)$$

for $k = 0, \dots, m-1$. By using matrix notations of equation (11.7), linear equations

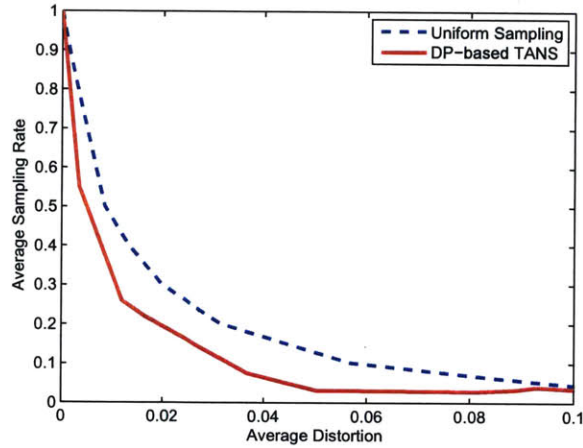


Figure 11-8: A rate-distortion plot of a dynamic programming-based TANS for an online source coding application explained in Section 11.4.1.

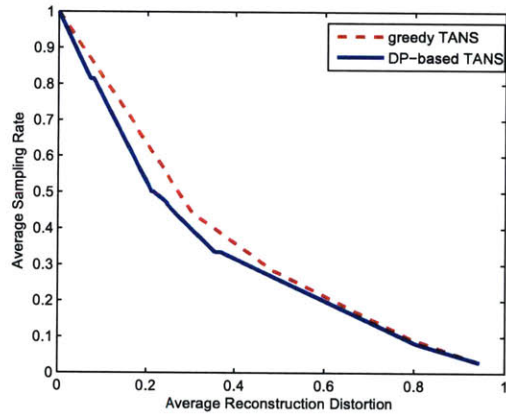


Figure 11-9: Comparison of a dynamic programming-based TANS with greedy TANS for a Markovian signal model.

of (11.30) can be written as follows:

$$\mathbf{p} = \mathbf{R}\mathbf{w}_\tau^*.$$

11.6.2 Proof of Theorem 11.3

Since $X(t)$ has zero mean and $X(t) = \hat{X}(t) + e^*(t)$, by using (11.29), we have

$$\sigma_{e^*}^2(S_i, T) = \sigma_X^2 - \sigma_{\hat{X}}^2 = r(0) - \sigma_{\hat{X}}^2, \quad (11.31)$$

where $\hat{X}(t) = (\mathbf{w}_\tau^*)^T(\mathbf{X}_{S_i})$. Therefore,

$$\begin{aligned}
 \sigma_{\hat{X}}^2 &= E[|\hat{X}(t)|^2] \\
 &= (\mathbf{w}_\tau^*)^T E[(\mathbf{X}_{S_i})(\mathbf{X}_{S_i})^T] \mathbf{w}_\tau^* \\
 &= (\mathbf{w}_\tau^*)^T \mathbf{R} \mathbf{w}_\tau^* \\
 &= \mathbf{p}^T \mathbf{w}_\tau^*.
 \end{aligned} \tag{11.32}$$

Equations (11.31) and (11.32) establish the theorem.

11.6.3 Proof of Theorem 11.4

By using the definition of an AR(1) signal, the sampling state cost can be written as follows:

$$c(S_i, T) = \sum_{j=1}^{T-1} (1 - \alpha^{2j}) + \frac{\rho}{T}. \tag{11.33}$$

Hence, an optimal greedy sampling solution for this sampling state cost can be computed as follows:

$$T_{\text{greedy}}^* = \arg \min_T \sum_{j=1}^{T-1} (1 - \alpha^{2j}) + \frac{\rho}{T}. \tag{11.34}$$

Now, we show that this optimal greedy TANS (which is a uniform sampling scheme) performs closely to an optimal sampling scheme, which may be nonuniform. A sampling state S_i is a uniform sampling state if signal samples within that state are distributed uniformly in time. Consider a uniform sampling state S_i with an inter-state sampling increment T_{greedy}^* . The sampling state cost of this uniform sampling scheme is referred to by c_{\min} (this sampling state cost can be achieved by having T_{greedy}^* as a solution of the optimization setup (11.8)). First, we show that this uniform sampling scheme satisfies the BFE (11.25). Since under this sampling scheme, $S_{i+1} = S_i$, therefore, the cost-to-go function at the sampling state S_i can be

written as follows:

$$\begin{aligned}
 J_f(S_i) &= c(S_i, T_i) + \beta c(S_{i+1}, T_{i+1}) + \dots \\
 &= c(S_i, T_{\text{greedy}}^*) (1 + \beta + \beta^2 + \dots) \\
 &= \frac{c_{\min}}{1 - \beta}.
 \end{aligned} \tag{11.35}$$

Now, consider the right hand side (RHS) of the BFE (11.25):

$$\begin{aligned}
 \text{RHS of BFE} &= \min_{T_i} (c(S_i, T_i) + \beta E[J_f(S_{i+1})]) \\
 &= c_{\min} + \beta \frac{c_{\min}}{1 - \beta} \\
 &= \frac{c_{\min}}{1 - \beta} \\
 &= J_f(S_i).
 \end{aligned} \tag{11.36}$$

Therefore, the above uniform sampling scheme satisfies the BFE when S_i (i.e., the initialization state) happens to be a uniform sampling state with an inter-state sampling step size T_{greedy}^* . If the process starts from another sampling state, we assume that there is always a way to reach to this uniform sampling state (e.g., take m uniform samples with a sampling increment T_{greedy}^*). Moreover, the mapping functions between policies and costs are continuous. Hence, a small difference in costs due to initialization effects has a small effect in sampling policies. Therefore, for an AR(1) signal model, the optimal greedy TANS (which is uniform) performs closely to an optimal sampling scheme (which may be nonuniform).

11.6.4 Proof of Corollary 11.5

Define $h(T) = (1 - \alpha^{2T}) - \rho/(T(T+1))$ for $T \in [1, \infty)$. Note that $h(T)$ is a continuous function over this interval. Since $h(1) = (1 - \alpha^2) - \rho/2 < 0$ and $\lim_{T \rightarrow \infty} h(T) = 1 > 0$, by using the mean value theorem, h has a root T_{root} in $[1, \infty)$. Since $dh(T)/dT > 0$, this root is unique, completing the proof of part (i).

For an AR(1) signal, we have:

$$h(T) = c(S_i, T + 1) - c(S_i, T) = (1 - \alpha^{2T}) - \frac{\rho}{T(T + 1)}.$$

Since $T^* = \arg \min_T c(S_i, T)$ and also T^* is an integer, either $T^* = \lfloor T_{\text{root}} \rfloor$ or $T^* = \lfloor T_{\text{root}} \rfloor + 1$, completing the proof of part (ii).

11.6.5 Proof of Theorem 11.8

Under conditions of equation (11.12), equation (11.14) shows the sampling state cost given $\hat{\theta}_{S_i} = 0$. By using the sampling scheme proposed in Algorithm 11.6, we have:

$$T_i = \arg \min_T c(S_i, T | \hat{\theta}_{S_i} = 0). \quad (11.37)$$

Since $\sigma_{\max}^2 \geq 1$, by having $P_e^{\text{low}} \leq P(S_i) \leq P_e^{\text{up}}$, given $\hat{\theta}_{S_i} = 0$, we have: $T_0^{\text{low}} \leq T_i \leq T_0^{\text{up}}$. Similarly, one can show that, given $\hat{\theta}_{S_i} = 1$, $T_1^{\text{low}} \leq T_i \leq T_1^{\text{up}}$. Since the underlying Markov chain is symmetric (i.e., $p_{01} = p_{10}$), in steady state, $\hat{\theta}_{S_i} = 0$ approximately half of the time and $\hat{\theta}_{S_i} = 1$ approximately half of the time. Hence, the number of samples taken in that state by using the proposed sampling scheme of Algorithm 11.6 is bounded between

$$\frac{T_{\text{tot}}}{2T_0^{\text{low}}} + \frac{T_{\text{tot}}}{2T_1^{\text{low}}} \quad \text{and} \quad \frac{T_{\text{tot}}}{2T_0^{\text{up}}} + \frac{T_{\text{tot}}}{2T_1^{\text{up}}}.$$

This demonstrates the sampling rate bounds. By using bounds on T_i and $P_e(S_i)$ in equation (11.14), deriving bounds on the expected reconstruction distortion is straightforward.

Chapter 12

Backward Adaptation for Power Efficient Sampling

For bandlimited signals, the Nyquist sampling theorem provides an appropriate uniform sampling setup which leads to a signal reconstruction with zero error [352]. Although the Nyquist sampling theorem defines a sufficient condition on the sampling rate for bandlimited signals, it does not necessarily provide an efficient representation because a signal may seem to have a lower bandwidth *locally* and thus sampling at the Nyquist rate seems to be a waste [338]. These extra samples not only increase the energy consumption of the sampling process, they also create some processing and transmission overheads. To avoid these energy and rate inefficiencies, sample redundancies can be removed using a compression method at the expense of increased computational costs.

An alternative to uniform sampling would be to take samples only when they are *innovative* for the relevant application. This leads to an adaptive nonuniform sampling scheme. For instance, a nonuniform sampling scheme based on level-crossings with iterative decoding is considered in [336], while reference [337] presents an approach based on level crossings with a filtering technique, which adapts the sampling rate and filter order by analyzing the input signal variations online. Two adaptive sampling schemes for bandlimited deterministic signals are also proposed in [338]. These schemes use definitions of local bandwidth, based on linear time-varying low pass filters [339] and time-warping of bandlimited signals [340]. Two main issues about these nonuniform sampling schemes make them difficult to be applied in practical applications: first, they are designed for specific signal models (i.e., they are not generic), and second, sampling times are required to be kept or transmitted in addition to sample values, in order to be used in the reconstruction process.

A relatively recent technique, called compressed sensing (CS) [307], has been used as a signal-agnostic acquisition method for highly sparse signals. In CS-based schemes, sampling is performed approximately at the information rate, rather than the Nyquist rate. Since CS has low encoding complexity, it has been adopted in several low power systems with the goal of reducing the power consumption [290, 353].

Here, we use a framework proposed in [354, 355] to address main limitations of previous sampling approaches. This new adaptive nonuniform sampling framework

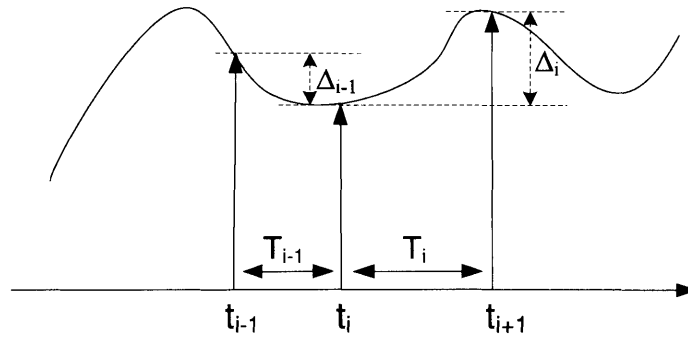


Figure 12-1: TANS setup

has two key characteristics: first, it continuously adapts the sampling rate according to the *local bandwidth* of the signal and consequently reduces the number of captured samples; second, unlike traditional nonuniform sampling procedures, sampling times do not need to be transmitted, since the receiver can recover them from the sample values themselves, thus avoiding the overhead of transmitting the sampling times which is required for traditional nonuniform sampling schemes. This sampling framework is called time-stampless adaptive nonuniform sampling (TANS).

Three TANS schemes are presented in this work for different signal models and sampling requirements. First, we use Taylor's theorem to derive a suitable sampling function for deterministic signals. We term this approach TANS by polynomial extrapolation (TPE). Next, for the same signal family, we present a TANS scheme by incremental variation (TIV), where the sampling time intervals are chosen from a lattice. Finally, we propose a TANS scheme which is constrained to a finite set of sampling rates (TFR) and illustrate its energy efficiency in the context of transmitting an electrocardiography (ECG) signal model from the MIT-BIH database [356]. By using detailed Verilog modeling of an ECG processor and transmitter, we show that, TANS provides significant improvements in terms of energy consumption compared to uniform sampling followed by a transformation-based compression, a nonuniform sampling scheme and compressed sensing.

12.1 Time-Stampless Adaptive Nonuniform Sampling

In this section, we review the TANS framework, proposed in [354, 355]. Consider a continuous-time signal $X(t)$. Suppose the i -th sample is taken at time t_i . Define $T_i \triangleq t_{i+1} - t_i$ and $\Delta_i \triangleq X(t_{i+1}) - X(t_i)$, as shown in Fig. 12-1. We take the $(i + 1)$ st sample after a time increment of

$$T_i \triangleq f(t_{i-m+1}, \dots, t_{i-1}, t_i, X(t_{i-m+1}), \dots, X(t_{i-1}), X(t_i)),$$

where f is called the sampling function, that allows us to reconstruct the time samples from the sample values themselves. Since the next sampling time step is a function of m most recently taken samples, we say the order of the sampling function $f(\cdot)$ is m . The sampling function has to be known at both the sampling and reconstruction sides.

This sampling structure is nonuniform except in trivial cases, when the sampling function is a constant-valued function. The key characteristic of our approach is that, unlike traditional nonuniform sampling procedures, keeping sampling times (time-stamps) is not necessary in our framework because these times can be recovered by using the sampling function and previously taken samples. In the above example, we have $t_{i+1} = t_i + f(\bigcup_{j=i-m+1}^{i-1} \{T_j, \Delta_j\})$. Note that the first m sampling times should be kept to initialize the process. However, the effect of those initial samples is negligible, on average, when the number of samples increases.

In the above case, TANS is causal because the next sampling time depends on samples taken before that time. In general, it can be designed to be non-causal. TANS is an adaptive process which learns from the signal through the taken samples, because the sampling function depends on local characteristics of the signal. Finding an appropriate sampling function depends on different applications with various requirements such as the sampling rate, the distortion requirement, the computational complexity, etc. Here, to maintain focus on sampling rate and adaptation of sampling increments, we do not explicitly include quantization effects.

12.2 TANS for Signals with Bounded Derivatives

In this section, we propose two TANS schemes for signals with bounded third derivative, by characterizing local variations of the signal by using taken samples. First, we propose a TANS scheme by polynomial extrapolation (TPE), which only assumes the third derivative of the signal is bounded, but requires no other specific knowledge of the signal. Then, we propose a TANS method by incremental variation (TIV), where the sampling time intervals are chosen from a lattice. Finally, we compare the performance of these two schemes with uniform sampling in terms of rate-distortion.

Suppose $X(t)$ is a continuous signal. We consider the following assumptions:

- (A₁) The order of the sampling function is two. In other words, the sampling function is a function of three previously taken samples ($T_i = f(\cup_{j=i-2}^{i-1} \{T_j, \Delta_j\})$).
- (A₂) The reconstruction method is connecting two adjacent samples by a line. Note that, until the receiver receives the sample $X(t_i)$, it only has an estimate of the signal on $(-\infty, t_{i-1}]$.

12.2.1 TANS by Polynomial Extrapolation

Various instantiations of TANS can be derived by using Taylor's theorem, depending on the signals, the desired complexity of the sampling function f , and the distortion requirement. In this section, as an example, we explain how a sampling function can be computed for a specific setup. Taylor's theorem is used to extrapolate the signal by a polynomial. However, since we don't actually have the first and second derivatives of the signal, we use the assumption of having bounded third derivative to bound the error incurred by using an approximation of the first and second derivatives.

Theorem 12.1. *In TPE, if $|X'''(t)|$ is uniformly bounded by a constant M , under assumptions A_1 - A_2 , the following sampling function satisfies the sampling requirement $|X(t) - \hat{X}(t)| < D$ for all t :*

$$T_i = \arg \max T$$

$$\text{s.t.} \quad (c_1 T + c_2) T^2 \leq D \quad (12.1)$$

where $\hat{X}(t)$ is the reconstructed signal, and c_1 and c_2 are constants defined as follows,

$$c_1 = \frac{2}{3} M$$

$$c_2 = \underbrace{\frac{|\frac{\Delta_{i-1}}{T_{i-1}} - \frac{\Delta_{i-2}}{T_{i-2}}|}{\frac{T_{i-1} + T_{i-2}}{2}}}_I + \frac{1}{3} M \frac{(T_{i-1} + T_{i-2})^2 + (T_{i-1})^2}{T_{i-2}}$$

We present the proof of Theorem 12.1 in Section 12.6.1. Note that the function $(c_1 T + c_2) T^2$ is increasing for positive T , and one can use intermediate value theorem to compute the solution of the optimization (12.1). In practice, T_i can be obtained by interpolating from a stored table. It is insightful to investigate the behavior of T_i , with respect to the parameters in (12.1).

- Increasing D leads to an increase in T_i . Intuitively, the higher the allowed distortion, the lower the sampling rate, and the larger the sampling time intervals.
- The first term of c_2 (part I) can be viewed as an approximation of $|X''(t)|$ at time t_i . Since the reconstruction method is a first-order linear filter, the higher the second derivative, the faster changes of the signal. Therefore, by increasing this term, T_i should decrease.

Define the reconstruction error signal $e(t) = X(t) - \hat{X}(t)$. If over the time interval $[t_i, t_{i+1})$, the normalized \mathcal{L}_1 norm of the error signal is considered as a distortion measure (i.e., $\frac{1}{T_i} \int_{t=t_i}^{t_{i+1}} |e(t)| dt \leq D$), the following sampling function can be derived:

Corollary 12.2. In TPE, if $|X'''(t)|$ is bounded by a constant M , under assumptions A_1 - A_2 , the following sampling function satisfies the normalized \mathcal{L}_1 norm of the error

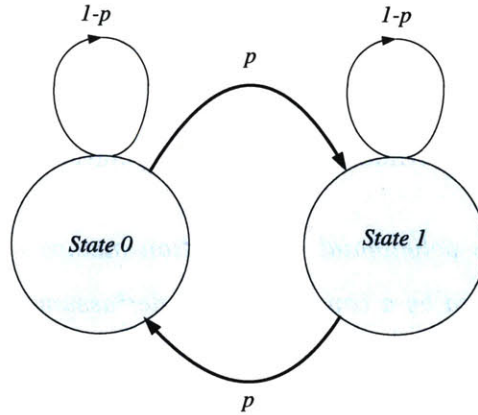


Figure 12-2: An underlying hidden Markov chain used for the signal creation.

signal over $[t_i, t_{i+1})$ to be less than D :

$$\begin{aligned}
 T_i &= \arg \max T \\
 \text{s.t.} \quad & (c_1 T + c_2) T^2 \leq D
 \end{aligned} \tag{12.2}$$

where c_1 and c_2 are constants, defined as follows,

$$\begin{aligned}
 c_1 &= \frac{M}{4}, \\
 c_2 &= \frac{5}{12} \left(\frac{|\frac{\Delta_{i-1}}{T_{i-1}} - \frac{\Delta_{i-2}}{T_{i-2}}|}{\frac{T_{i-1} + T_{i-2}}{2}} + \frac{1}{3} M \frac{(T_{i-1} + T_{i-2})^2 + (T_{i-1})^2}{T_{i-2}} \right).
 \end{aligned}$$

Proof of this corollary is presented in Section 12.6.2.

A counterpart of this method can be derived for discrete signals. For example, suppose $X[n]$ is a discrete signal, where n is a non-negative integer. The i^{th} sample is taken at time n_i . Let us define $N_i \triangleq n_{i+1} - n_i$ and $\Delta_i \triangleq X[n_{i+1}] - X[n_i]$. Then, the $(i + 1)^{\text{th}}$ sample is taken after $N_i = f(\bigcup_{j=i-2}^{i-1} \{N_j, \Delta_j\})$ samples. Let us define $X'[n] \triangleq X[n] - X[n - 1]$, $X''[n] \triangleq X[n] - 2X[n - 1] + X[n - 2]$ and $X'''[n] \triangleq X[n] - 3X[n - 1] + 3X[n - 2] - X[n - 3]$. Suppose $|X'''[n]|$ is uniformly bounded by a constant M . Also, consider the following assumptions as counterparts for assumptions A_1 and A_2 :

(A₄) The sampling function is a function of three previously taken samples (i.e., $N_i = f(\bigcup_{j=i-2}^{i-1} \{N_j, \Delta_j\})$).

(A₅) The reconstruction method is a linear interpolation among taken samples.

Theorem 12.3. *In the polynomial extrapolation method of TANS for discrete signals, if $|X'''[n]|$ is bounded by a constant M , under assumptions A₄-A₅, the following sampling function satisfies the sampling requirement $|X[n] - \hat{X}[n]| < D$:*

$$\begin{aligned} N_i &= \arg \max [N] \\ \text{s.t.} \quad &(c_1 N + c_2) N^2 \leq D \end{aligned} \quad (12.3)$$

where $\hat{X}[n]$ is the reconstructed signal, and c_1 and c_2 are constants, defined as follows,

$$\begin{aligned} c_1 &= \frac{2}{3} M \\ c_2 &= \frac{|\frac{\Delta_{i-1}}{N_{i-1}} - \frac{\Delta_{i-2}}{N_{i-2}}|}{\frac{N_{i-1} + N_{i-2}}{2}} + \frac{1}{3} M \frac{(N_{i-1} + N_{i-2})^2 + (N_{i-1})^2}{N_{i-2}}. \end{aligned}$$

Proof of this theorem is similar to Theorem 12.1.

While our proof has considered deterministic signals, since the proposed sampling method is blind, it renders this method a well-suited choice for uncertain signals. Hence, we empirically evaluate the performance of our scheme when the signal is deterministic over stochastically varying time intervals. To do this, consider an underlying hidden Markov chain (MC) with a transition probability p , depicted in Fig. 12-2. At each state, the signal is drawn by a third order spline, up-sampled by a factor of u_s , where s represents the state of MC. Thus, if the state of MC is s at time n , $X[n]$ and $X[n + u_s + 1]$ are drawn from a uniform distribution over $[-A, A]$. Signal values at other times are determined by a third order spline interpolation. Fig. 12-3(a) shows a sample of this signal family, where $p = 0.001$, $u_0 = 1$, $u_1 = 20$ and $A = 0.5$.

Note that spline interpolation is performed over the entire signal, rather than indi-

vidual states, so that the signal has continuous derivatives everywhere. The sampling function used is given by (12.2). Since our sampling scheme is blind, in our simulations, the signal model is unknown and we are only allowed to use the sampling function of (12.2). To do this, we need M , an upper bound for the third derivative of the signal. From taken samples, we approximate the third derivative as follows. First, we define,

$$w(t_i) \triangleq \frac{|\frac{\Delta_{i-1}}{T_{i-1}} - \frac{\Delta_{i-2}}{T_{i-2}}|}{\frac{T_{i-1} + T_{i-2}}{2}}. \quad (12.4)$$

In addition, $\zeta(t_i) \triangleq \frac{|w(t_i) - w(t_{i-1})|}{(T_{i-1} + T_{i-2} + T_{i-3})/3}$ is used as an approximation of the third derivative of the signal by Taylor's theorem. At each time, M is chosen as the maximum of $\{\zeta(t_i), \zeta(t_{i-1}), \dots, \zeta(t_{i-W})\}$. In our simulations, we pick $W = 10$.

Fig. 12-3(b) shows the number of samples (i) versus time (n_i) when $D = 0.2$. At each time n_i , the slope of the tangent line on the curve represents the sampling rate at that time. We draw state transitions of the underlying hidden Markov chain in Fig. 12-3(a) to clarify the rate adaptation in different states. In state 0, in which the signal model is an up-sampled third-order spline by a factor $u_0 = 1$, the sampling rate is greater than the one of state 1, where $u_1 = 20$. The overall sampling rate is 0.6 with the normalized \mathcal{L}_1 distortion 0.001 for this signal sample. Note that, although the sampling rate is fairly stable within each state, it has some fluctuations.

12.2.2 TANS by Incremental Variation

In this section, we introduce a family of TANS where sampling time intervals are discrete values. In this method, the sampling rate is being adapted based on an estimation of local variations of the signal. This method can be applied to both continuous and discrete signals. In this section, we consider its continuous time version; however, its discrete time version can be derived similarly to Theorem 12.3.

If $|X'''(t)|$ is uniformly bounded by a constant M , under assumptions C_1 and C_2 mentioned in Section 12.2.1, an estimate of local variation of the signal at time t_i , by using previously taken samples, is given by Eq. (12.4).

By a similar argument as the (12.30), the error of this estimation can be bounded

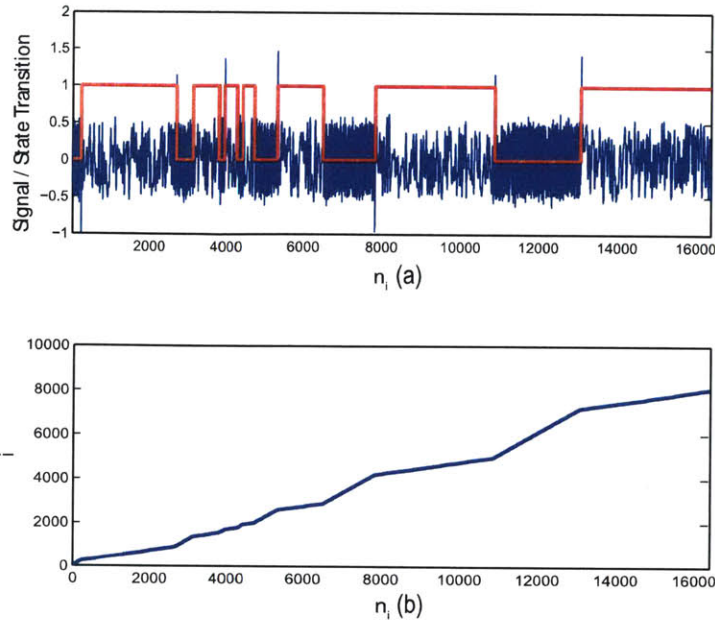


Figure 12-3: (a) A sample signal and its MC state transitions used in Section (12.2.1). (b) Sampling rate versus time for this signal with TPE method.

as follows,

$$|w(t_i) - |X''(t_i)|| \leq \frac{2}{3} M \frac{(T_{i-1} + T_{i-2})^2 + (T_{i-1})^2}{T_{i-2}}. \quad (12.5)$$

Consider the following heuristic sampling function,

$$T_i = \begin{cases} f_1(T_{i-1}) & w(t_i) < th_1 \text{ and } T_i > T^{min} \\ T_{i-1} & th_1 \leq w(t_i) < th_2 \\ f_2(T_{i-1}) & w(t_i) \geq th_2 \text{ and } T_i < T^{max} \end{cases} \quad (12.6)$$

where $f_1(T_{i-1}) > T_{i-1}$ and $f_2(T_{i-1}) < T_{i-1}$. T^{min} and T^{max} are *a priori* selected time bounds. Thresholds th_1 and th_2 depend on signal characteristics and sampling requirements. If $w(t_i)$ is smaller than a threshold, the signal's slope variations are small and we can decrease the sampling rate in this case, since $f_1(T_{i-1}) > T_{i-1}$. An analogous argument can be expressed when $w(t_i)$ is greater than a threshold.

An example of a sampling function (12.6) with linear increase or decrease of T_i can be expressed as follows,

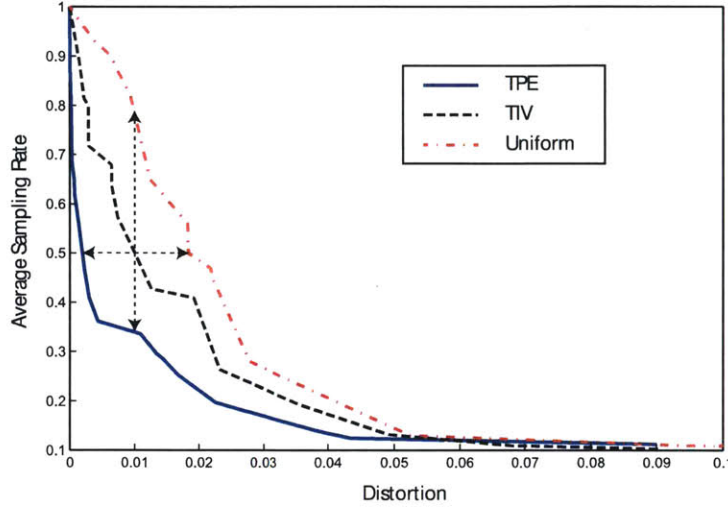


Figure 12-4: A rate-distortion comparison among TPE, TIV and uniform sampling.

$$T_i = \begin{cases} T_{i-1} + \epsilon_1 & w(t_i) < th_1 \text{ and } T_i > T^{min} \\ T_{i-1} & th_1 \leq w(t_i) < th_2 \\ T_{i-1} - \epsilon_2 & w(t_i) \geq th_2 \text{ and } T_i < T^{max} \end{cases} \quad (12.7)$$

where ϵ_1 and ϵ_2 are positive constants. Note that, given T_{i-1} , there are only three possibilities for T_i , so the sampling time intervals are discrete.

12.2.3 Comparison with Synthetic Signals

In this subsection, the performance of proposed TANS schemes is compared against uniform sampling using a synthetic signal. In a uniform sampling scheme with a given sampling rate $R = \frac{1}{N_s}$, where N_s is not necessarily an integer number, the i^{th} sample is taken at time $n_i = Ri$ from the underlying continuous spline signal. This enables us to capture all possible sampling rates for a full performance assessment.

We use the signal model described in Section 12.2.1. For the TPE method, we use a setup explained in Section 12.2.1. For the TIV method, we set $\epsilon_1 = \epsilon_2 = 1$, $T^{min} = 1$ and $T^{max} = 10$. The reconstruction method used for all methods is the first order linear filter. The distortion measure is the \mathcal{L}_1 norm of the error.

Fig. 12-4 shows the average sampling rate versus distortion for different methods. As illustrated in this figure, TPE method outperforms both TIV and uniform methods, where TIV method outperforms uniform as well. For instance, at distortion 0.01, both TPE and TIV methods outperform uniform sampling with a rate factor of approximately 2.8 and 1.7, respectively (e.g., TPE requires 2.8 times less samples to have the same distortion, compared to uniform sampling). For fixed average sampling rate of 0.5, TPE and TIV outperform uniform sampling by a distortion factor of approximately 8.7 and 1.8, respectively. Note that, non-smoothness of the curves is due to experimental artifacts.

12.3 TANS with Finite Sampling Rates

In this section, we introduce a TANS scheme where sampling rates cannot be arbitrary as in the ones of TPE, but should be chosen from a finite set $\{R_1, R_2, \dots, R_k\}$, created according to a specific signal model. In other words:

$$T_i = \frac{1}{R_{s_i}}, \quad (12.8)$$

where s_i is the sampling index which belongs to $\{1, 2, \dots, k\}$. This scheme is called TANS with finite sampling rates (TFR). Similarly to a general TANS framework, the sampling index at time t_i (s_i) should only depend on previously taken samples so that storing/transmitting sampling time-stamps would be unnecessary. Computing this sampling index function depends on the signal model and sampling characteristics. For instance, consider a signal $X(t)$ with k states which are repeated periodically with a known order. The expected length of state j is denoted by l_j . We use TFR with k rates $\{R_1, \dots, R_k\}$ to take measurements from this signal. Suppose N_i is the number of taken samples until time t_i from the current period of the signal. Then, the following is a TFR sampling scheme based on a mean-field approximation:

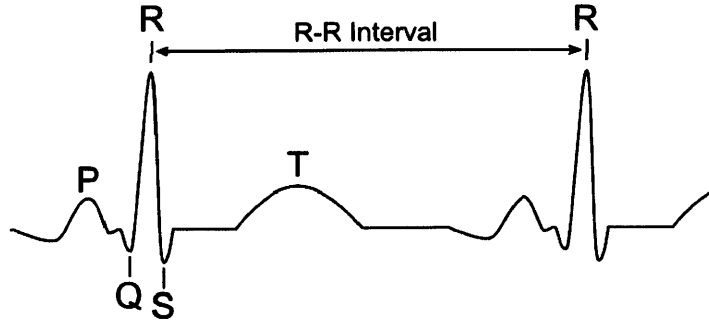


Figure 12-5: An ECG signal modeled as a P-QRS-T complex.

$$\begin{aligned}
 s_i &= \arg \min_{1 \leq r \leq k} r \\
 \text{s.t.} \quad & \sum_{s=1}^r (l_s R_s - N_i) \geq 0.
 \end{aligned} \tag{12.9}$$

This sampling function uses N_i to estimate the state of the signal at time t_i , and then takes samples with the corresponding rate of that state. There are some points to be noted about this sampling function:

- The TFR sampling function of Equation (12.9) assumes that the beginning of each signal period is known or can be detected using the taken samples. This assumption can be satisfied for some signal models. For example, ECG signals have a distinguished peak in each of their periods that can be used to detect the index of the signal period. We will discuss about it later in this section.
- Several parameters of this sampling function should be learned from the signal during a training period. Specifically, the number of states of the signal, their expected lengths and rates should be learned if they are not given.
- Suppose $D_i(R)$ is the expected distortion of the state i of the signal given the sampling rate is R . If $D_i(R) > D_j(R)$ for all R , then $R_i > R_j$. In other words, a higher sampling rate is used if the signal is more variable in a state.

One practical example of such a signal model is an ECG signal model depicted in

Fig. 12-5. We call the peak region around the point R, state 1, while state 2 covers the area close to the peak between points Q and S, and state 3 covers the rest of the ECG signal period. We learn TFR sampling parameters of Equation (12.9) using a training dataset once, where $R_1 = 180$, $R_2 = 360$, and $R_3 = 36$ samples/sec and $l_1 = 6$, $l_2 = 23$ and $l_3 = 33$, on average. Note that, state 2 has the highest sampling rate among others since detecting a transition from state 2 to 1 at the end of each signal period is crucial. Specifically, if $X(t_i) > th$ (in our example, $th = 50$), we say a transition occurs from state 2 to 1 and a new signal cycle has been started. State 3 is assigned to the lowest sampling rate since there is not much information in this region of the signal.

12.4 Performance Evaluation of TANS Schemes on Real Datasets

In this section, we examine the performance of the three TANS schemes presented earlier against state-of-the-art sampling schemes in a real application. Although the comparison is performed in the context of a health monitoring application, the proposed approach and results are generally applicable to a wide range of signals and applications. Fig. 12-6(a) shows an ECG signal from the MIT-BIH database¹ [356], which has been sampled with a frequency of 360 *Hz* and quantized to 10 bits. Using this signal, the performance of different sampling and compression schemes will be examined and compared based on their rate-distortion curves generated through simulations. The best scheme would have the lowest rate (*i.e.* highest compression ratio) for the same distortion level, or equivalently, the lowest distortion for the same rate.

Since the energy consumption of the sampling process is approximately proportional to its sampling frequency, it is desirable from an energy perspective to take as few samples as possible. However, the limit on the sampling rate imposed by the Nyquist theorem guarantees that distortion will be introduced if a signal is sampled

¹The DC component of the signal has been removed.

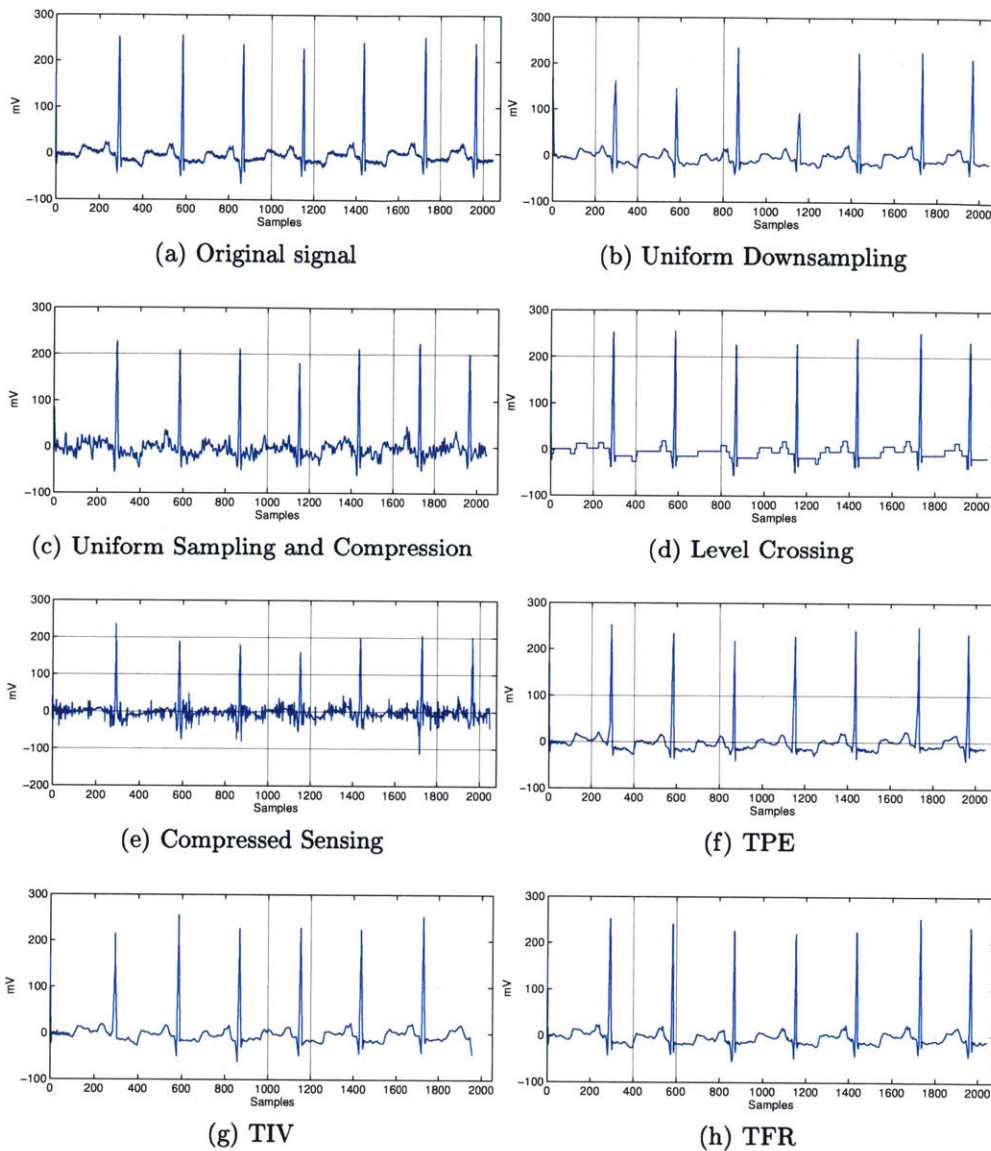


Figure 12-6: (a) A 6-sec ECG signal from MIT-BIH database, with the DC component removed. (b-h) Reconstructed ECG signals with an effective compression ratio of 6, using different sampling-compression schemes: (b) Uniform Downsampling, (c) Uniform Sampling followed by a DCT-based compression scheme. (d) Level-crossing Nonuniform Sampling. (e) Compressed sensing. (f) TANS by polynomial extrapolation (TPE). (g) TANS by incremental variations (TIV). (h) TANS with finite sampling rates (TFR).

below that limit. Fig. 12-6(b) shows the ECG signal after having been downsampled by a factor of 6 and reconstructed with linear interpolation. As it can be seen

in this figure, many of the features of the initial signal are missing. Thus, although this simplistic technique lowers the energy consumption of sampling and transmission processes by approximately a factor of 6, it does not provide satisfactory results in terms of the distortion.

Other techniques have been proposed in the literature which have better rate-distortion performance at the expense of higher complexity. In our analysis, we compare the performance of a few widely used schemes against the proposed TANS schemes. Fig. 12-6 (c)-(h) show the reconstructed signals of the original ECG signal after having been compressed by a factor of 6 using different schemes. In more details, the compared approaches are:

(i) **Uniform sampling and compression:** One widely used technique to minimize the amount of captured data sampled uniformly at Nyquist rate is to use compression algorithms, either lossless or lossy, where data is mapped to a more efficient representation before transmission. Majority of lossy compression algorithms proposed in the literature are transformation-based methods, such as wavelet- and DCT-based algorithms. These methods project signals to a certain domain with a more compact representation and, exploiting signal energy concentration, discard smallest coefficients through thresholding.

Depending on the application, different transformations may be preferable. In our comparisons, we use a DCT-based compression algorithm where DCT coefficients with magnitude smaller than a threshold are discarded. We choose this method mainly because it exhibits a better rate-distortion trade-off than using a wavelet transformation [357] and has practical implementation benefits as well [358]. Fig. 12-6(c) shows the reconstructed ECG signal after having been compressed with a DCT based algorithm with block size of 2048 samples, resulting in a compression ratio of 6. Adjusting the threshold of the encoding process balances between the required number of transmitted signal samples and the reconstruction distortion.

The main energy inefficiency of this approach is that, samples are taken at the acquisition stage at a relatively high frequency and then, some of these samples are discarded during the compression stage, resulting in extra energy consumption. For

this reason, several nonuniform schemes have been proposed in the literature, in which the sampling moments are triggered only when the signal crosses some predefined thresholds and consecutive samples close to the same value are not taken.

(ii) **Level crossing sampling:** A common nonuniform sampling technique is the level-crossing sampling, where the signal is compared with its previously taken value continuously and a new sample is taken when the difference exceeds a threshold [359]. However, apart from the samples themselves, time-stamps have to be transmitted since they are required in the reconstruction process. The reconstructed ECG signal using level-crossing sampling is shown in Fig. 12-6(d).

(iii) **Compressed Sensing (CS):** CS is an efficient signal-agnostic acquisition method for sparse signals [307]. The encoding process of a signal $X \in \mathbb{R}^N$, being k -sparse in some domain, is:

$$Y = \Phi X, \quad (12.10)$$

where $\Phi \in \mathbb{R}^{M \times N}$ and is called measurement matrix. If $M \sim O(k \log(N/k))$ measurements are taken from a signal X , CS exactly recovers the initial signal with high probability [360] if Φ satisfies the RIP condition:

$$(1 - \delta_k) \|s\|_2 \leq \|\Phi s\| \leq (1 + \delta_k) \|s\|_2, \quad (12.11)$$

for all k -sparse signals s , where δ_k is a constant $\in (0, 1)$. Signal reconstruction may be performed through an optimization problem which can be formulated as:

$$\min_{\hat{s} \in \mathbb{R}^n} \|\hat{s}\|_{l_1}, \text{ s.t. } Y = \Phi \Psi \hat{S}, \quad (12.12)$$

where $\Psi \in \mathbb{R}^{N \times N}$ is a matrix representing the synthesis operator of an orthonormal basis where the signal X has a sparse representation. The initial signal is recovered by multiplying the solution of the optimization problem with the inverse basis matrix:

$$\hat{X} = \Psi \hat{S}. \quad (12.13)$$

Alternatives to this basis pursuit (BP) formulation [361] include lasso [331], the Dantzig selector [362], greedy methods [363, 364] and iterative methods related to belief propagation [365]. Exact asymptotic performance under Bayesian formulations can be computed using the methods of [366]. In our simulations, a BP reconstruction algorithm is used, implemented with the CVX optimization package [367]. A Gaussian matrix is used as the measurement matrix with entries from $\mathcal{N}(0, 1/N)$, and a wavelet basis (*coiflet 4*) is used as the sparsifying domain. The result of the reconstruction process is shown in Fig. 12-6(e) for a compression ratio of 6.

(iv) **TPE, TIV and TFR:** These TANS schemes are adaptive sampling schemes and their sampling moments are calculated based only on previously taken samples through their sampling functions, without the need for tracking the current value of the signal continuously and for transmitting time-stamps. The sampling function of TFR used in the experiments is described in Section 12.3. Fig. 12-6(f-h) shows the reconstructed signals using these three adaptive schemes where most of the original signal features are preserved.

Although subjective distortion metrics have been proposed in the literature specifically designed for ECG signals, here we use two distortion metrics, MSE and PRD, defined as follows:

$$MSE = \frac{1}{N} \sqrt{\sum_{n=1}^N (X(n) - \hat{X}(n))^2},$$

$$PRD(\%) = \sqrt{\frac{\sum_{n=1}^N (X(n) - \hat{X}(n))^2}{\sum_{n=1}^N X^2(n)}} \times 100\%.$$

Different distortion levels for uniform sampling and DCT-based compression is performed by adjusting the threshold of the encoding process, resulting in a variable number of DCT coefficients, while for the nonuniform scheme, the rate distortion curve is derived by using different number of quantization levels. For TFR, the curve is derived by modifying the rates of each sampling region while for TPE and TIV, D and $(th1, th2)$ parameters are tuned respectively. For the rate-distortion comparison

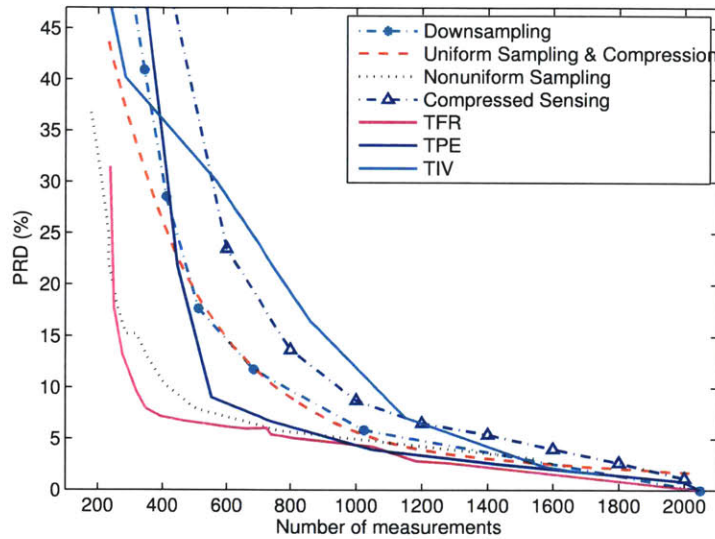


Figure 12-7: PRD in terms of the transmitted number of samples for each technique.

of these schemes, several ECG signals were used, each of approximately an hour long duration. As illustrated in Fig. 12-7 and 12-8, TFR consistently outperforms the other techniques in terms of distortion across all rate regimes. For instance, at a PRD of 10%, the approximate rate improvement of TFR is a factor of 1.2, 2.5, 2.5 and 3 compared to a uniform downsampling, a DCT-based compression scheme, a level-crossing nonuniform sampling scheme and compressed sensing, respectively.

12.5 Energy Analysis of TANS

12.5.1 Implementation Considerations of TANS

A high-level block diagram of a system implementing TANS is shown in Fig. 12-9. It is composed of a nonuniform ADC and a block implementing the sampling function, called timing control logic (TCL). The nonuniform ADC shown in the block diagram is a standard ADC with the only difference that the sample spacing is not even but rather is dynamically adapted. Thus, triggering the ADC to capture a sample is not performed by an oscillating clock, as in a uniform ADC, but by the TCL block. This means that at the acquisition stage, TANS does not need extra circuitry, such as

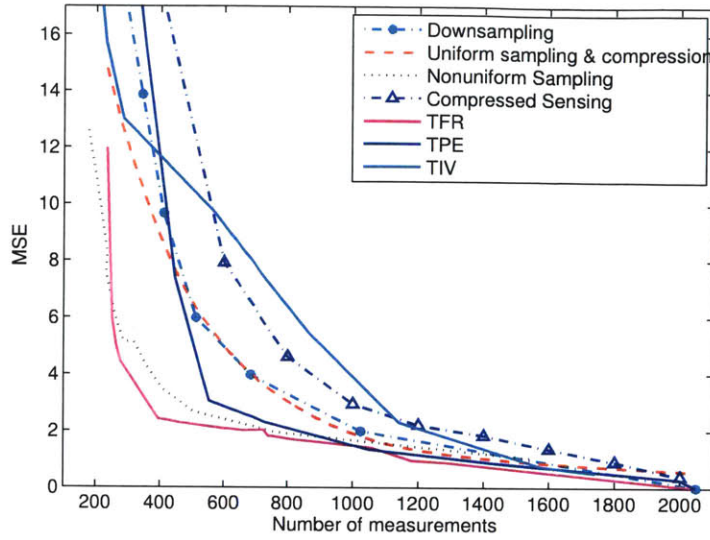


Figure 12-8: MSE in terms of transmitted number of samples for each technique.

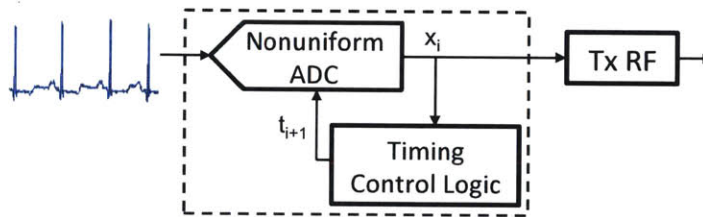


Figure 12-9: High-level block diagram of a system implementing TANS framework. A nonuniform ADC is controlled by a timing control logic, implementing the sampling function depending on the signal and sampling model.

input voltage comparators, integrators and time-to-voltage converters which are used in nonuniform sampling schemes such as level crossing. Thus, no additional energy dissipation is consumed at the acquisition stage.

Depending on the memory (m) of the sampling function, the TCL block stores the m most recently taken samples and determines the next sampling time. The block diagram of Fig. 12-9 represents the general case of implementing the proposed adaptive sampling framework where different variations of the TCL block are used for different TANS schemes. For instance, the TCL block in the case of TANS with incremental variations (TIV) calculates $w(t_i)$ using an adder and a divider, compares T_i with two predefined thresholds and adjusts its value accordingly, as in Equation (8). In the case of TANS with constrained finite rates (TFR), the TCL block implements a

finite state machine using $\log_2(k)$ registers, indicating the current sampling rate (R_i). Three accumulators and two comparators are used to calculate/store the current values of the length of each state (l_i) and to trigger a state transition, respectively.

12.5.2 Energy Comparison and Evaluation

In this part, we use a detailed system-level analysis in order to examine and compare the energy efficiency of different sampling and compression schemes. In our comparisons, state-of-the-art published results are used for each individual component of these schemes, such as ADC, compression algorithms, *etc.* Although our analysis is not an exhaustive comparison of all sampling schemes that have been proposed so far in the literature, it illustrates the advantages of TANS in terms of its energy efficiency.

In this section, TANS with finite rates (TFR) is compared against the following methods:

- **Uniform sampling and DCT-based compression:** A high-level block diagram of this scheme is shown in Fig. 12-10(a), where uniformly sampled data at fixed sampling frequency F_s are compressed by a DSP block implementing the DCT transform and the thresholding process, described in Section IV-A.
- **Level-crossing nonuniform sampling:** Fig. 12-10(b) shows the block diagram of a level-crossing sampling system using a nonuniform ADC, which can be viewed as a joint sampling and compression stage. When the local bandwidth or activity of the signal is low, automatically less samples are taken. To achieve this, the voltage difference of the input signal with the previously taken sample is compared continuously by an analog-front-end (AFE) circuit. When the difference exceeds a threshold value, the ADC is triggered and a sample (x_i) is taken, along with its time-stamp (t_i). This extra circuitry which tracks the input signal continuously, needs to be powered on all the time which consumes power constantly. Implementation challenges as well as simulation results are described in detail in [368], where a low power, low sampling frequency level

crossing ADC is designed and optimized.

- **Compressed Sensing (CS):** Two main approaches exist for the implementation of the CS block shown in Fig. 12-10(c). According to the first, a signal is sampled with a rate above the Nyquist limit by a standard uniform ADC and then multiplied by the CS measurement matrix, as explained in Section 12.4, resulting in the transmission of a number of samples ($M\lambda N$) corresponding to sub-Nyquist sampling. This approach uses CS as a signal-agnostic compression technique and, because of its low complexity requirements, has been used in the design of low power systems [353, 369].

The second approach targets directly capturing signals at sub-Nyquist rates. Examples of this approach include the analog-to-information converters (AIC) and random modulators [370, 371]. As explained in [372], for low sampling frequency applications, similar to the one considered in our case, these schemes consume higher power mainly because of their implementation requirements of multiple parallel branches with mixers and amplifiers, which makes them more suitable for high-speed applications. Demonstration of a low power implementation of these schemes, competitive to the original CS approaches, is still an ongoing effort. Thus, relying on the arguments of references [369] and [372], although the sampling rate is lower compared to the previously described CS method, the overall energy consumption tends to be higher and therefore it is not included in our comparisons.

In our analysis, we assume that all schemes sample the same ECG signal, use a compression ration (CR) which results in a target PRD value of 10%, and use the same transmission block². The total energy consumption for each scheme is:

$$E_{tot} = E_{samp} + E_{compr} + E_{TX}, \quad (12.14)$$

where E_{samp} , E_{compr} and E_{TX} are the energy consumption of the sampling, compres-

²As a reference, we assume that transmitting a bit of information consumes: $E_{bit} = 10pJ/bit$ in the transmission block [373] and each transmitted sample is quantized using: $B_{quant} = 10$ bits.

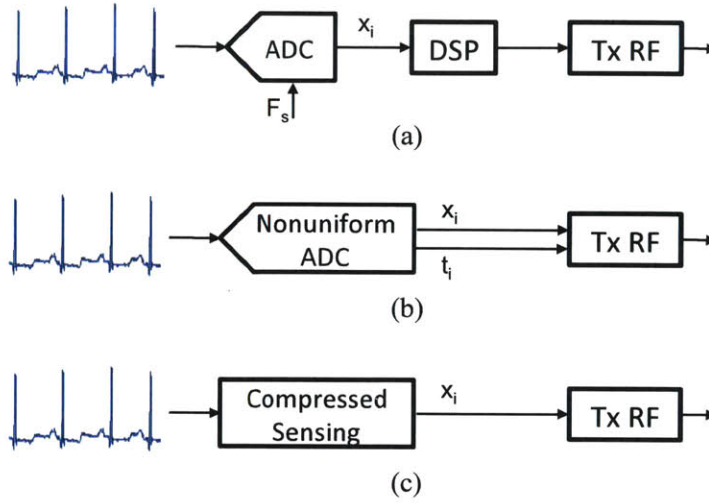


Figure 12-10: (a) A system using uniform sampling followed by transformation-based compression. (b) A level-crossing nonuniform sampling system. (c) A compressed sensing system.

sion and transmission stages, respectively. In addition, we assume that the energy consumption at the acquisition stage is approximately proportional to the sampling frequency, ignoring any leakage component. So, E_{samp} can be written as:

$$E_{samp} = R \cdot E_{samp}^{ref} = R \cdot P_{samp}^{ref} \cdot t, \quad (12.15)$$

where R is the sampling rate and $E_{samp}^{ref}/P_{samp}^{ref}$ is a reference energy/power consumption. Based on an ultra low power ADC optimized for low sampling rate applications [374], we assume that a reference ADC power consumption is 1.6uW for 1kHz sampling frequency. For the purpose of this comparison, we assume that the time period (t) equals 6 seconds, which corresponds to the signal duration of Fig. 12-6.

The uniform sampling followed by DCT-based compression and CS schemes sample the signal at 360 samples/sec, thus their energy consumption is:

$$E_{samp}^{unif} = E_{samp}^{CS} = (0.36 \cdot 1.6 \cdot 6)\mu J = 3.45\mu J. \quad (12.16)$$

Using reported results from [368], the power consumed at the acquisition stage of a low-power level-crossing nonuniform scheme is approximately 23uW, therefore its

consumed energy at the acquisition stage is:

$$E_{samp}^{lev} = (23 \cdot 6)\mu J = 138\mu J. \quad (12.17)$$

TFR achieves the same PRD by taking fewer samples. According to Fig. 12-7, its sampling rate for the target PRD is approximately 50 samples/sec, since a compression ratio (CR) of approximately 7 is achieved, thus its consumed sampling energy is:

$$E_{samp}^{TFR} = (0.05 \cdot 1.6 \cdot 6)\mu J = 0.48\mu J. \quad (12.18)$$

The data compression blocks of the considered schemes have different computational requirements and, therefore, different energy consumptions (E_{compr}). Based on [375], a low-power implementation of a DCT-based compression algorithm consumes approximately $166\mu W$, while in [290] the power consumption of CS block is reported to be around $1.9\mu W$. Thus, their energy consumption for our ECG signal comparison would be:

$$E_{compr}^{unif} = 1mJ \text{ and } E_{compr}^{CS} = 11.4\mu J \quad (12.19)$$

Compression in TANS schemes is performed by the sampling function which calculates the next sampling time. We use a hardware description language (HDL) to model the sampling function described in Section V-A and, using the standard hardware design flow³, we measure the energy consumption of the TCL block of TFR to be:

$$E_{compr}^{TFR} = 6\mu J. \quad (12.20)$$

Finally, the transmission energy is given by:

$$E_{TX} = E_{bit} \cdot \left(\frac{R}{CR} \cdot B_{quant} \cdot t + H \right), \quad (12.21)$$

³We model the sampling function using Verilog and then we proceed with synthesis, place-and-route, extraction and post-layout simulation. CAD tools from Cadence and Synopsys are used, as well as libraries from a 65nm TSMC process.

Table 12.1: A energy comparison of 3 different schemes for sampling, compressing and transmitting a 6-sec ECG signal.

	Uniform Sampling and Compression	Nonuniform Sampling	Compressed Sensing	TFR
ADC	3.45 μ J	138 μ J	3.45 μ J	0.48 μ J
Compression	1 mJ	0	11.4 μ J	6 μ J
Transmission	80 μ J	44 μ J	80 μ J	30 μ J
Total	1083.45 μ J	182 μ J	102.45 μ J	36.48 μ J

where E_{bit} is the energy per bit, B_{quant} is the number of quantization bits and H is the overhead information in bits used in the DCT thresholding process indicating the indexes of non-zero coefficients, and in the level-crossing scheme representing the time-stamps required for the reconstruction process. For the CS and TFR schemes, H equals zero.

Table I shows the energy components for each sampling scheme as well as their total energy consumption. It can be seen that, TFR outperforms all approaches considered in the analysis. Specifically, TFR provides a benefit of at least an order of 3 in terms of the total energy consumption compared to state-of-the-art sampling approaches. This energy efficiency benefit is achieved by adapting the sampling rate without requiring extra circuitry at the acquisition stage and by avoiding the transmission of the time-stamps. In addition, the energy overhead of the TCL block, which controls the ADC, is significantly lower when compared to standard block-based compression algorithms.

12.6 Proofs

12.6.1 Proof of Theorem 12.1

By Taylor's theorem, for each t in $[t_i, t_{i+1})$, there exists s_1 in (t_i, t) such that,

$$\begin{aligned}
 X(t) &= X(t_i) + X'(t_i)(t - t_i) + \frac{1}{2}X''(t_i)(t - t_i)^2 \\
 &\quad + \frac{1}{6}X'''(s_1)(t - t_i)^3.
 \end{aligned}
 \tag{12.22}$$

Since $|X'''(t)| \leq M$ for all t , we have $X_l(t) \leq X(t) \leq X_u(t)$ where,

$$\begin{aligned} X_u(t) &= X(t_i) + X'(t_i)(t - t_i) + \frac{1}{2}X''(t_i)(t - t_i)^2 \\ &\quad + \frac{1}{6}M(t - t_i)^3 \\ X_l(t) &= X(t_i) + X'(t_i)(t - t_i) + \frac{1}{2}X''(t_i)(t - t_i)^2 \\ &\quad - \frac{1}{6}M(t - t_i)^3. \end{aligned} \tag{12.23}$$

For the reconstruction, according to assumption A_2 , we connect adjacent samples by a line. Hence, for $t \in [t_i, t_{i+1})$, we have,

$$\hat{X}(t) = X(t_i) + (t - t_i)\frac{\Delta_i}{T_i}. \tag{12.24}$$

By Taylor's theorem, there exists $s_2 \in (t_i, t_{i+1})$ such that,

$$\begin{aligned} X(t_{i+1}) &= X(t_i) + X'(t_i)(t_{i+1} - t_i) \\ &\quad + \frac{1}{2}X''(t_i)(t_{i+1} - t_i)^2 + \frac{1}{6}X'''(s_2)(t_{i+1} - t_i)^3. \end{aligned}$$

Hence,

$$\frac{\Delta_i}{T_i} = X'(t_i) + \frac{1}{2}X''(t_i)T_i + \frac{1}{6}X'''(s_2)T_i^2. \tag{12.25}$$

Then,

$$\begin{aligned}
 e(t) &\triangleq X(t) - \hat{X}(t) \\
 &\stackrel{(a)}{=} (t - t_i)(X'(t_i) - \frac{\Delta_i}{T_i}) + \frac{1}{2}X''(t_i)(t - t_i)^2 \\
 &\quad + \frac{1}{6}X'''(s_1)(t - t_i)^3 \\
 &\stackrel{(b)}{=} \frac{1}{2}X''(t_i)(t - t_i)(-T_i + t - t_i) \\
 &\quad + \frac{1}{6}(t - t_i)[X'''(s_1)(t - t_i)^2 - X'''(s_2)T_i^2].
 \end{aligned}$$

Equality (a) follows from (12.22) and (12.24). Equality (b) follows from (12.25). Then, by using the bounded third derivative assumption, for $t \in [t_i, t_{i+1})$, we have,

$$\begin{aligned}
 |e(t)| &\leq \frac{1}{2}|X''(t_i)|(t - t_i)(t - t_i + T_i) \\
 &\quad + \frac{1}{3}M(t - t_i)((t - t_i)^2 + T_i^2)
 \end{aligned} \tag{12.26}$$

The right-hand side of (12.26) is an increasing function of t for $t \in [t_i, t_{i+1})$. Hence, its supremum is achieved at $t = t_{i+1}$. Thus, we have,

$$(|X''(t_i)| + \frac{2}{3}MT_i)T_i^2 \leq D. \tag{12.27}$$

By using Taylor's theorem at times $t = t_{i-1}$ and $t = t_{i-2}$, centered at $t = t_i$, there exist s_3 and s_4 such that $t_{i-1} < s_3 < t_i$ and $t_{i-2} < s_4 < t_i$, and,

$$\begin{aligned}
 \frac{X(t_{i-1}) - X(t_i)}{T_{i-1}} &= -X'(t_i) + \frac{1}{2}X''(t_i)T_{i-1} \\
 &\quad - \frac{1}{6}X'''(s_3)T_{i-1}^2
 \end{aligned} \tag{12.28}$$

$$\begin{aligned} \frac{X(t_{i-2}) - X(t_i)}{T_{i-1} + T_{i-2}} &= -X'(t_i) + \frac{1}{2}X''(t_i)(T_{i-1} + T_{i-2}) \\ &\quad - \frac{1}{6}X'''(s_4)(T_{i-1} + T_{i-2})^2. \end{aligned} \quad (12.29)$$

Hence,

$$\begin{aligned} |X''(t_i)| &\leq \frac{\left| \frac{\Delta_{i-1}}{T_{i-1}} - \frac{\Delta_{i-2} + \Delta_{i-1}}{T_{i-2} + T_{i-1}} \right|}{\frac{T_{i-2}}{2}} \\ &\quad + \frac{1}{3}M \frac{(T_{i-1} + T_{i-2})^2 + (T_{i-1})^2}{T_{i-2}} \\ &= \frac{\left| \frac{\Delta_{i-1}}{T_{i-1}} - \frac{\Delta_{i-2}}{T_{i-2}} \right|}{\frac{T_{i-1} + T_{i-2}}{2}} \\ &\quad + \frac{1}{3}M \frac{(T_{i-1} + T_{i-2})^2 + (T_{i-1})^2}{T_{i-2}}. \end{aligned} \quad (12.30)$$

Combining (12.27) and (12.30) completes the proof.

12.6.2 Proof of Theorem 12.2

By using (12.26), we can write,

$$\begin{aligned} \int_{t=t_i}^{t_{i+1}} |e(t)| dt &\leq \frac{1}{2}|X''(t_i)| \int_{t=t_i}^{t_{i+1}} (t - t_i)(t - t_i + T_i) dt \\ &\quad + \frac{1}{3}M \int_{t=t_i}^{t_{i+1}} (t - t_i)((t - t_i)^2 + T_i^2) dt \\ &= \left(\frac{5}{12}|X'''(t_i)| + \frac{1}{4}MT_i\right)T_i^3. \end{aligned} \quad (12.31)$$

Combining (12.31) and (12.30) completes the proof.

Chapter 13

Conclusion and Future Directions

Network science is an interdisciplinary research area that gives a theoretical approach to applied problems faced by practitioners in biological and other sciences. It is well suited to make significant progress in bringing together advanced mathematical tools and biological/experimental challenges that will leave a lasting impact on both fields. The focus of this thesis was to develop and apply methods for better understanding key challenges in inference and analysis of complex networks.

In Chapters 2-4, we considered network inference challenges where we aim to learn dependencies among variables represented by the network topology:

- In Chapter 2, we proposed Network Maximal Correlation (NMC) as a measure to capture nonlinear associations among variables in large datasets. We showed that NMC extends the standard bivariate Maximal Correlation to the case of having large number of variables, by assigning each variable to a single transformation function, thus avoiding over-fitting issues of using multiple MC optimizations over variable pairs. We also introduced a regularized NMC optimization which penalizes total distances of inferred transformed variables from the original ones. One can use other standard regularization techniques to further restrict inferred nonlinear functions in practical applications.

One of the main contributions of this work is providing a unifying framework to compute NMC (and therefore, MC) for both discrete and continuous variables using projections over appropriate Hilbert spaces. Using this framework, we established a connection between the NMC optimization with the MCP and MEP for discrete random variables, and with the Max-Cut problem for jointly Gaussian variables. Using these relationships, we provided efficient algorithms to compute NMC in different cases. To compute NMC for continuous random variables with general distributions, one can use the proposed optimization framework by choosing appropriate orthonormal basis for Hilbert spaces. For example, we used projections over Hermite-Chebyshev polynomials to characterize an optimal solution of the NMC optimization for jointly Gaussian variables. Finally note that the inferred, possibly nonlinear, functions of the NMC

optimization can be used in other applications such as nonlinear regression.

- In Chapter 3, we developed a network integration framework to infer regulatory interactions in human, fly and worm organisms by integrating genome-wide functional (gene expression) and physical (motif and ChIP) regulatory evidences. We integrated edge weights of feature-specific networks using a rank-based integration approach which combines edge ranks of input networks. We found that inferred interactions have significant overlap with known benchmarks, especially in human and fly, indicating the accuracy and robustness of the integration framework used. Moreover, we found that many biological processes and pathways are significantly enriched and conserved over the inferred networks.

The inferred interactions in human, fly and worm can be used as a potential guide to design informative high-throughput experiments to infer accurate regulatory interactions in different contexts and reduce their false positive edges. For instance, our inferred interactions can guide expensive ChIP assay experiments in order to assess more reliable physical links between TFs and target genes, eventually leading to construct higher-quality regulatory networks, possibly cell-type specific. Finally, inferred regulatory interactions, specially conserved ones, provide useful resources to understand regulatory roles of individual genes, pathways and modules in human diseases, cancer and drug designs.

- In Chapter 4, we introduce network deconvolution which provides a framework for computing direct dependencies in a network by use of observed similarities. It can recognize and remove spurious transitive edges due to indirect effects, decrease edge weights that are overestimated due to indirect relationships, and assign edge weights corresponding to direct dependencies to the remaining edges. Thereby, network deconvolution can improve the quality of a broad range of observed networks that are tainted by indirect edge weights due to transitive effects. We introduced an efficient and scalable algorithm for deconvolving an observed network based on a nonlinear filter computing the inverse of a Tay-

lor series expansion over each eigenvalue. We demonstrated the effectiveness of network deconvolution for gene regulatory network inference, protein folding based on protein sequence alignment, and social networks of co-authorship information.

The problem of indirect spurious edges has been widely recognized in network inference, but characterized mostly at the local level. In particular, even top-performing network inference methods have been shown to contain many false transitive edges in cascade network motifs, and efforts to remedy this situation lead to incorrect removal of true edges in feed-forward loops. At this local level, we have shown that network deconvolution has the ability to correctly remove spurious transitive edges in true cascade network motifs, while maintaining true feed-forward edges in feed-forward network motifs. In contrast to previous methods that make well-documented tradeoffs in sensitivity vs. specificity for these transitive edges, network deconvolution reduces the number of false positives on indirect interactions, while maintaining true positives in feed-forward loops.

In contrast to previous approaches that study local patterns of dependencies to recognize potential indirect edges, network deconvolution takes a global approach by directly inverting the transitive closure of the true network. Previous algorithms [135] have sought local approximations to the removal of indirect effects which have been limited to indirect paths of only limited lengths (typically of length 2), due to the computational complexity of enumerating and evaluating all higher-order paths, and the lack of a systematic way to compute their combined effects. By exploiting eigenvector decomposition and Taylor series closed form solutions, network deconvolution provides four advantages over local approaches: (1) it leads to a much more computationally efficient solution, (2) it has the power to remove indirect effects over paths with arbitrary lengths, (3) it can remove the combined effects of arbitrarily many indirect paths between two nodes, and (4) it eliminates the need for iterative network

refinement. These advantages are due to the fact that network deconvolution is essentially a single global operation to subtract the transitive effects of all powers of an adjacency matrix, rather than testing only pair-wise relationships or small network motifs one at a time.

Moreover, we showed that network deconvolution can be applied to networks with very different properties. The networks used were of different size, density, clustering coefficient, or network centrality, showing that network deconvolution is robust to these parameters. The input networks were also based on different properties, including mutual information and correlation that network deconvolution was designed for, but also networks based on regression, tree-based ensemble methods, feature selection approaches, and other non-linear similarity metrics. We also applied network deconvolution to both weighted and unweighted networks, and used the results both for re-weighting of edges and for edge classification, demonstrating the discrete and continuous applications of the approach. More generally, network deconvolution is not just about edge inclusion or removal, but about probabilistic weighing of individual edges to reveal direct interactions based on observed relationships across the complete network.

In Chapters 5-8, we considered network analysis challenges where we aim to learn from network topologies:

- In Chapter 5, we made a principled connection between spectral network alignment techniques and relaxations of the network alignment optimization (quadratic assignment problem). Using such connection, we proposed a network alignment algorithm which employs an orthogonal relaxation of the underlying QAP in a maximum weight bipartite matching optimization. Our method simplifies, in a principled way, the network alignment optimization to simultaneous alignment of eigenvectors of (transformations of) adjacency graphs scaled by corresponding eigenvalues. Our framework also provides a theoretical justification for other existing heuristic spectral network alignment methods. Our proposed

framework advances existing network alignment methods not only in algorithmic aspects, but also in qualitative terms of the network alignment objective function. Our formulation considers both matched and mismatched interactions in its optimization and therefore, it is effective in aligning networks even with low similarity. This is critical in applications where networks have low similarities such as comparative analysis of biological networks of distal species. This idea can also be adapted to existing network alignment packages.

For Erdős-Rényi graphs, we proved that the proposed algorithm is asymptotically optimal with high probability, under some general conditions. Through simulations, we compared the performance of our algorithm with the one of existing network alignment methods based on belief propagation (NetAlign), spectral decomposition (IsoRank), Lagrange relaxation (Klau optimization), and a SDP-based method. Our simulations illustrated the effectiveness of our proposed algorithm in aligning various network structures such as Erdős-Rényi, power law, regular, and stochastic block structures, under different noise models. For modular graph structure, we proposed a framework that uses spectral network alignment to split the large network alignment optimization into small subproblems. This enables the use of computationally expensive, but tight, semidefinite programming relaxations over each subproblem. This hybrid method has high performance similar to SDP, with significantly less computational complexity. Designing other SDP-based network alignment methods with low computational complexity for various network models remains a promising direction for future work. Moreover, in our network alignment framework one can use other relaxations of QAP (for example by considering the alignment matrix in the intersection of the sets of orthogonal and stochastic matrices) to obtain tighter bounds.

We also considered two real-data applications. We applied our proposed network alignment algorithm to compare gene regulatory networks across human, fly and worm species which we inferred by integrating genome-wide functional

and physical genomics datasets from ENCODE and modENCODE consortia. Our method inferred conserved regulatory interactions across these species despite large evolutionary distances spanned. Moreover, we found strong conservation of centrally-connected genes and biological pathways, especially for human-fly comparisons. In a second application, we used network alignment in user de-anonymization over twitter follower subgraphs sampled in two different years. This application illustrates the extent of personal information that can be retrieved from network structures, and raises additional considerations that need to be addressed in different privacy-related applications.

- In Chapter 6, we considered the source inference problem which aims to infer source nodes in a network by merely knowing the network structure and observing the information spread at single or multiple snapshots. The key idea of our proposed source inference method is designing diffusion processes that closely approximate the observed diffusion pattern, while leading to tractable source inference methods for large complex networks. Our proposed path-based network diffusion kernel well-approximates a standard SI diffusion model particularly over sparse networks. The key intuition is that to infer the source node in the network, full characterization of diffusion dynamics, in many cases, may not be necessary. One advantage of our proposed path-based network diffusion kernel is its tunability, that it can consider different number of shortest paths in kernel computation. This resembles, vaguely, a Taylor-expansion of network topology to form a diffusion kernel with different orders of expansion. One can extend this key idea to design other network diffusion kernels to approximate other general diffusion models such as SIR (susceptible-infected-recovered) [238], or to design network-specific diffusion kernels considering different topological properties such as their symmetry, degree distribution, etc.

The NI framework infers source nodes using the network topology and snapshots from the infection spread. On the other hand, [376, 377] have considered the problem of network inference given propagation pathways over the network.

If the network topology is partially unknown or has some errors with either false positive or false negative edges, one can devise a *joint network inference-network infusion* framework where in one direction, the pattern of infection spread is used to learn the network topology by denoising false positive and false negative edges, while in the other direction, the topology of the network is used to infer pathways of infection spread over the network.

- In Chapter 7, we developed a spectral inference method over logistic Random Dot Product Graphs (RDPGs), and showed that the proposed method is asymptotically equivalent to the maximum likelihood latent-position inference. Previous justifications for spectral clustering have usually been either consistency results [195,288] or partial-recovery results [282,300]; to the best of our knowledge, our likelihood-based justification is the first of its kind for a spectral method. This type of justification is satisfying because maximum likelihood inference methods can generally be expected to have optimal asymptotic performance characteristics; for example, it is known that maximum likelihood estimators are consistent over the SBM [378,379]. It remains an important future direction to characterize the asymptotic performance of the MLE over the Logistic RDPG. We have focused in this chapter on the network clustering problem; however, latent space models such as the Logistic RDPG can be viewed as a more general tool for exploring and analyzing network structures. They can be used for visualization [380,381] and for inference of partial-membership type structures, similar to the mixed-membership stochastic blockmodel [266]. Our approach can also be generalized to multi-edge graphs, in which the number of edges between two nodes is binomially distributed. Such data is emerging in areas including systems biology, in the form of cell type and tissue specific networks [382].
- In Chapter 8, we have proposed a new biclustering technique called *Biclustering Using Message Passing* that, unlike existent methods, infers a globally optimal collection of biclusters rather than a collection of locally optimal ones. This

distinction is especially relevant in the presence of overlapping clusters, which are common in most applications. Such overlaps can be of importance if one is interested in the relationships among biclusters. We showed through simulations that our proposed method outperforms two popular existent methods, ISA and LAS, in both Bernoulli and Gaussian noise models, when the planted biclusters were overlapping. We also found that BCMP performed well when applied to gene expression datasets. The same optimization strategy used in this chapter for biclustering can also be applied to perform clustering, generalizing the graph-partitioning problem by allowing nodes to be in zero or several clusters. We believe that the flexibility of our framework to fit various statistical and heuristic models will allow BCMP to be used in diverse clustering and biclustering applications.

Finally in Chapters 9-12, we considered the design challenges of systems and algorithms for engineered networks using Compressive Sensing, Signal Processing and Information Theory:

- In Chapter 9, we demonstrated some applications of compressive sensing over networks. We made a connection between compressive sensing and traditional information theoretic techniques in source coding and channel coding in order to provide mechanisms to explicitly trade-off between the decoding complexity and the rate. Although optimal decoders to recover the original signal compressed by source coding have high complexity, the compressive sensing decoder is a linear or convex optimization. First, we investigated applications of compressive sensing on distributed compression of correlated sources. For a depth one tree network, references ([317]- [322]) provide some low complexity decoding techniques for Slepian-Wolf compression. However, none of these methods can explicitly provide a trade-off between the rate and the decoding complexity. For a general multicast network, reference [315] extended Slepian-Wolf compression for a multicast problem. It is showed in [323] that, there is no separation between Slepian-Wolf compression and network coding. For this network, a low

complexity decoder for an optimal compression scheme has not been proposed yet. Here, by using compressive sensing, we proposed a compression scheme for a family of correlated sources with a modularized decoder providing a trade-off between the compression rate and the decoding complexity. We called this scheme *Sparse Distributed Compression*. We used this compression scheme for a general multicast network with correlated sources. Here, we first decoded some of the sources by a network decoding technique and then, we used a compressive sensing decoder to obtain the whole sources. Next, we investigated applications of compressive sensing on channel coding. We proposed a coding scheme that combines compressive sensing and random channel coding for a high-SNR point-to-point Gaussian channel. We called this scheme *Sparse Channel Coding*. Our coding scheme provides a modularized decoder to have a trade-off between the capacity loss and the decoding complexity. The idea is to add intentionally some correlation to transmitted signals in order to decrease the decoding complexity. At the decoder side, first we used a compressive sensing decoder to get an estimate of the original sparse signal, and then we used a channel coding decoder in the subspace of non-zero coordinates.

- In Chapter 10, we proposed a joint source-channel-network coding scheme, based on compressive sensing principles, for wireless networks with AWGN channels (that may include multiple access and broadcast). A key idea of our proposed scheme is to reshape inherent redundancies among source samples (temporal or spatial) in joint source-channel-network coding to have a power/rate efficient sensing communication scheme. By characterizing the performance of the proposed scheme, we showed that, our proposed scheme (1) is power and rate efficient, that is, although it performs closely to optimal rates, it provides power efficiencies, (2) has low decoding complexity where decoders are convex optimizations, (3) is broadly structure independent, that is, sensing, coding and communication processes in sensors are application independent, (4) and has continuous rate-distortion performance, that is, if channel qualities change, the

reconstruction quality will change in a continuous manner. In this framework, we perform joint source-channel coding at each source by randomly projecting source values to a lower dimensional space. We consider sources that spatially satisfy the restricted eigenvalue (RE) condition and may use sparse networks as well as more general sources which use the randomness of the network to allow mapping to lower dimensional spaces. The receiver uses compressive sensing decoders to reconstruct source signals.

- In Chapter 11, we introduced a new framework for an adaptive nonuniform sampling scheme called *time-stampless adaptive nonuniform sampling* (TANS). The key idea of this framework is that, *time increments between samples are computed by using a function of previously taken samples*. Therefore, keeping sampling times (time stamps), except initialization times, is not necessary. We introduced two methods to design sampling functions for discrete-time stochastic signals: a greedy method, and a method based on dynamic programming. We analyzed the performances of these schemes by computing (or bounding) their trade-offs between sampling rate and expected reconstruction distortion for autoregressive and Markovian signals. We showed that, by being time-stampless and opportunistically adapting to local signal characteristics, TANS can provide significant rate-distortion gains, which can be translated to improved power efficiency in some applications.
- In Chapter 12, we analyzed rate and energy efficiencies of a recently proposed sampling framework, TANS, for different signal models. This framework reduces sampling rate and therefore energy consumption of the sampling process. In this framework, time intervals between samples can be computed by using a function of previously taken samples, called sampling function. We proposed three schemes of TANS for different signal models. The first two are designed for signals with bounded derivatives, TANS by polynomial extrapolation and TANS by incremental variation methods. Then, we introduced a scheme of TANS with constrained finite rates. We compared the performance in terms

of MSE and PRD of TANS with uniform sampling, level-crossing nonuniform sampling and compressed sensing. We showed that our approach outperforms all the compared approaches we considered, both in terms of rate and energy consumed.

Bibliography

- [1] H. Liu, F. Han, M. Yuan, J. Lafferty, L. Wasserman *et al.*, “High-dimensional semiparametric gaussian copula graphical models,” *The Annals of Statistics*, vol. 40, no. 4, pp. 2293–2326, 2012.
- [2] M. E. Newman, “Modularity and community structure in networks,” *Proceedings of the National Academy of Sciences*, vol. 103, no. 23, pp. 8577–8582, 2006.
- [3] D. Marbach, J. C. Costello, R. Küffner, N. M. Vega, R. J. Prill, D. M. Camacho, K. R. Allison, M. Kellis, J. J. Collins, G. Stolovitzky *et al.*, “Wisdom of crowds for robust gene network inference,” *Nature methods*, vol. 9, no. 8, pp. 796–804, 2012.
- [4] R. De Smet and K. Marchal, “Advantages and limitations of current network inference methods,” *Nature Reviews Microbiology*, vol. 8, no. 10, pp. 717–729, 2010.
- [5] T. Ideker, T. Galitski, and L. Hood, “A new approach to decoding life: systems biology,” *Annual review of genomics and human genetics*, vol. 2, no. 1, pp. 343–372, 2001.
- [6] B. J. Frey and D. Dueck, “Clustering by passing messages between data points,” *science*, vol. 315, no. 5814, pp. 972–976, 2007.
- [7] M. E. Newman, “The structure and function of complex networks,” *SIAM review*, vol. 45, no. 2, pp. 167–256, 2003.

BIBLIOGRAPHY

- [8] R. Koetter and M. Medard, "An algebraic approach to network coding," *IEEE/ACM Trans. on networking*, vol. 11, no. 5, pp. 782–795, 2003.
- [9] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger, "Factor graphs and the sum-product algorithm," *Information Theory, IEEE Transactions on*, vol. 47, no. 2, pp. 498–519, 2001.
- [10] R. Castro, M. Coates, G. Liang, R. Nowak, and B. Yu, "Network tomography: recent developments," *Statistical science*, pp. 499–517, 2004.
- [11] J.-F. Rual, K. Venkatesan, T. Hao, T. Hirozane-Kishikawa, A. Dricot, N. Li, G. F. Berriz, F. D. Gibbons, M. Dreze, N. Ayivi-Guedehoussou *et al.*, "Towards a proteome-scale map of the human protein–protein interaction network," *Nature*, vol. 437, no. 7062, pp. 1173–1178, 2005.
- [12] A. L. Hopkins, "Network pharmacology: the next paradigm in drug discovery," *Nature chemical biology*, vol. 4, no. 11, pp. 682–690, 2008.
- [13] V. R. Chintapalli, J. Wang, and J. A. Dow, "Using flyatlas to identify better drosophila melanogaster models of human disease," *Nature genetics*, vol. 39, no. 6, pp. 715–720, 2007.
- [14] P. Csermely, V. Agoston, and S. Pongor, "The efficiency of multi-target drugs: the network approach might help drug design," *Trends in Pharmacological Sciences*, vol. 26, no. 4, pp. 178–182, 2005.
- [15] H.-Y. Chuang, E. Lee, Y.-T. Liu, D. Lee, and T. Ideker, "Network-based classification of breast cancer metastasis," *Molecular systems biology*, vol. 3, no. 1, 2007.
- [16] M. Granovetter, *The strength of weak ties: A network theory revisited*. JSTOR, 1981.
- [17] H. Liu, J. Lafferty, and L. Wasserman, "The nonparanormal: Semiparametric estimation of high dimensional undirected graphs," *The Journal of Machine Learning Research*, vol. 10, pp. 2295–2328, 2009.

BIBLIOGRAPHY

- [18] R. Albert, “Network inference, analysis, and modeling in systems biology,” *The Plant Cell*, vol. 19, no. 11, pp. 3327–3338, 2007.
- [19] M. S. Handcock, *Statistical models for social networks: Inference and degeneracy*. National Academies Press, 2002.
- [20] A. G. Haldane, “Rethinking the financial network,” *Speech delivered at the Financial Student Association, Amsterdam, April*, pp. 1–26, 2009.
- [21] K. Pearson, “Note on regression and inheritance in the case of two parents,” *Proceedings of the Royal Society of London*, pp. 240–242, 1895.
- [22] B. Thompson, “Canonical correlation analysis,” *Encyclopedia of statistics in behavioral science*, 2005.
- [23] L. L. Scharf and C. T. Mullis, “Canonical coordinates and the geometry of inference, rate, and capacity,” 2000.
- [24] C. Spearman, “The proof and measurement of association between two things,” *The American journal of psychology*, vol. 15, no. 1, pp. 72–101, 1904.
- [25] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York, NY: John Wiley and Sons, 1991.
- [26] D. N. Reshef, Y. A. Reshef, H. K. Finucane, S. R. Grossman, G. McVean, P. J. Turnbaugh, E. S. Lander, M. Mitzenmacher, and P. C. Sabeti, “Detecting novel associations in large data sets,” *Science*, vol. 334, no. 6062, pp. 1518–1524, 2011.
- [27] C. Chow and C. Liu, “Approximating discrete probability distributions with dependence trees,” *Information Theory, IEEE Transactions on*, vol. 14, no. 3, pp. 462–467, 1968.
- [28] D. Lopez-Paz, P. Hennig, and B. Schölkopf, “The randomized dependence coefficient,” in *Advances in Neural Information Processing Systems*, 2013, pp. 1–9.

BIBLIOGRAPHY

- [29] H. O. Hirschfeld, "A connection between correlation and contingency," in *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 31, no. 04, 1935, pp. 520–524.
- [30] H. Gebelein, "Das statistische problem der korrelation als variations-und eigenwertproblem und sein zusammenhang mit der ausgleichsrechnung," *ZAMM-Journal of Applied Mathematics and Mechanics*, vol. 21, no. 6, pp. 364–379, 1941.
- [31] O. Sarmanov, "Maximum correlation coefficient (nonsymmetric case)," *Selected Translations in Mathematical Statistics and Probability*, vol. 2, pp. 207–210, 1962.
- [32] A. Rényi, "On measures of dependence," *Acta Mathematica Hungarica*, vol. 10, no. 3, pp. 441–451, 1959.
- [33] H. S. Witsenhausen, "On sequences of pairs of dependent random variables," *SIAM Journal on Applied Mathematics*, vol. 28, no. 1, pp. 100–113, 1975.
- [34] R. Ahlswede and P. Gács, "Spreading of sets in product spaces and hypercontraction of the markov operator," *The Annals of Probability*, pp. 925–939, 1976.
- [35] H. O. Lancaster, "Some properties of the bivariate normal distribution considered in the form of a contingency table," *Biometrika*, vol. 44, no. 1-2, pp. 289–292, 1957.
- [36] L. Breiman and J. H. Friedman, "Estimating optimal transformations for multiple regression and correlation," *Journal of the American Statistical Association*, vol. 80, no. 391, pp. 580–598, 1985.
- [37] A. Buja, "Remarks on functional canonical variates, alternating least squares methods and ace," *The Annals of Statistics*, pp. 1032–1069, 1990.
- [38] V. Anantharam, A. Gohari, S. Kamath, and C. Nair, "On maximal correlation, hypercontractivity, and the data processing inequality studied by erkip and cover," *arXiv preprint:1304.6133*, 2013.

BIBLIOGRAPHY

- [39] Y. Polyanskiy, “Hypothesis testing via a comparator and hypercontractivity,” *Prob. Peredachi Inform*, 2013.
- [40] M. Raginsky, “Strong data processing inequalities and ϕ -sobolev inequalities for discrete channels,” *arXiv preprint:1411.3575*, 2014.
- [41] T. Courtade, “Outer bounds for multiterminal source coding via a strong data processing inequality,” in *IEEE International Symposium on Information Theory*, 2013, pp. 559–563.
- [42] F. Calmon, M. Varia, and M. Médard, “An exploration of the role of principal inertia components in information theory,” in *IEEE Information Theory Workshop (ITW)*, 2014, pp. 252–256.
- [43] F. Calmon, M. Varia, and M. Medard, “On information-theoretic metrics for symmetric-key encryption and privacy,” in *Allerton Conference on Communication, Control, and Computing*, 2014, pp. 889–894.
- [44] F. Calmon, M. Varia, M. Médard, M. M. Christiansen, K. R. Duffy, and S. Tessaro, “Bounds on inference,” in *Allerton Conference on Communication, Control, and Computing*, 2013, pp. 567–574.
- [45] S. Beigi and A. Gohari, “On the duality of additivity and tensorization,” *arXiv:1502.00827*, 2015.
- [46] O. Etesami and A. Gohari, “Maximal rank correlation,” *arXiv:1511.04191*, 2015.
- [47] F. Farnia, M. Razaviyayn, S. Kannan, and D. Tse, “Minimum HGR correlation principle: From marginals to joint distribution,” *arXiv:1504.06010*, 2015.
- [48] M. Razaviyayn, F. Farnia, and D. Tse, “Discrete rényi classifiers,” in *Advances in Neural Information Processing Systems*, 2015, pp. 3276–3284.
- [49] A. Makur and L. Zheng, “Bounds between contraction coefficients,” *arXiv:1510.01844*, 2015.

BIBLIOGRAPHY

- [50] M. T. Chu and J. L. Watterson, “On a multivariate eigenvalue problem, part I: Algebraic theory and a power method,” *SIAM Journal on Scientific Computing*, vol. 14, no. 5, pp. 1089–1106, 1993.
- [51] H. Hotelling, “The most predictable criterion.” *Journal of Educational Psychology*, vol. 26, no. 2, p. 139, 1935.
- [52] —, “Relations between two sets of variates,” *Biometrika*, pp. 321–377, 1936.
- [53] L.-H. Zhang and M. T. Chu, “Computing absolute maximum correlation,” *IMA Journal of Numerical Analysis*, vol. 32, no. 1, pp. 163–184, 2012.
- [54] P. Klein, S. A. Plotkin, and S. Rao, “Excluded minors, network decomposition, and multicommodity flow,” in *Proceedings of the twenty-fifth annual ACM symposium on Theory of computing*, 1993, pp. 682–690.
- [55] S. Rao, “Small distortion and volume preserving embeddings for planar and Euclidean metrics,” in *Proceedings of the fifteenth annual symposium on Computational geometry*, 1999, pp. 300–306.
- [56] D. R. Karger and M. Ruhl, “Finding nearest neighbors in growth-restricted metrics,” in *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, 2002, pp. 741–750.
- [57] M. R. Garey, D. S. Johnson, and L. Stockmeyer, “Some simplified NP-complete graph problems,” *Theoretical computer science*, vol. 1, no. 3, pp. 237–267, 1976.
- [58] M. X. Goemans and D. P. Williamson, “Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming,” *Journal of the ACM (JACM)*, vol. 42, no. 6, pp. 1115–1145, 1995.
- [59] M. J. Wainwright and M. I. Jordan, “Graphical models, exponential families, and variational inference,” *Foundations and Trends in Machine Learning*, vol. 1, no. 1-2, pp. 1–305, 2008.

BIBLIOGRAPHY

- [60] J. Neter, M. H. Kutner, C. J. Nachtsheim, and W. Wasserman, *Applied linear statistical models*. Irwin Chicago, 1996, vol. 4.
- [61] K. Baba, R. Shibata, and M. Sibuya, "Partial correlation and conditional correlation as measures of conditional independence," *Australian and New Zealand Journal of Statistics*, vol. 46, no. 4, pp. 657–664, 2004.
- [62] O. Banerjee, L. El Ghaoui, and A. d'Aspremont, "Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data," *The Journal of Machine Learning Research*, vol. 9, pp. 485–516, 2008.
- [63] G. A. Calin and C. M. Croce, "MicroRNA signatures in human cancers," *Nature Reviews Cancer*, vol. 6, no. 11, pp. 857–866, 2006.
- [64] Y. Li and Z. Zhang, "Potential microrna-mediated oncogenic intercellular communication revealed by pan-cancer analysis," *Scientific reports*, vol. 4, 2014.
- [65] E. M. Stein and R. Shakarchi, *Real analysis: measure theory, integration, and Hilbert spaces*. Princeton University Press, 2009.
- [66] W. Bryc and A. Dembo, "On the maximum correlation coefficient," *Theory of Probability and its Applications*, vol. 49, no. 1, pp. 132–138, 2005.
- [67] I. Goodman and D. H. Johnson, "Multivariate dependence and the sarmanov-lancaster expansion," *J. Multivariate Analysis*, 2005.
- [68] D. P. Bertsekas, *Nonlinear programming*. Athena scientific, 1999.
- [69] P. Horst, *Factor analysis of data matrices*. Holt, Rinehart and Winston Publisher, 1965.
- [70] —, "Relations among sets of measures," *Psychometrika*, vol. 26, no. 2, pp. 129–149, 1961.
- [71] L.-H. Zhang and L.-Z. Liao, "An alternating variable method for the maximal correlation problem," *Journal of Global Optimization*, vol. 54, no. 1, pp. 199–218, 2012.

BIBLIOGRAPHY

- [72] G. H. Golub and C. F. Van Loan, *Matrix Computations*. JHU Press, 2012, vol. 3.
- [73] K. Jung and D. Shah, "Local algorithms for approximate inference in minor-excluded graphs," in *Advances in Neural Information Processing Systems*, 2008, pp. 729–736.
- [74] F. Hadlock, "Finding a maximum cut of a planar graph in polynomial time," *SIAM Journal on Computing*, vol. 4, no. 3, pp. 221–225, 1975.
- [75] J. Whittaker, *Graphical models in applied multivariate statistics*. Wiley Publishing, 2009.
- [76] M. Newman, *Networks: an introduction*. Oxford University Press, 2010.
- [77] E. L. Kaplan and P. Meier, "Nonparametric estimation from incomplete observations," *Journal of the American statistical association*, vol. 53, no. 282, pp. 457–481, 1958.
- [78] J. M. Bland and D. G. Altman, "The logrank test," *BMJ*, vol. 328, no. 7447, p. 1073, 2004.
- [79] G. S. Bailey, A. P. Reddy, C. B. Pereira, U. Harttig, W. Baird, J. M. Spitsbergen, J. D. Hendricks, G. A. Orner, D. E. Williams, and J. A. Swenberg, "Nonlinear cancer response at ultralow dose: a 40800-animal ed001 tumor and biomarker study," *Chemical research in toxicology*, vol. 22, no. 7, pp. 1264–1276, 2009.
- [80] X. Huang, W. Qian, I. H. El-Sayed, and M. A. El-Sayed, "The potential use of the enhanced nonlinear properties of gold nanospheres in photothermal cancer therapy," *Lasers in surgery and medicine*, vol. 39, no. 9, pp. 747–753, 2007.
- [81] J. S. Lowengrub, H. B. Frieboes, F. Jin, Y. Chuang, X. Li, P. Macklin, S. Wise, and V. Cristini, "Nonlinear modelling of cancer: bridging the gap between cells and tumours," *Nonlinearity*, vol. 23, no. 1, 2010.

BIBLIOGRAPHY

- [82] L. Nagel and R. Rohrer, "Computer analysis of nonlinear circuits, excluding radiation (cancer)," *IEEE Journal of Solid-State Circuits*, vol. 6, no. 4, pp. 166–182, 1971.
- [83] J. G. Wagner, J. W. Gyves, P. L. Stetson, S. C. Walker-Andrews, I. S. Wollner, M. K. Cochran, and W. D. Ensminger, "Steady-state nonlinear pharmacokinetics of 5-fluorouracil during hepatic arterial and intravenous infusions in cancer patients," *Cancer research*, vol. 46, no. 3, pp. 1499–1506, 1986.
- [84] B. W. Stewart, P. Kleihues, and et al., *World cancer report*. IARC press Lyon, 2003, vol. 57.
- [85] J. K. Pritchard, M. Stephens, and P. Donnelly, "Inference of population structure using multilocus genotype data," *Genetics*, vol. 155, no. 2, pp. 945–959, 2000.
- [86] W. Huber, A. Von Heydebreck, H. Sülthmann, A. Poustka, and M. Vingron, "Variance stabilization applied to microarray data calibration and to the quantification of differential expression," *Bioinformatics*, vol. 18, pp. S96–S104, 2002.
- [87] P. Good, *Permutation tests: a practical guide to resampling methods for testing hypotheses*. Springer Science and Business Media, 2013.
- [88] Y. Benjamini and Y. Hochberg, "Controlling the false discovery rate: a practical and powerful approach to multiple testing," *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 289–300, 1995.
- [89] L. Devroye, "The equivalence of weak, strong and complete convergence in l_1 for kernel density estimates," *The Annals of Statistics*, pp. 896–904, 1983.
- [90] D. Berend and A. Kontorovich, "On the convergence of the empirical distribution," *arXiv preprint:1205.6711*, 2012.
- [91] V. Jeyakumar, A. M. Rubinov, and Z. Wu, "Non-convex quadratic minimization problems with quadratic constraints: global optimality conditions," *Mathematical Programming*, vol. 110, no. 3, pp. 521–541, 2007.

BIBLIOGRAPHY

- [92] E. Segal, M. Shapira, A. Regev, D. Pe'er, D. Botstein, D. Koller, and N. Friedman, "Module networks: identifying regulatory modules and their condition-specific regulators from gene expression data," *Nature Genetics*, vol. 34, no. 2, pp. 166–176, 2003.
- [93] Z. Bar-Joseph, G. K. Gerber, T. I. Lee, N. J. Rinaldi, J. Y. Yoo, F. Robert, D. B. Gordon, E. Fraenkel, T. S. Jaakkola, R. A. Young *et al.*, "Computational discovery of gene modules and regulatory networks," *Nature biotechnology*, vol. 21, no. 11, pp. 1337–1342, 2003.
- [94] D. Marbach, S. Roy, F. Ay, P. E. Meyer, R. Candeias, T. Kahveci, C. A. Bristow, and M. Kellis, "Predictive regulatory models in drosophila melanogaster by integrative inference of transcriptional networks," *Genome research*, vol. 22, no. 7, pp. 1334–1349, 2012.
- [95] R. Sharan and T. Ideker, "Modeling cellular machinery through biological network comparison," *Nature biotechnology*, vol. 24, no. 4, pp. 427–433, 2006.
- [96] S. A. McCarroll, C. T. Murphy, S. Zou, S. D. Pletcher, C.-S. Chin, Y. N. Jan, C. Kenyon, C. I. Bargmann, and H. Li, "Comparing genomic expression patterns across species identifies shared transcriptional profile in aging," *Nature genetics*, vol. 36, no. 2, pp. 197–204, 2004.
- [97] J. O. Woods, U. M. Singh-Blom, J. M. Laurent, K. L. McGary, and E. M. Marcotte, "Prediction of gene-phenotype associations in humans, mice, and plants using phenologs," *BMC bioinformatics*, vol. 14, no. 1, p. 203, 2013.
- [98] P. Kheradpour, A. Stark, S. Roy, and M. Kellis, "Reliable prediction of regulator targets using 12 drosophila genomes," *Genome research*, vol. 17, no. 12, pp. 1919–1931, 2007.
- [99] J. J. Faith, B. Hayete, J. T. Thaden, I. Mogno, J. Wierzbowski, G. Cottarel, S. Kasif, J. J. Collins, and T. S. Gardner, "Large-scale mapping and validation

BIBLIOGRAPHY

- of escherichia coli transcriptional regulation from a compendium of expression profiles,” *PLoS biology*, vol. 5, no. 1, 2007.
- [100] A. Irrthum, L. Wehenkel, P. Geurts *et al.*, “Inferring regulatory networks from expression data using tree-based methods,” *PloS One*, vol. 5, no. 9, 2010.
- [101] S. Roy, J. Ernst, P. V. Kharchenko, P. Kheradpour, N. Negre, M. L. Eaton, J. M. Landolin, C. A. Bristow, L. Ma, M. F. Lin *et al.*, “Identification of functional elements and regulatory circuits by drosophila modencode,” *Science*, vol. 330, no. 6012, pp. 1787–1797, 2010.
- [102] S. Feizi, D. Marbach, M. Médard, and M. Kellis, “Network deconvolution as a general method to distinguish direct dependencies in networks,” *Nature biotechnology*, 2013.
- [103] D. J. Reiss, N. S. Baliga, and R. Bonneau, “Integrated biclustering of heterogeneous genome-wide datasets for the inference of global regulatory networks,” *BMC bioinformatics*, vol. 7, no. 1, p. 280, 2006.
- [104] A. Greenfield, A. Madar, H. Ostrer, and R. Bonneau, “Dream4: Combining genetic and dynamic information to identify biological networks and dynamical models,” *PloS one*, vol. 5, no. 10, 2010.
- [105] M. Mendoza, “Exploring ensemble learning techniques to optimize the reverse engineering of gene regulatory networks,” *Ph.D thesis, Instituto de Inform tica, Universidade Federal do Rio Grande do Sul, Brazil*, 2014.
- [106] D. Marbach, R. J. Prill, T. Schaffter, C. Mattiussi, D. Floreano, and G. Stolovitzky, “Revealing strengths and weaknesses of methods for gene network inference,” *Proceedings of the National Academy of Sciences*, vol. 107, no. 14, pp. 6286–6291, 2010.
- [107] R. Bonneau, D. J. Reiss, P. Shannon, M. T. Facciotti, L. Hood, and *et al.*, “The inferelator: an algorithm for learning parsimonious regulatory networks from systems-biology data sets,” *Genome Biology*, vol. 7, 2006.

BIBLIOGRAPHY

- [108] N. Friedman, M. Linial, I. Nachman, and D. Pe'er, "Using bayesian networks to analyze expression data," *Journal of computational biology*, vol. 7, no. 3-4, pp. 601–620, 2000.
- [109] D. Maglott, J. Ostell, K. D. Pruitt, and T. Tatusova, "Entrez gene: gene-centered information at ncbi," *Nucleic acids research*, vol. 33, pp. D54–D58, 2005.
- [110] E. Wingender, X. Chen, R. Hehl, H. Karas, I. Liebich, V. Matys, T. Meinhardt, M. Prüß, I. Reuter, and F. Schacherer, "Transfac: an integrated system for gene expression regulation," *Nucleic acids research*, vol. 28, no. 1, pp. 316–319, 2000.
- [111] S. M. Gallo, D. T. Gerrard, D. Miner, M. Simich, B. Des Soye, C. M. Bergman, and M. S. Halfon, "Redfly v3. 0: toward a comprehensive database of transcriptional regulatory elements in drosophila," *Nucleic acids research*, vol. 39, pp. D118–D123, 2011.
- [112] M. I. Barrasa, P. Vaglio, F. Cavasino, L. Jacotot, and A. J. Walhout, "Edgedb: a transcription factor-dna interaction database for the analysis of c. elegans differential gene expression," *BMC genomics*, vol. 8, no. 1, p. 21, 2007.
- [113] I. H. Witten and E. Frank, *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, 2005.
- [114] D. Marbach, T. Schaffter, C. Mattiussi, and D. Floreano, "Generating realistic *in silico* gene networks for performance assessment of reverse engineering methods," *Journal of Computational Biology*, vol. 16, no. 2, pp. 229–239, 2009.
- [115] S. D. Dunn, L. M. Wahl, and G. B. Gloor, "Mutual information without the influence of phylogeny or entropy dramatically improves residue contact prediction," *Bioinformatics*, vol. 24, no. 3, pp. 333–340, 2008.
- [116] L. Burger and E. Van Nimwegen, "Disentangling direct from indirect co-evolution of residues in protein alignments," *PLoS Computational Biology*, vol. 6, no. 1, 2010.

BIBLIOGRAPHY

- [117] B. Giraud, J. M. Heumann, and A. S. Lapedes, “Superadditive correlation,” *Physical Review E*, vol. 59, no. 5, p. 4983, 1999.
- [118] A. S. Lapedes, B. G. Giraud, L. Liu, and G. D. Stormo, “Correlated mutations in models of protein sequences: phylogenetic and structural effects,” *Lecture Notes-Monograph Series*, pp. 236–256, 1999.
- [119] N. E. Friedkin, “Information flow through strong and weak ties in intraorganizational social networks,” *Social networks*, vol. 3, no. 4, pp. 273–285, 1982.
- [120] A. De La Fuente, N. Bing, I. Hoeschele, and P. Mendes, “Discovery of meaningful associations in genomic data using partial correlation coefficients,” *Bioinformatics*, vol. 20, no. 18, pp. 3565–3574, 2004.
- [121] C. Hemelrijk, “A matrix partial correlation test used in investigations of reciprocity and other social interaction patterns at group level,” *Journal of theoretical Biology*, vol. 143, no. 3, pp. 405–420, 1990.
- [122] D. Veiga, F. Vicente, M. Grivet, A. De la Fuente, and A. Vasconcelos, “Genome-wide partial correlation analysis of escherichia coli microarray data,” *Genet Mol Res*, vol. 6, no. 4, pp. 730–742, 2007.
- [123] D. S. Marks, L. J. Colwell, R. Sheridan, T. A. Hopf, A. Pagnani, R. Zecchina, and C. Sander, “Protein 3d structure computed from evolutionary sequence variation,” *PloS one*, vol. 6, no. 12, 2011.
- [124] T. A. Hopf, L. J. Colwell, R. Sheridan, B. Rost, C. Sander, and D. S. Marks, “Three-dimensional structures of membrane proteins from genomic sequencing,” *Cell*, vol. 149, no. 7, pp. 1607–1621, 2012.
- [125] M. Weigt, R. A. White, H. Szurmant, J. A. Hoch, and T. Hwa, “Identification of direct residue contacts in protein–protein interaction by message passing,” *Proceedings of the National Academy of Sciences*, vol. 106, no. 1, pp. 67–72, 2009.

BIBLIOGRAPHY

- [126] A. Seth, “Granger causality,” *Scholarpedia*, vol. 2, no. 7, p. 1667, 1667.
- [127] C. J. Quinn, T. P. Coleman, N. Kiyavash, and N. G. Hatsopoulos, “Estimating the directed information to infer causal relationships in ensemble neural spike train recordings,” *Journal of computational neuroscience*, vol. 30, no. 1, pp. 17–44, 2011.
- [128] M. Kamiński, M. Ding, W. A. Truccolo, and S. L. Bressler, “Evaluating causal relations in neural systems: Granger causality, directed transfer function and statistical assessment of significance,” *Biological cybernetics*, vol. 85, no. 2, pp. 145–157, 2001.
- [129] J. Pearl, “Causality: models, reasoning and inference,” *Economet. Theory*, vol. 19, pp. 675–685, 2003.
- [130] N. Meinshausen and P. Bühlmann, “High-dimensional graphs and variable selection with the lasso,” *The annals of statistics*, vol. 34, pp. 1436–1462, 2006.
- [131] A. Pinna, N. Soranzo, and A. De La Fuente, “From knockouts to networks: establishing direct cause-effect relationships through graph analysis,” *PloS one*, vol. 5, no. 10, 2010.
- [132] N. Friedman, M. Linial, I. Nachman, and D. Pe’er, “Using bayesian networks to analyze expression data,” *Journal of computational biology*, vol. 7, no. 3-4, pp. 601–620, 2000.
- [133] N. Friedman, “Inferring cellular networks using probabilistic graphical models,” *Science*, vol. 303, no. 5659, pp. 799–805, 2004.
- [134] A. J. Hartemink, D. K. Gifford, T. S. Jaakkola, R. A. Young *et al.*, “Using graphical models and genomic expression data to statistically validate models of genetic regulatory networks.” in *Pacific symposium on biocomputing*, vol. 6, no. 422-433, 2001.

BIBLIOGRAPHY

- [135] A. A. Margolin, I. Nemenman, K. Basso, C. Wiggins, G. Stolovitzky, R. D. Favera, and A. Califano, "Aracne: an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context," *BMC bioinformatics*, vol. 7, 2006.
- [136] D. S. Marks, T. A. Hopf, and C. Sander, "Protein structure prediction from sequence variation," *Nature biotechnology*, vol. 30, no. 11, pp. 1072–1080, 2012.
- [137] D. T. Jones, D. W. Buchan, D. Cozzetto, and M. Pontil, "Psicov: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments," *Bioinformatics*, vol. 28, no. 2, pp. 184–190, 2012.
- [138] Y. Saad, *Numerical methods for large eigenvalue problems*, 1992, vol. 158.
- [139] R. B. Lehoucq and D. C. Sorensen, "Deflation techniques for an implicitly restarted arnoldi iteration," *SIAM Journal on Matrix Analysis and Applications*, vol. 17, no. 4, pp. 789–821, 1996.
- [140] D. C. Sorensen, "Implicit application of polynomial filters in ak-step arnoldi method," *Siam journal on matrix analysis and applications*, vol. 13, no. 1, pp. 357–385, 1992.
- [141] D. J. Reiss, N. S. Baliga, and R. Bonneau, "Integrated biclustering of heterogeneous genome-wide datasets for the inference of global regulatory networks," *BMC bioinformatics*, vol. 7, no. 1, p. 280, 2006.
- [142] N. Friedman, L. M., I. Nachman, and D. Pe'er, "Using bayesian networks to analyze expression data," *J Comput Biol.*, vol. 7, no. 3-4, pp. 601–620, 2000.
- [143] A. A. Margolin, K. Wang, W. K. Lim, M. Kustagi, I. Nemenman, and A. Califano, "Reverse engineering cellular networks," *Nature Protocols*, vol. 1, pp. 662–671, 2006.
- [144] A. J. Butte and I. S. Kohane, "Mutual information relevance networks: functional genomic clustering using pairwise entropy measurements," in *Pac Symp Biocomput*, vol. 5, 2000, pp. 418–429.

BIBLIOGRAPHY

- [145] R. Küffner, T. Petri, P. Tavakkolkhah, L. Windhager, and R. Zimmer, “Inferring gene regulatory networks by anova,” *Bioinformatics*, vol. 28, no. 10, pp. 1376–1382, 2012.
- [146] V. Huynh-Thu, A. Irrthum, L. Wehenkel, and P. Geurts, “Inferring regulatory networks from expression data using tree-based methods,” *PLoS One*, vol. 5, p. e12776.
- [147] A.-C. Haury, F. Mordelet, P. Vera-Licona, and J.-P. Vert, “Tigress: trustful inference of gene regulation using stability selection,” *BMC systems biology*, vol. 6, no. 1, p. 145, 2012.
- [148] R. J. Prill, D. Marbach, J. Saez-Rodriguez, P. K. Sorger, L. G. Alexopoulos, and *et al*, “Towards a rigorous assessment of systems biology models: the dream3 challenges,” *PloS ONE*, 2010.
- [149] R. Kohavi *et al.*, “A study of cross-validation and bootstrap for accuracy estimation and model selection,” in *Ijcai*, vol. 14, no. 2, 1995, pp. 1137–1145.
- [150] D. Altschuh, A. Lesk, A. Bloomer, and A. Klug, “Correlation of co-ordinated amino acid substitutions with function in viruses related to tobacco mosaic virus,” *Journal of molecular biology*, vol. 193, no. 4, pp. 693–707, 1987.
- [151] U. Göbel, C. Sander, R. Schneider, and A. Valencia, “Correlated mutations and residue contacts in proteins,” *Proteins: Structure, Function, and Bioinformatics*, vol. 18, no. 4, pp. 309–317, 1994.
- [152] E. Neher, “How frequent are correlated changes in families of protein sequences?” *Proceedings of the National Academy of Sciences*, vol. 91, no. 1, pp. 98–102, 1994.
- [153] T. Nugent and D. T. Jones, “Accurate de novo structure prediction of large transmembrane protein domains using fragment-assembly and correlated mutation analysis,” *Proceedings of the National Academy of Sciences*, vol. 109, no. 24, pp. E1540–E1547, 2012.

BIBLIOGRAPHY

- [154] F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D. S. Marks, C. Sander, R. Zecchina, J. N. Onuchic, T. Hwa, and M. Weigt, “Direct-coupling analysis of residue coevolution captures native contacts across many protein families,” *Proceedings of the National Academy of Sciences*, vol. 108, no. 49, pp. E1293–E1301, 2011.
- [155] A. Lapedes, B. Giraud, and C. Jarzynski, “Using sequence alignments to predict protein structure and stability with high accuracy,” *arXiv:1207.2484*, 2012.
- [156] M. Ekeberg, C. Lökvist, Y. Lan, M. Weigt, and E. Aurell, “Improved contact prediction in proteins: using pseudolikelihoods to infer potts models,” *Physical Review E*, vol. 87, no. 1, 2013.
- [157] M. E. Newman, “Finding community structure in networks using the eigenvectors of matrices,” *Physical review E*, vol. 74, no. 3, 2006.
- [158] E. Ferrara, P. De Meo, G. Fiumara, and A. Provetti, “The role of strong and weak ties in facebook: a community structure perspective,” *arXiv:1203.0535*, 2012.
- [159] J. Tang, J. Sun, C. Wang, and Z. Yang, “Social influence analysis in large-scale networks,” in *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2009, pp. 807–816.
- [160] J. A. Bondy and U. S. R. Murty, *Graph theory with applications*. Macmillan London, 1976, vol. 6.
- [161] R. Singh, J. Xu, and B. Berger, “Global alignment of multiple protein interaction networks with application to functional orthology detection,” *Proceedings of the National Academy of Sciences*, vol. 105, no. 35, pp. 12 763–12 768, 2008.
- [162] C.-S. Liao, K. Lu, M. Baym, R. Singh, and B. Berger, “Isorankn: spectral methods for global alignment of multiple protein networks,” *Bioinformatics*, vol. 25, no. 12, pp. i253–i258, 2009.

BIBLIOGRAPHY

- [163] J. Flannick, A. Novak, B. S. Srinivasan, H. H. McAdams, and S. Batzoglou, “Graemlin: general and robust alignment of multiple large interaction networks,” *Genome research*, vol. 16, no. 9, pp. 1169–1181, 2006.
- [164] M. Zaslavskiy, F. Bach, and J.-P. Vert, “Global alignment of protein–protein interaction networks by graph matching methods,” *Bioinformatics*, vol. 25, no. 12, pp. i259–1267, 2009.
- [165] B. P. Kelley, B. Yuan, F. Lewitter, R. Sharan, B. R. Stockwell, and T. Ideker, “Pathblast: a tool for alignment of protein interaction networks,” *Nucleic acids research*, vol. 32, pp. W83–W88, 2004.
- [166] M. Kalaei, M. Smoot, T. Ideker, and R. Sharan, “Networkblast: comparative analysis of protein networks,” *Bioinformatics*, vol. 24, no. 4, pp. 594–596, 2008.
- [167] D. Conte, P. Foggia, C. Sansone, and M. Vento, “Thirty years of graph matching in pattern recognition,” *International journal of pattern recognition and artificial intelligence*, vol. 18, no. 03, pp. 265–298, 2004.
- [168] C. Schellewald and C. Schnörr, “Probabilistic subgraph matching based on convex relaxation,” in *Energy minimization methods in computer vision and pattern recognition*, 2005, pp. 171–186.
- [169] S. Lacoste-Julien, B. Taskar, D. Klein, and M. I. Jordan, “Word alignment via quadratic assignment,” in *Proceedings of the main conference on Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics*, 2006, pp. 112–119.
- [170] S. Melnik, H. Garcia-Molina, and E. Rahm, “Similarity flooding: A versatile graph matching algorithm and its application to schema matching,” in *IEEE International Conference on Data Engineering*, 2002, pp. 117–128.
- [171] A. Narayanan and V. Shmatikov, “De-anonymizing social networks,” in *IEEE Symposium on Security and Privacy*, 2009, pp. 173–187.

- [172] R. E. Burkard, *Quadratic assignment problems*. Springer, 2013.
- [173] K. Makarychev, R. Manokaran, and M. Sviridenko, “Maximum quadratic assignment problem: Reduction from maximum label cover and lp-based approximation algorithm,” in *Automata, Languages and Programming*. Springer, 2010, pp. 594–604.
- [174] M. Bazaraa and O. Kirca, “A branch-and-bound-based heuristic for solving the quadratic assignment problem,” *Naval research logistics quarterly*, vol. 30, no. 2, pp. 287–304, 1983.
- [175] M. S. Bazaraa and H. D. Sherali, “On the use of exact and heuristic cutting plane methods for the quadratic assignment problem,” *Journal of the Operational Research Society*, pp. 991–1003, 1982.
- [176] E. L. Lawler, “The quadratic assignment problem,” *Management science*, vol. 9, no. 4, pp. 586–599, 1963.
- [177] L. Kaufman and F. Broeckx, “An algorithm for the quadratic assignment problem using bender’s decomposition,” *European Journal of Operational Research*, vol. 2, no. 3, pp. 207–211, 1978.
- [178] A. Frieze and J. Yadegar, “On the quadratic assignment problem,” *Discrete applied mathematics*, vol. 5, no. 1, pp. 89–98, 1983.
- [179] W. P. Adams and T. A. Johnson, “Improved linear programming-based lower bounds for the quadratic assignment problem,” *DIMACS series in discrete mathematics and theoretical computer science*, vol. 16, pp. 43–75, 1994.
- [180] G. Finke, R. E. Burkard, and F. Rendl, “Quadratic assignment problems,” *North-Holland Mathematics Studies*, vol. 132, pp. 61–82, 1987.
- [181] S. Hadley, F. Rendl, and H. Wolkowicz, “A new lower bound via projection for the quadratic assignment problem,” *Mathematics of Operations Research*, vol. 17, no. 3, pp. 727–739, 1992.

- [182] K. Anstreicher and H. Wolkowicz, "On lagrangian relaxation of quadratic matrix constraints," *SIAM Journal on Matrix Analysis and Applications*, vol. 22, no. 1, pp. 41–55, 2000.
- [183] K. M. Anstreicher and N. W. Brixius, "Solving quadratic assignment problems using convex quadratic programming relaxations," *Optimization Methods and Software*, vol. 16, no. 1-4, pp. 49–68, 2001.
- [184] Q. Zhao, S. E. Karisch, F. Rendl, and H. Wolkowicz, "Semidefinite programming relaxations for the quadratic assignment problem," *Journal of Combinatorial Optimization*, vol. 2, no. 1, pp. 71–109, 1998.
- [185] J. Peng, H. Mittelman, and X. Li, "A new relaxation framework for quadratic assignment problems based on matrix splitting," *Mathematical Programming Computation*, vol. 2, no. 1, pp. 59–77, 2010.
- [186] J. T. Vogelstein, J. M. Conroy, V. Lyzinski, L. J. Podrazik, S. G. Kratzer, E. T. Harley, D. E. Fishkind, R. J. Vogelstein, and C. E. Priebe, "Fast approximate quadratic programming for large (brain) graph matching," *arXiv:1112.5507*, 2011.
- [187] M. Kolář, J. Meier, V. Mustonen, M. Lässig, and J. Berg, "Graphalignment: Bayesian pairwise alignment of biological networks," *BMC systems biology*, vol. 6, no. 1, p. 144, 2012.
- [188] M. Bayati, D. F. Gleich, A. Saberi, and Y. Wang, "Message-passing algorithms for sparse network alignment," *ACM Transactions on Knowledge Discovery from Data (TKDD)*, vol. 7, no. 1, p. 3, 2013.
- [189] M. Leordeanu and M. Hebert, "A spectral technique for correspondence problems using pairwise constraints," in *IEEE International Conference on Computer Vision*, vol. 2, 2005, pp. 1482–1489.
- [190] M. Carcassoni and E. R. Hancock, "Alignment using spectral clusters." in *BMVC*, 2002, pp. 1–10.

BIBLIOGRAPHY

- [191] E. Kazemi, H. S Hamed, and M. Grossglauser, “Growing a graph matching from a handful of seeds,” in *Proceedings of the Vldb Endowment International Conference on Very Large Data Bases*, vol. 8, 2015.
- [192] C. Clark and J. Kalita, “A multiobjective memetic algorithm for ppi network alignment,” *Bioinformatics*, p. btv063, 2015.
- [193] N. Malod-Dognin and N. Pržulj, “L-graal: Lagrangian graphlet-based network aligner,” *Bioinformatics*, p. btv130, 2015.
- [194] E. M. Loiola, N. M. M. de Abreu, P. O. Boaventura-Netto, P. Hahn, and T. Querido, “A survey for the quadratic assignment problem,” *European Journal of Operational Research*, vol. 176, no. 2, pp. 657–690, 2007.
- [195] D. L. Sussman, M. Tang, D. E. Fishkind, and C. E. Priebe, “A consistent adjacency spectral embedding for stochastic blockmodel graphs,” *Journal of the American Statistical Association*, vol. 107, no. 499, pp. 1119–1128, 2012.
- [196] T. Qin and K. Rohe, “Regularized spectral clustering under the degree-corrected stochastic blockmodel,” in *Advances in Neural Information Processing Systems*, 2013, pp. 3120–3128.
- [197] A. Athreya, C. Priebe, M. Tang, V. Lyzinski, D. Marchette, and D. Sussman, “A limit theorem for scaled eigenvectors of random dot product graphs,” *Sankhya A*, pp. 1–18, 2013.
- [198] A. Saade, F. Krzakala, and L. Zdeborová, “Spectral clustering of graphs with the bethe hessian,” in *Advances in Neural Information Processing Systems*, 2014, pp. 406–414.
- [199] P. Erdős and A. Rényi, “On the strength of connectedness of a random graph,” *Acta Mathematica Hungarica*, vol. 12, no. 1, pp. 261–267, 1961.
- [200] T. Czajka and G. Pandurangan, “Improved random graph isomorphism,” *Journal of Discrete Algorithms*, vol. 6, no. 1, pp. 85–92, 2008.

BIBLIOGRAPHY

- [201] L. Babai, “Graph isomorphism in quasipolynomial time [extended abstract],” in *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing*, 2016, pp. 684–697.
- [202] W. Ali and C. M. Deane, “Functionally guided alignment of protein interaction networks for module detection,” *Bioinformatics*, vol. 25, no. 23, pp. 3166–3173, 2009.
- [203] M. E. Newman, “Communities, modules and large-scale structure in networks,” *Nature Physics*, vol. 8, no. 1, pp. 25–31, 2012.
- [204] G. W. Klau, “A new graph-based method for pairwise global network alignment,” *BMC bioinformatics*, vol. 10, 2009.
- [205] P. Schweitzer, “Problems of unknown complexity: graph isomorphism and ramsey theoretic numbers,” Ph.D. dissertation, Saarbrücken, Univ., Diss., 2009, 2009.
- [206] L. Babai and L. Kucera, “Canonical labelling of graphs in linear average time,” in *IEEE Symposium on Foundations of Computer Science*, 1979, pp. 39–46.
- [207] L. Babai, P. Erdős, and S. M. Selkow, “Random graph isomorphism,” *SIAM Journal on Computing*, vol. 9, no. 3, pp. 628–635, 1980.
- [208] D. B. West, *Introduction to graph theory*. Prentice hall Upper Saddle River, 2001, vol. 2.
- [209] J. Kuczynski and H. Wozniakowski, “Estimating the largest eigenvalue by the power and lanczos algorithms with a random start,” *SIAM journal on matrix analysis and applications*, vol. 13, no. 4, pp. 1094–1122, 1992.
- [210] P.-A. Wedin, “Perturbation bounds in connection with singular value decomposition,” *Informationsbehandlung (BIT)*, vol. 12, pp. 99–111, 1972.
- [211] E. W. Weisstein, “Gershgorin circle theorem,” 2003.

BIBLIOGRAPHY

- [212] S. Umeyama, “An eigendecomposition approach to weighted graph matching problems,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 10, no. 5, pp. 695–703, 1988.
- [213] S. Tu and J. Wang, “Practical first order methods for large scale semidefinite programming,” 2014.
- [214] S. Burer and R. D. Monteiro, “A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization,” *Mathematical Programming*, vol. 95, no. 2, pp. 329–357, 2003.
- [215] Z. Wen, D. Goldfarb, S. Ma, and K. Scheinberg, “Row by row methods for semidefinite programming,” *Industrial Engineering*, pp. 1–21, 2009.
- [216] C. Helmberg and F. Rendl, “A spectral bundle method for semidefinite programming,” *SIAM Journal on Optimization*, vol. 10, no. 3, pp. 673–696, 2000.
- [217] J. Renegar, “Efficient first-order methods for linear programming and semidefinite programming,” *arXiv:1409.5832*, 2014.
- [218] G. Evanno, S. Regnaut, and J. Goudet, “Detecting the number of clusters of individuals using the software structure: a simulation study,” *Molecular ecology*, vol. 14, no. 8, pp. 2611–2620, 2005.
- [219] G. W. Milligan and M. C. Cooper, “An examination of procedures for determining the number of clusters in a data set,” *Psychometrika*, vol. 50, no. 2, pp. 159–179, 1985.
- [220] A. Bobbio and K. S. Trivedi, “An aggregation technique for the transient analysis of stiff markov chains,” *IEEE Transactions on Computers*, vol. 100, no. 9, pp. 803–814, 1986.
- [221] W. Aiello, F. Chung, and L. Lu, “A random graph model for power law graphs,” *Experimental Mathematics*, vol. 10, no. 1, pp. 53–66, 2001.

BIBLIOGRAPHY

- [222] R. Kohavi *et al.*, “A study of cross-validation and bootstrap for accuracy estimation and model selection,” vol. 14, no. 2, 1995, pp. 1137–1145.
- [223] Y.-C. Wu, M. S. Bansal, M. D. Rasmussen, J. Herrero, and M. Kellis, “Phylogenetic identification and functional validation of orthologous genes across human, mouse, fly, worm, yeast,” *submitted to Genome Research, available on arXiv*, 2014.
- [224] P. D. Thomas, V. Wood, C. J. Mungall, S. E. Lewis, J. A. Blake, G. O. Consortium *et al.*, “On the use of gene ontology annotations to assess functional similarity among orthologs and paralogs: a short report,” *PLoS computational biology*, vol. 8, no. 2, 2012.
- [225] L. d. F. Costa, F. A. Rodrigues, G. Travieso, and P. R. Villas Boas, “Characterization of complex networks: A survey of measurements,” *Advances in Physics*, vol. 56, no. 1, pp. 167–242, 2007.
- [226] L. C. Freeman, “A set of measures of centrality based on betweenness,” *Sociometry*, pp. 35–41, 1977.
- [227] M. P. Joy, A. Brock, D. E. Ingber, and S. Huang, “High-betweenness proteins in the yeast protein interaction network,” *BioMed Research International*, no. 2, pp. 96–103, 2005.
- [228] B. Zhou, J. Pei, and W. Luk, “A brief survey on anonymization techniques for privacy preserving publishing of social network data,” *ACM SIGKDD Explorations Newsletter*, vol. 10, no. 2, pp. 12–22, 2008.
- [229] M. F. Schwartz, A. R. Brecher, J. Whyte, and M. G. Klein, “A patient registry for cognitive rehabilitation research: a strategy for balancing patients’ privacy rights with researchers’ need for access,” *Archives of physical medicine and rehabilitation*, vol. 86, no. 9, pp. 1807–1814, 2005.
- [230] E. Deutsch and M. Neumann, “On the first and second order derivatives of the perron vector,” *Linear algebra and its applications*, vol. 71, pp. 57–76, 1985.

BIBLIOGRAPHY

- [231] R. Pastor-Satorras and A. Vespignani, "Epidemic spreading in scale-free networks," *Physical review letters*, vol. 86, no. 14, p. 3200, 2001.
- [232] D. Acemoglu, A. Ozdaglar, and A. ParandehGheibi, "Spread of (mis) information in social networks," *Games and Economic Behavior*, vol. 70, no. 2, pp. 194–227, 2010.
- [233] D. Hirshleifer and S. Hong Teoh, "Herd behaviour and cascading in capital markets: A review and synthesis," *European Financial Management*, vol. 9, no. 1, pp. 25–66, 2003.
- [234] M. T. Maurano, R. Humbert, E. Rynes, R. E. Thurman, E. Haugen, H. Wang, A. P. Reynolds, R. Sandstrom, H. Qu, J. Brody *et al.*, "Systematic localization of common disease-associated variation in regulatory dna," *Science*, vol. 337, no. 6099, pp. 1190–1195, 2012.
- [235] A. B. Glinskii, J. Ma, S. Ma, D. Grant, C.-U. Lim, S. Sell, and G. V. Glinsky, "Identification of intergenic trans-regulatory rnas containing a disease-linked snp sequence and targeting cell cycle progression/differentiation pathways in multiple common human disorders," *Cell Cycle*, vol. 8, no. 23, pp. 3925–3942, 2009.
- [236] T. W. Valente, "Network models of the diffusion of innovations," *Computational and Mathematical Organization Theory*, vol. 2, no. 2, pp. 163–164, 1996.
- [237] D. Achlioptas, R. M. D'Souza, and J. Spencer, "Explosive percolation in random networks," *Science*, vol. 323, no. 5920, pp. 1453–1455, 2009.
- [238] N. T. Bailey *et al.*, *The mathematical theory of infectious diseases and its applications*. Charles Griffin., 1975.
- [239] M. E. Newman, "Spread of epidemic disease on networks," *Physical review E*, vol. 66, no. 1, p. 016128, 2002.

BIBLIOGRAPHY

- [240] C. Moore and M. E. Newman, “Epidemics and percolation in small-world networks,” *Physical Review E*, vol. 61, no. 5, p. 5678, 2000.
- [241] A. Ganesh, L. Massoulié, and D. Towsley, “The effect of network topology on the spread of epidemics,” in *24th Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM)*, vol. 2, 2005, pp. 1455–1466.
- [242] N. Demiris and P. D. O’neill, “Bayesian inference for epidemics with two levels of mixing,” *Scandinavian journal of statistics*, vol. 32, no. 2, pp. 265–280, 2005.
- [243] P. D. O’Neill, “A tutorial introduction to bayesian inference for stochastic epidemic models using markov chain monte carlo methods,” *Mathematical Biosciences*, vol. 180, no. 1, pp. 103–114, 2002.
- [244] H. Okamura, K. Tateishi, and T. Dohi, “Statistical inference of computer virus propagation using non-homogeneous poisson processes,” in *IEEE International Symposium on Software Reliability*, 2007, pp. 149–158.
- [245] S. Mostafavi, A. Goldenberg, and Q. Morris, “Labeling nodes using three degrees of propagation,” *PloS one*, vol. 7, no. 12, p. e51947, 2012.
- [246] Y. Bengio, O. Delalleau, and N. Le Roux, “Label propagation and quadratic criterion,” *Semi-supervised learning*, pp. 193–216, 2006.
- [247] M. E. Newman, “Mixing patterns in networks,” *Physical Review E*, vol. 67, no. 2, p. 026126, 2003.
- [248] —, “Scientific collaboration networks. ii. shortest paths, weighted networks, and centrality,” *Physical review E*, vol. 64, no. 1, p. 016132, 2001.
- [249] E. Sefer and C. Kingsford, “Diffusion archaeology for diffusion progression history reconstruction,” *IEEE International Conference on Data Mining*, pp. 530–539, 2014.

BIBLIOGRAPHY

- [250] M. Farajtabar, M. Gomez-Rodriguez, N. Du, M. Zamani, H. Zha, and L. Song, "Back to the past: Source identification in diffusion networks from partially observed cascades," *arXiv preprint arXiv:1501.06582*, 2015.
- [251] B. A. Prakash, J. Vreeken, and C. Faloutsos, "Spotting culprits in epidemics: How many and which ones?" in *International Conference on Data Mining*, 2012, pp. 11–20.
- [252] D. Shah and T. Zaman, "Rumors in a network: Who's the culprit?" *IEEE Transactions on Information Theory*, vol. 57, no. 8, pp. 5163–5181, 2011.
- [253] T. Lappas, E. Terzi, D. Gunopulos, and H. Mannila, "Finding effectors in social networks," in *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2010, pp. 1059–1068.
- [254] N. Du, L. Song, H. Woo, and H. Zha, "Uncover topic-sensitive information diffusion networks," in *Proceedings of the sixteenth international conference on artificial intelligence and statistics*, 2013, pp. 229–237.
- [255] S. Goel, D. J. Watts, and D. G. Goldstein, "The structure of online diffusion networks," in *Proceedings of the 13th ACM Conference on Electronic Commerce*, 2012, pp. 623–638.
- [256] K. Lerman and R. Ghosh, "Information contagion: An empirical study of the spread of news on digg and twitter social networks." *ICWSM*, vol. 10, pp. 90–97, 2010.
- [257] J. Hershberger, M. Maxel, and S. Suri, "Finding the k shortest simple paths: A new algorithm and its implementation," *ACM Transactions on Algorithms*, vol. 3, no. 4, p. 45, 2007.
- [258] C. De Bacco, S. Franz, D. Saad, and C. H. Yeung, "Shortest node-disjoint paths on random graphs," *arXiv:1401.8096*, 2014.

BIBLIOGRAPHY

- [259] D. P. Bertsekas and J. N. Tsitsiklis, *Introduction to Probability*. Athena Scientific, 2002.
- [260] M. Girvan and M. E. Newman, “Community structure in social and biological networks,” *Proceedings of the National Academy of Sciences*, vol. 99, no. 12, pp. 7821–7826, 2002.
- [261] S. Butenko, *Clustering challenges in biological networks*. World Scientific, 2009.
- [262] N. Mishra, R. Schreiber, I. Stanton, and R. E. Tarjan, “Clustering social networks,” in *Algorithms and Models for the Web-Graph*, 2007, pp. 56–67.
- [263] C.-H. Lee, M. N. Hoehn-Weiss, and S. Karim, “Grouping interdependent tasks: Using spectral graph partitioning to study system modularity and performance,” *Available at SSRN*, 2014.
- [264] S. E. Schaeffer, “On the np-completeness of some graph cluster measures,” *SOFSEM 2006: Theory and Practice of Computer Science*, pp. 530–537, 2006.
- [265] A. Decelle, F. Krzakala, C. Moore, and L. Zdeborová, “Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications,” *Physical Review E*, vol. 84, no. 6, 2011.
- [266] E. M. Airoldi, D. M. Blei, S. E. Fienberg, and E. P. Xing, “Mixed membership stochastic blockmodels,” in *Advances in Neural Information Processing Systems*, 2009, pp. 33–40.
- [267] T. A. Snijders and K. Nowicki, “Estimation and prediction for stochastic blockmodels for graphs with latent block structure,” *Journal of classification*, vol. 14, no. 1, pp. 75–100, 1997.
- [268] B. Hajek, Y. Wu, and J. Xu, “Achieving exact cluster recovery threshold via semidefinite programming: Extensions,” *arXiv:1502.07738*, 2015.
- [269] A. A. Amini and E. Levina, “On semidefinite relaxations for the block model,” *arXiv:1406.5647*, 2014.

BIBLIOGRAPHY

- [270] A. Y. Ng, M. I. Jordan, Y. Weiss *et al.*, “On spectral clustering: Analysis and an algorithm,” *Advances in Neural Information Processing Systems*, vol. 2, pp. 849–856, 2002.
- [271] U. Von Luxburg, “A tutorial on spectral clustering,” *Statistics and Computing*, vol. 17, no. 4, pp. 395–416, 2007.
- [272] B. Mohar, Y. Alavi, G. Chartrand, and O. Oellermann, “The laplacian spectrum of graphs,” *Graph theory, combinatorics, and applications*, vol. 2, no. 871-898, p. 12, 1991.
- [273] F. R. Chung, *Spectral graph theory*. American Mathematical Soc., 1997, vol. 92.
- [274] J. Shi and J. Malik, “Normalized cuts and image segmentation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 8, pp. 888–905, 2000.
- [275] P. W. Holland, K. B. Laskey, and S. Leinhardt, “Stochastic blockmodels: First steps,” *Social networks*, vol. 5, no. 2, pp. 109–137, 1983.
- [276] S. J. Young and E. R. Scheinerman, “Random dot product graph models for social networks,” in *Algorithms and models for the web-graph*, 2007, pp. 138–149.
- [277] P. D. Hoff, A. E. Raftery, and M. S. Handcock, “Latent space approaches to social network analysis,” *Journal of the american Statistical association*, vol. 97, no. 460, pp. 1090–1098, 2002.
- [278] L. Hagen and A. B. Kahng, “New spectral methods for ratio cut partitioning and clustering,” *IEEE Transactions on Computer-aided Design of Integrated Circuits and Systems*, vol. 11, no. 9, pp. 1074–1085, 1992.
- [279] E. Mossel, J. Neeman, and A. Sly, “Stochastic block models and reconstruction,” *arXiv:1202.1499*, 2012.
- [280] ———, “A proof of the block model threshold conjecture,” *arXiv:1311.4115*, 2013.

BIBLIOGRAPHY

- [281] L. Massoulié, “Community detection thresholds and the weak ramanujan property,” in *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*, 2014, pp. 694–703.
- [282] F. Krzakala, C. Moore, E. Mossel, J. Neeman, A. Sly, L. Zdeborová, and P. Zhang, “Spectral redemption in clustering sparse networks,” *Proceedings of the National Academy of Sciences*, vol. 110, no. 52, pp. 20 935–20 940, 2013.
- [283] E. Abbe, A. S. Bandeira, and G. Hall, “Exact recovery in the stochastic block model,” *IEEE Transactions on Information Theory*, vol. 62, no. 1, pp. 471–487, 2016.
- [284] S. Shortreed, M. S. Handcock, and P. Hoff, “Positional estimation within a latent space model for networks,” *Methodology*, vol. 2, no. 1, pp. 24–33, 2006.
- [285] M. S. Handcock, A. E. Raftery, and J. M. Tantrum, “Model-based clustering for social networks,” *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, vol. 170, no. 2, pp. 301–354, 2007.
- [286] M. Salter-Townshend and T. B. Murphy, “Variational bayesian inference for the latent position cluster model for network data,” *Computational Statistics and Data Analysis*, vol. 57, no. 1, pp. 661–671, 2013.
- [287] N. Friel, C. Ryan, and J. Wyse, “Bayesian model selection for the latent position cluster model for social networks,” *arXiv:1308.4871*, 2013.
- [288] K. Rohe, S. Chatterjee, and B. Yu, “Spectral clustering and the high-dimensional stochastic blockmodel,” *The Annals of Statistics*, pp. 1878–1915, 2011.
- [289] L. Zelnik-Manor and P. Perona, “Self-tuning spectral clustering,” in *Advances in neural information processing systems*, 2004, pp. 1601–1608.
- [290] F. Chen, A. P. Chandrakasan, and V. Stojanović, “A signal-agnostic compressed sensing acquisition system for wireless and implantable sensors,” in *IEEE Custom Integrated Circuits Conference*, 2010, pp. 1–4.

BIBLIOGRAPHY

- [291] P. Dao, R. Colak, R. Salari, F. Moser, E. Davicioni, A. Schönhuth, and M. Ester, “Inferring cancer subnetwork markers using density-constrained biclustering,” *Bioinformatics*, vol. 26, no. 18, pp. i625–i631, 2010.
- [292] G. Bisson and F. Hussain, “Chi-sim: A new similarity measure for the co-clustering task,” in *International Conference on Machine Learning and Applications*, 2008, pp. 211–217.
- [293] S. Bergmann, J. Ihmels, and N. Barkai, “Iterative signature algorithm for the analysis of large-scale gene expression data,” *Physical review E*, vol. 67, no. 3, p. 031902, 2003.
- [294] A. A. Shabalín, V. J. Weigman, C. M. Perou, and A. B. Nobel, “Finding large average submatrices in high dimensional data,” *The Annals of Applied Statistics*, pp. 985–1012, 2009.
- [295] A. Prelić, S. Bleuler, P. Zimmermann, A. Wille, P. Bühlmann, W. Gruissem, L. Hennig, L. Thiele, and E. Zitzler, “A systematic comparison and evaluation of biclustering methods for gene expression data,” *Bioinformatics*, vol. 22, no. 9, pp. 1122–1129, 2006.
- [296] A. Tanay, R. Sharan, and R. Shamir, “Discovering statistically significant biclusters in gene expression data,” *Bioinformatics*, vol. 18, pp. S136–S144, 2002.
- [297] S. Hochreiter, U. Bodenhofer, M. Heusel, A. Mayr, A. Mitterecker, A. Kasim, T. Khamiakova, S. Van Sanden, D. Lin, W. Talloen *et al.*, “Fabia: factor analysis for bicluster acquisition,” *Bioinformatics*, vol. 26, no. 12, pp. 1520–1527, 2010.
- [298] L. Li, Y. Guo, W. Wu, Y. Shi, J. Cheng, and S. Tao, “A comparison and evaluation of five biclustering algorithms by quantifying goodness of biclusters for gene expression data,” *BioData mining*, vol. 5, no. 1, p. 1, 2012.
- [299] K. Eren, M. Deveci, O. Küçüktunç, and Ü. V. Çatalyürek, “A comparative analysis of biclustering algorithms for gene expression data,” *Briefings in bioinformatics*, vol. 14, no. 3, pp. 279–292, 2013.

BIBLIOGRAPHY

- [300] R. R. Nadakuditi and M. E. Newman, “Graph spectra and the detectability of community structure in networks,” *Physical review letters*, vol. 108, no. 18, p. 188701, 2012.
- [301] D. Dueck, B. J. Frey, N. Jojic, V. Jojic, G. Giaever, A. Emili, G. Musso, and R. Hegele, “Constructing treatment portfolios using affinity propagation,” in *Research in Computational Molecular Biology*, 2008, pp. 360–371.
- [302] G. Govaert and M. Nadif, “Block clustering with bernoulli mixture models: Comparison of different approaches,” *Computational Statistics and Data Analysis*, vol. 52, no. 6, pp. 3233–3245, 2008.
- [303] J. S. Yedidia, W. T. Freeman, and Y. Weiss, “Constructing free-energy approximations and generalized belief propagation algorithms,” *IEEE Transactions on Information Theory*, vol. 51, no. 7, pp. 2282–2312, 2005.
- [304] A. P. Dempster, N. M. Laird, and D. B. Rubin, “Maximum likelihood from incomplete data via the em algorithm,” *Journal of the royal statistical society. Series B (methodological)*, pp. 1–38, 1977.
- [305] E. Candès, J. Romberg, and T. Tao, “Stable signal recovery from incomplete and inaccurate measurements,” *Communications on Pure and Applied Mathematics*, vol. 59, no. 8, pp. 1207–1223, 2006.
- [306] E. Candès and J. Romberg, “Sparsity and incoherence in compressive sampling,” *Inverse Problems*, vol. 23, p. 969, 2007.
- [307] D. L. Donoho, “Compressed sensing,” *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [308] M. Akçakaya and V. Tarokh, “Shannon-theoretic limits on noisy compressive sampling,” *IEEE Transactions on Information Theory*, vol. 56, no. 1, pp. 492–504, 2010.

BIBLIOGRAPHY

- [309] S. Sarvotham, D. Baron, and R. Baraniuk, “Measurements vs. bits: Compressed sensing meets information theory,” in *Proceedings of 44th Allerton Conf. Comm., Ctrl., Computing*, 2006.
- [310] Y. Jin, Y.-H. Kim, and B. D. Rao, “Performance tradeoffs for exact support recovery of sparse signals,” in *IEEE International Symposium on Information Theory, ISIT*, 2010, pp. 1558–1562.
- [311] E. Candes and T. Tao, “Decoding by linear programming,” *IEEE Transactions on Information Theory*, vol. 51, no. 12, pp. 4203–4215, 2005.
- [312] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin, “A simple proof of the restricted isometry property for random matrices,” *Constructive Approximation*, vol. 28, no. 3, pp. 253–263, 2008.
- [313] D. Achlioptas, “Database-friendly random projections: Johnson-Lindenstrauss with binary coins,” *Journal of Computer and System Sciences*, vol. 66, no. 4, pp. 671–687, 2003.
- [314] D. Slepian and J. K. Wolf, “Noiseless coding of correlated information sources,” *IEEE Transactions on Information Theory*, vol. 19, no. 4, pp. 471–480, 1973.
- [315] T. Ho, M. Médard, R. Koetter, D. R. Karger, M. Effros, J. Shi, and B. Leong, “A random linear network coding approach to multicast,” *IEEE Transactions on Information Theory*, vol. 52, no. 10, pp. 4413–4430, 2006.
- [316] I. Csiszar and J. Körner, in *Information Theory: Coding Theorems for Discrete Memoryless Systems*. Academic Press, New York, 1981.
- [317] S. Pradhan and K. Ramchandran, “Distributed source coding using syndromes (DISCUS): Design and construction,” in *DCC*, 1999, p. 158.
- [318] T. Coleman, A. Lee, M. Médard, and M. Effros, “Low-complexity approaches to Slepian–Wolf near-lossless distributed data compression,” *IEEE Transactions on Information Theory*, vol. 52, no. 8, 2006.

BIBLIOGRAPHY

- [319] R. Gallager, “Low-density parity-check codes,” *IEEE Transactions on Information Theory*, vol. 8, no. 1, pp. 21–28, 1962.
- [320] T. Tian, J. Garcia-Frias, and W. Zhong, “Compression of correlated sources using LDPC codes,” in *Data Compression Conference, DCC*, 2003.
- [321] C. Berrou, A. Glavieux, P. Thitimajshima *et al.*, “Near Shannon limit error-correcting coding and decoding,” in *Proceedings of ICC*, vol. 93, 1993, pp. 1064–1070.
- [322] J. Garcia-Frias, “Compression of correlated binary sources using turbo codes,” *IEEE Communications Letters*, vol. 5, no. 10, pp. 417–419, 2001.
- [323] A. Ramamoorthy, K. Jain, P. Chou, and M. Effros, “Separating distributed source coding from network coding,” *IEEE Transactions on Information Theory*, vol. 52, no. 6, pp. 2785–2795, 2006.
- [324] G. Maierbacher, J. Barros, and M. Médard, “Practical source-network decoding,” in *6th International Symposium on Wireless Communication Systems, ISWCS*, 2009, pp. 283–287.
- [325] R. Ahlswede, N. Cai, S.-Y. R. Li, and R. W. Yeung, “Network information flow,” *IEEE Transactions on Information Theory*, vol. 46, pp. 1204–1216, 2000.
- [326] G. Forney, *Concatenated codes*. Citeseer, 1966.
- [327] S. Feizi, M. Médard, and M. Effros, “Compressive sensing over networks,” in *Communication, Control, and Computing (Allerton), 2010 48th Annual Allerton Conference on*, 2010, pp. 1129–1136.
- [328] I. Maric, A. Goldsmith, and M. Médard, “Analog network coding in the high-snr regime,” in *IEEE Wireless Network Coding Workshop*, 2010.
- [329] D. Donoho, I. Johnstone, A. Maleki, and A. Montanari, “Compressed sensing over l_p -balls: Minimax mean square error,” *arXiv:1103.1943v2*, 2011.

BIBLIOGRAPHY

- [330] G. Raskutti, M. Wainwright, and B. Yu, "Restricted eigenvalue properties for correlated gaussian designs," *The Journal of Machine Learning Research*, vol. 11, pp. 2241–2259, 2010.
- [331] R. Tibshirani, "Regression shrinkage and selection via the lasso," *J. Royal. Statist. Soc. B*, vol. 58, no. 1, pp. 267–288, 1996.
- [332] A. Cohen, W. Dahmen, and R. DeVore, "Compressed sensing and best k-term approximation," *American Mathematical Society*, vol. 22, no. 1, pp. 211–231, 2009.
- [333] F. Chen, A. Chandrakasan, and V. Stojanovic, "A signal-agnostic compressed sensing acquisition system for wireless and implantable sensors," in *IEEE Custom Integrated Circuits Conference*. IEEE, pp. 1–4.
- [334] S. Katti, S. Gollakota, and D. Katabi, "Embracing wireless interference: Analog network coding," in *Proceedings of the conference on Applications, technologies, architectures, and protocols for computer communications*, 2007, pp. 397–408.
- [335] R. H. Walden, "Analog-to-digital converter survey and analysis," *IEEE Journal on Selected Areas in Communications*, vol. 17, no. 4, pp. 539–550, 1999.
- [336] F. Marvasti, *Nonuniform sampling: theory and practice*. Springer Science and Business Media, 2012.
- [337] S. M. Qaisar, L. Fesquet, and M. Renaudin, "Computationally efficient adaptive rate sampling and filtering," in *15th European Signal Processing Conference*, 2007, pp. 2139–2143.
- [338] D. Wei and A. V. Oppenheim, "Sampling based on local bandwidth," in *Asilomar Conference on Signals, Systems and Computers, 2007. ACSSC 2007*, 2007, pp. 1103–1107.
- [339] T. A. Claasen and W. F. Mecklenbräuker, "On stationary linear time-varying systems," *IEEE Transactions on Circuits and Systems*, vol. 29, no. 3, pp. 169–184, 1982.

BIBLIOGRAPHY

- [340] K. Horiuchi, "Sampling principle for continuous signals with time-varying bands," *Information and Control*, vol. 13, no. 1, pp. 53–61, 1968.
- [341] M. Unser and J. Zerubia, "A generalized sampling theory without band-limiting constraints," *IEEE Transactions on Circuits and Systems II: Analog and Digital Signal Processing*, vol. 45, no. 8, pp. 959–969, 1998.
- [342] P. Vaidyanathan, "Sampling theorems for nonbandlimited signals," in *Sampling, wavelets, and tomography*, 2004, pp. 115–136.
- [343] J. D. Gibson, "Sequentially adaptive backward prediction in adpcm speech coders," *IEEE Transactions on Communications*, vol. 26, no. 1, pp. 145–150, 1978.
- [344] N. S. Jayant, "Adaptive quantization with a one-word memory," *The Journal of the Acoustical Society of America*, vol. 54, no. 1, pp. 340–340, 1973.
- [345] J. D. Gibson, "Adaptive prediction in speech differential encoding systems," *Proceedings of the IEEE*, vol. 68, no. 4, pp. 488–525, 1980.
- [346] A. Ortega and M. Vetterli, "Adaptive scalar quantization without side information," *IEEE Transactions on Image Processing*, vol. 6, no. 5, pp. 665–676, 1997.
- [347] V. K. Goyal, J. Zhuang, and M. Vetterli, "Transform coding with backwards adaptive updates," *IEEE Transactions on Information Theory*, vol. 46, pp. 1623–1633, 2000.
- [348] T. Lai, "Sequential analysis: some classical problems and new challenges," in *Statistica Sinica*, vol. 11, no. 2, 2001, pp. 302–350.
- [349] R. M. Castro and R. D. Nowak, "Upper and lower error bounds for active learning," in *The 44th Annual Allerton Conference on Communication, Control and Computing*, vol. 2, no. 2.1, 2006, p. 1.
- [350] S. S. Haykin, *Introduction to adaptive filters*. Macmillan, 1984.

BIBLIOGRAPHY

- [351] D. P. Bertsekas, D. P. Bertsekas, D. P. Bertsekas, and D. P. Bertsekas, *Dynamic programming and optimal control*, 1995, vol. 1, no. 2.
- [352] H. Nyquist, "Certain topics in telegraph transmission theory," 1928.
- [353] A. M. Dixon, E. G. Allstot, D. Gangopadhyay, and D. J. Allstot, "Compressed sensing system considerations for ecg and emg wireless biosensors," *IEEE Transactions on Biomedical Circuits and Systems*, vol. 6, no. 2, pp. 156–166, 2012.
- [354] S. Feizi, V. K. Goyal, and M. Médard, "Time-stampless adaptive nonuniform sampling for stochastic signals," *IEEE Transactions on Signal Processing*, vol. 60, no. 10, pp. 5440–5450, 2012.
- [355] S. Feizi, V. Goyal, and M. Médard, "Locally adaptive sampling," in *Allerton Conference on Communication, Control, and Computing*, 2010, pp. 152–159.
- [356] A. L. Goldberger, L. A. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, and H. E. Stanley, "Physiobank, physiotoolkit, and physionet components of a new research resource for complex physiologic signals," *Circulation*, vol. 101, no. 23, pp. e215–e220, 2000.
- [357] R. Benzid, A. Messaoudi, and A. Boussaad, "Constrained ecg compression algorithm using the block-based discrete cosine transform," *Digital Signal Processing*, vol. 18, no. 1, pp. 56–64, 2008.
- [358] B. Yu, L. Yang, and C.-C. Chong, "Ecg monitoring over bluetooth: data compression and transmission," in *IEEE Wireless Communications and Networking Conference*, 2010, pp. 1–5.
- [359] K. Kozmin, J. Johansson, and J. Delsing, "Level-crossing adc performance evaluation toward ultrasound application," *Circuits and Systems I: Regular Papers, IEEE Transactions on*, vol. 56, no. 8, pp. 1708–1719, 2009.
- [360] E. J. Candes and T. Tao, "Near-optimal signal recovery from random projections: Universal encoding strategies?" *Information Theory, IEEE Transactions on*, vol. 52, no. 12, pp. 5406–5425, 2006.

BIBLIOGRAPHY

- [361] S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM review*, vol. 43, no. 1, pp. 129–159, 2001.
- [362] E. Candes and T. Tao, "The dantzig selector: statistical estimation when p is much larger than n ," *The Annals of Statistics*, pp. 2313–2351, 2007.
- [363] J. A. Tropp and A. C. Gilbert, "Signal recovery from random measurements via orthogonal matching pursuit," *IEEE Transactions on Information Theory*, vol. 53, no. 12, pp. 4655–4666, 2007.
- [364] D. Needell and J. A. Tropp, "Cosamp: Iterative signal recovery from incomplete and inaccurate samples," *Applied and Computational Harmonic Analysis*, vol. 26, no. 3, pp. 301–321, 2009.
- [365] D. L. Donoho, A. Maleki, and A. Montanari, "Message-passing algorithms for compressed sensing," *Proceedings of the National Academy of Sciences*, vol. 106, no. 45, pp. 18 914–18 919, 2009.
- [366] S. Rangan, V. Goyal, and A. K. Fletcher, "Asymptotic analysis of map estimation via the replica method and compressed sensing," in *Advances in Neural Information Processing Systems*, 2009, pp. 1545–1553.
- [367] M. Grant, S. Boyd, and Y. Ye, "Cvx: Matlab software for disciplined convex programming," 2008.
- [368] R. Grimaldi, S. Rodriguez, and A. Rusu, "A 10-bit 5khz level-crossing adc," in *20th European Conference on Circuit Theory and Design*, 2011, pp. 564–567.
- [369] F. Chen, A. P. Chandrakasan, and V. M. Stojanović, "Design and analysis of a hardware-efficient compressed sensing architecture for data compression in wireless sensors," *IEEE Journal of Solid-State Circuits*, vol. 47, no. 3, pp. 744–756, 2012.
- [370] M. Mishali and Y. C. Eldar, "From theory to practice: Sub-nyquist sampling of sparse wideband analog signals," *IEEE Journal of Selected Topics in Signal Processing*, vol. 4, no. 2, pp. 375–391, 2010.

BIBLIOGRAPHY

- [371] S. Kirolos, J. Laska, M. Wakin, M. Duarte, D. Baron, T. Ragheb, Y. Masoud, and R. Baraniuk, "Analog-to-information conversion via random demodulation," in *IEEE Dallas/CAS Workshop on Design, Applications, Integration and Software*, 2006, pp. 71–74.
- [372] O. Abari, F. Chen, F. Lim, and V. Stojanović, "Performance trade-offs and design limitations of analog-to-information converter front-ends," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2012, pp. 5309–5312.
- [373] P. P. Mercier, D. C. Daly, and A. P. Chandrakasan, "An energy-efficient all-digital uwb transmitter employing dual capacitively-coupled pulse-shaping drivers," *IEEE Journal of Solid-State Circuits*, vol. 44, no. 6, pp. 1679–1688, 2009.
- [374] S. C. Jocke, J. F. Bolus, S. N. Wooters, A. Jurik, A. Weaver, T. Blalock, and B. Calhoun, "A 2.6- μ w sub-threshold mixed-signal ecg soc," in *Symposium on VLSI Circuits*, 2009, pp. 60–61.
- [375] E. Chua, C.-C. Fu, W.-C. Fang *et al.*, "Vlsi implementation of a mixed bio-signal lossless data compressor for portable brain-heart monitoring systems," in *IEEE International Conference on Consumer Electronics*, 2011, pp. 557–558.
- [376] M. G. Rodriguez, J. Leskovec, D. Balduzzi, and B. Schölkopf, "Uncovering the structure and temporal dynamics of information propagation," *Network Science*, vol. 2, no. 01, pp. 26–65, 2014.
- [377] E. Sefer and C. Kingsford, "Convex risk minimization to infer networks from probabilistic diffusion data at multiple scales," in *IEEE International Conference on Data Engineering*, 2015.
- [378] P. Bickel, D. Choi, X. Chang, H. Zhang *et al.*, "Asymptotic normality of maximum likelihood and its variational approximation for stochastic blockmodels," *The Annals of Statistics*, vol. 41, no. 4, pp. 1922–1943, 2013.

BIBLIOGRAPHY

- [379] A. Celisse, J.-J. Daudin, L. Pierre *et al.*, “Consistency of maximum-likelihood and variational estimators in the stochastic block model,” *Electronic Journal of Statistics*, vol. 6, pp. 1847–1899, 2012.
- [380] K. M. Hall, “An r-dimensional quadratic placement algorithm,” *Management science*, vol. 17, no. 3, pp. 219–229, 1970.
- [381] Y. Koren, “Drawing graphs by eigenvectors: theory and practice,” *Computers and Mathematics with Applications*, vol. 49, no. 11, pp. 1867–1888, 2005.
- [382] S. Neph, A. B. Stergachis, A. Reynolds, R. Sandstrom, E. Borenstein, and J. A. Stamatoyannopoulos, “Circuitry and dynamics of human transcription factor regulatory networks,” *Cell*, vol. 150, no. 6, pp. 1274–1286, 2012.