

Identifying Prognostic Biomarkers for Cancer Using Gene Expression Data

by

Maha Shady

B.S., Massachusetts Institute of Technology (2016)

Submitted to the Department of Electrical Engineering and Computer
Science

in partial fulfillment of the requirements for the degree of

Master of Engineering in Computer Science and Molecular Biology

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2019

© Massachusetts Institute of Technology 2019. All rights reserved.

Author
Department of Electrical Engineering and Computer Science
May 24, 2019

Certified by.....
Manolis Kellis
Professor
Thesis Supervisor

Accepted by
Katrina LaCurts
Chairman, Department Committee on Graduate Theses

Identifying Prognostic Biomarkers for Cancer Using Gene Expression Data

by

Maha Shady

Submitted to the Department of Electrical Engineering and Computer Science
on May 24, 2019, in partial fulfillment of the
requirements for the degree of
Master of Engineering in Computer Science and Molecular Biology

Abstract

With the advancements in next-generation sequencing technologies, there is an unprecedented amount of tumor genomic and transcriptomic data available for research. This data availability opens the door for using data-driven techniques to identify biomarkers that can provide useful information about a patient's tumor, such as progression risk and potential sensitivity to treatment. Given the heterogeneity of the disease, the availability of such patient-specific information to clinicians will enable personalized care and improve clinical outcomes. In this thesis we utilize gene expression data from The Cancer Genome Atlas to identify prognostic biomarkers and predictors of potential response to immunotherapy in lung adenocarcinoma (LUAD) and metastatic skin cutaneous melanoma (SKCM). We utilize various statistical and machine learning techniques in the identification of such biomarker genes, and we analyze the final candidate genes for their utility in characterizing the immune response to the tumor and in stratifying patients by survival time. For both LUAD and SKCM we identified genes which can be reliable predictors of patient survival and whose enrichment correlates with immune cell abundance in the tumor micro-environment, specifically with CD4 and CD8 T-cells, and macrophages. We also found that some of the genes we identified represent known cancer hallmarks. Finally, a subset of the genes we identified have been shown to be predictive of response to immunotherapy in melanoma patients in a previous study. These results show that the identified genes may have biological and clinical significance in the context of cancer. Furthermore, the methodology we used to identify those biomarkers was not dependent on any disease-specific information, and can be adjusted for utility in other cancer types or even other diseases.

Thesis Supervisor: Manolis Kellis
Title: Professor

Acknowledgments

The results shown here are in whole or part based upon data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>.

First, I would like to thank Prof. Manolis Kellis for welcoming me into his lab, and for his mentorship and support throughout this project. For the past two years, my interactions with him and his lab members taught me a tremendous amount about the field of Computational Biology, and about what it means to be a successful scientist. I am inspired by his enthusiasm for research and scientific discovery.

I would especially like to thank Vinod Yadav, a post-doctoral fellow in the Kellis lab for his close mentorship and guidance through out this project. I could not have taken this project as far as I did without Vinod's guidance and expertise. His patient mentorship style helped me get through the many obstacles I encountered throughout the project. He helped me continue to see the project's potential which motivated me to move forward.

I would also like to thank other Kellis lab members, especially Yongjin Park and Alivin Shi for thoughtful conversations that helped my project. Also, Nicola Rinaldi, the lab project manager, for helping me adjust to the lab and for her support throughout my time in the lab. And Nathan Higgins, the lab administrative assistant, for making all administrative matters run seamlessly.

Last, but not least, I would like to thank my family and friends for continuous support, and for providing me outlets that helped me stay sane when things got stressful.

I am especially grateful beyond words for my parents. My mom who continues to be the biggest inspiration in my life, and for teaching me how to be a scientist and what it means to be successful in all aspects of life. And my dad, who always believes in me and pushes me to shoot for the stars.

Contents

1	Introduction	17
2	Data	21
2.1	Data Sources	21
2.1.1	The Cancer Genome Atlas (TCGA)	21
2.1.2	Genotype-Tissue Expression Project (GTEx)	21
2.2	Data Modalities	22
2.2.1	Gene Expression	22
2.3	Data Processing	22
2.3.1	Gene Expression Data Processing	22
2.3.2	Clinical Data Processing	23
3	Gene Set Selection Methods Background and Implementation	25
3.1	Feature Selection Methods	26
3.1.1	Tumor-Normal Differential Expression Analysis	26
3.1.2	Minimum Redundancy Maximum Relevance	27
3.2	Dimensionality Reduction Methods	28
3.2.1	Principal Components Analysis	28
3.3	Prediction Methods	30
3.3.1	Logistic Regression	30
3.4	Work-flow Procedure and Evaluation	31
4	Results	35

4.1	Lung Adenocarcinoma (LUAD)	35
4.1.1	LUAD Gene Set Selection	35
4.1.2	Genes in Cancer Hallmark Pathways	38
4.1.3	Association with Immune Cell Signatures	43
4.2	Skin Cutaneous Melanoma (SKCM)	46
4.2.1	SKCM Gene Set Selection	46
4.2.2	Genes in Cancer Hallmark Pathways	47
4.2.3	Association with Immune Cell Signature	50
4.2.4	Differentiating Immunotherapy Responders and Non-Responders	53
5	Discussion	57
5.1	Motivation and Contributions	57
5.2	Limitations and Future Work	58

List of Figures

2-1	Procedure for obtaining patient survival time	24
3-1	Work-flow comparison procedure	33
4-1	Expression levels of the identified genes are associated with survival. (a) Heatmap showing 1-year survival rate for each gene divided by expression quartiles. The depth of color of each cell represents the percentage of patients with at least 1 year survival. There are some genes in the gene-set for which expression in the lower quartile is asso- ciated with high survival and in the upper quartile is associated with low survival, and others for which the trend is flipped. (b) GSEA plot showing that high enrichment of some of the genes in the gene-set is associated with high survival, and high enrichment for others is as- sociated with low survival.(c) Kaplan-Meier plot comparing survival function for patients with above-median enrichment for the gene-set, and patients with up-to-median enrichment for the gene-set.	37

4-4	Kaplan-Meier plots comparing high vs. low expression of EIF5A. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations.	44
4-5	Analysis of LUAD gene-set correlation with estimated immune cell abundance in the tumor environment	45
4-6	Expression levels of the identified genes are associated with survival. (a) Heatmap showing 3-year survival rate for each gene divided by expression quartiles. The depth of color of each cell represents the percentage of patients with at least 3 year survival. For most genes in the gene-set expression in the upper quartile is associated with high survival and expression in the lower quartile is associated with low survival. There are a few genes for which the trend is reversed. (b) GSEA plot showing that high enrichment of most of the genes in the gene-set is associated with high survival, and high enrichment for a few others is associated with low survival. (c) Kaplan-Meier plot comparing survival function for patients with above-median enrichment for the gene-set, and patients with up-to-medain enrichment for the gene-set.	48

4-7 Kaplan-Meier plots comparing high vs. low expression of DDX60 and PARP9. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations. 51

4-8 Kaplan-Meier plots comparing high vs. low expression of DRAM1. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations. 52

4-9 Analysis of SKCM gene-set correlation with estimated immune cell abundance in the tumor environment 54

4-10 Kaplan-Meier plots comparing high vs. low expression of AFTPH and PARAP11. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations. . 55

List of Tables

4.1	Validation Metrics for Best Performing Models for LUAD	38
4.2	List of Cancer Hallmark Pathways Represented in LUAD Gene Set .	39
4.3	SKCM Sample Breakdown	46
4.4	Best Performing Predictive Model for SKCM	47
4.5	List of Cancer Hallmark Pathways Represented in SKCM Gene Set .	49

Chapter 1

Introduction

Despite decades of research, cancer remains a leading cause of death worldwide. It is a disease which is caused by uncontrolled cell growth due to molecular alterations in genes associated with cell cycle pathways [13]. The recent advances in next-generation sequencing (NGS) technologies and the declining costs of genomic and transcriptomic sequencing has led to the collection of massive amounts of data that can be used for the molecular analysis of tumors. Such analyses has led to the characterization of the genomic landscape of several cancer types, and to uncovering trends in cancer progression and tumor heterogeneity. Due to the heterogeneity of the disease, patients with the same cancer type may respond differently to different drugs and their disease may be characterized by alterations in different pathways. This motivated the emergence of the field of 'precision oncology', where patients would receive treatments tailored to their disease [7]. To achieve the goal of precision oncology, it is important to identify biomarkers for disease progression and potential response to treatment. In this work, we utilize a work-flow that includes various statistical and machine learning techniques in the identification of such biomarkers in lung adenocarcinoma (LUAD) and metastatic skin cutaneous melanoma (SKCM). In those diseases, we are interested in biomarkers for patient prognosis and potential response to treatment. We are specifically interested in patients' immune response to tumors and how it informs prognosis and potential response to immunotherapy. We use survival time as a proxy for our phenotypes of interest. In the remainder of this introduction, we provide a

more detailed background and motivation for the problems we are addressing and for our approach.

Immunotherapy has shown promising results in the treatment of many cancer types [33, 35], and continues to be heavily researched as a cancer therapeutic. Since immunotherapy works by stimulating a patient’s immune system to fight the tumor, it should theoretically be a treatment option for any cancer patient. However, immunotherapy can have toxic side effects and immune-related complications [33, 1]. The identification of features that can predict a patient’s likelihood of responding to immunotherapy can help guide clinical decision-making, and can help in the understanding of the immunologic response to cancer. Researchers have identified various biomarkers that correlate with response to immunotherapy. Mutational burden [22] and PD-L1 expression [33] are among the biomarkers previously identified. Studies have also demonstrated the potential of gene expression profiling in predicting response to immunotherapy [33]. Each of these biomarkers alone has limited predictive power, and researchers hypothesize that a combined approach would be most informative [12].

Here, we focus specifically on lung adenocarcinoma (LUAD) and metastatic skin cutaneous melanoma (SKCM). LUAD is a subtype of non-small cell lung cancer, which is the most common type of lung cancer. Lung cancer has poor prognosis and accounts for most cancer deaths world-wide [16]. This makes it especially important to identify clinically relevant biomarkers for LUAD progression and therapy-response prediction. Immunotherapy is commonly used in treating melanoma, however, there is still a subset of patients who do not respond to treatment [20], motivating the need for immunotherapy response prediction biomarkers in SKCM.

We will use patient survival as a proxy for anti-tumor immune response and potential immunotherapy response. It has been shown that high immune cell infiltration in the tumor micro-environment, specifically T-cell infiltration, is a predictor for immunotherapy response [36, 31]. High immune-cell infiltration in the tumor micro-environment has also been shown to associate with improved survival [11]. As such, we hypothesize that survival can be a reliable proxy for immunotherapy response

prediction, as high survival can indicate higher immune-cell infiltrate in the tumor environment, which in turn can indicate potential immunotherapy response.

Our approach for gene selection will be unbiased and extendable to other cancer types or even other diseases. We will not use any disease-specific knowledge in our gene-selection workflows or to pre-select candidate genes. Rather, we will implement various statistical and machine learning techniques into four different workflows that attempt to extract signal from whole-transcriptome datasets to produce sets of candidate genes. We will then evaluate the predictive potential of each of those candidate gene-sets to select the genes to proceed with in downstream analysis and to further study with respect to their role in the immunologic response to tumors.

In this thesis, we implement various statistical and machine learning based feature selection methods for the unbiased identification of genes whose expression levels are important in the context of cancer, starting from whole-transcriptome data-sets. We specifically identify genes that may act as prognostic biomarkers and immunotherapy response predictors in lung adenocarcinoma and metastatic skin cutaneous melanoma patients.

Chapter 2

Data

2.1 Data Sources

2.1.1 The Cancer Genome Atlas (TCGA)

Tumor data from The Cancer Genome Atlas (TCGA) was used in this thesis for model training as well as downstream analysis. TCGA is a cancer genomics program led by the National Cancer Institute (NCI) Center for Genomics and the National Human Genome Research Institute (<https://www.cancer.gov/tcga>). It includes molecular profiling of over 20,000 primary tumor and normal tissue biopsies from 33 cancer types. It also includes clinical data for each patient, including survival status and time. We obtained gene expression and clinical data for patients with Lung Adenocarcinoma (LUAD) and Skin Cutaneous Melanoma (SKCM).

We obtained gene expression data from the mRNAseq analysis pipeline from the Broad Institute TCGA Genomic Data Analysis Center (GDAC). Clinical data was also obtained from GDAC. Details of our processing methodology for both gene expression and clinical data are available in the data processing section of this thesis.

2.1.2 Genotype-Tissue Expression Project (GTEx)

Gene expression data from normal non-disease tissue samples from the Broad Institute Genotype-Tissue Expression project (GTEx) was used for differential expression

analysis. RNA-seq data for differential expression analysis was obtained from the Memorial Sloan Kettering Cancer Center RNAseqDB [30] (<https://github.com/mskcc/RNAseqDB>). This database includes RNAseq data from TCGA and GTEx that has been processed by the same pipeline, which allowed us to perform accurate tumor-normal differential expression analysis without needing to process raw sequencing files.

2.2 Data Modalities

2.2.1 Gene Expression

Gene expression is measured by RNA sequencing, and provides a snapshot of the transcript abundance of different gene transcripts in the cell. This gives a closer approximation of which proteins are actually produced in the cell than could be approximated by DNA measurements, which is one of the advantages of using gene expression values to identify relevant biomarkers to disease [6].

2.3 Data Processing

In this section I will discuss the steps taken to format downloaded data and produce the data files used in feature selection and analysis.

2.3.1 Gene Expression Data Processing

For feature extraction and predictive modelling, we used normalized log₂ transformed transcript abundance values obtained using RSEM [2] from GDAC. We performed a filtering step where we removed all genes for which there were missing values in some patients. For differential expression analysis, we obtained normalized tumor TCGA data and normal GTEx data from RNAseqDB. We used the normalized FPKM (Fragments per Kilobase of transcript per Million) data files obtained using RSEM [3].

2.3.2 Clinical Data Processing

We obtained clinical data from GDAC and processed it using a procedure similar to the one used by downstream analysis performed on the GDAC pipeline. We considered each patient's `vital_status`, `days_to_death`, and `days_to_last_followup` to determine their survival time. Vital status refers to a flag of whether the patient is dead or alive. If a patient's `vital_status` was 'alive' but either their `days_to_death` field had a number of days or their `days_to_last_followup` field was empty, that patient was removed as this data was considered to be either contradictory or insufficient for determining survival time. If a patient's `vital_status` field was empty, or their vital status was 'dead', but their days to death field was empty, that patient was also removed as we did not have sufficient information to determine survival time. Furthermore, if any patient was marked as 'alive' and their `days_to_last_followup` field was 30 days or less, that patient was removed from the analysis as it is likely the patient moved to a different center, their diagnosis was a miss-diagnosis, or the field was incorrectly filled. For all remaining patients, we set a patient's survival time to be the value in their `days_to_death` field if their `vital_status` was 'dead', or to the value in their `days_to_last_followup` field if their `vital_status` was 'alive'. Figure 2-1 shows a schematic of the steps for obtaining patient survival time as described here. Finally, we created flags indicating whether each patient survived past 6 months, 1 year, 2 years, and 3 years. These additional flags were useful in downstream prediction.

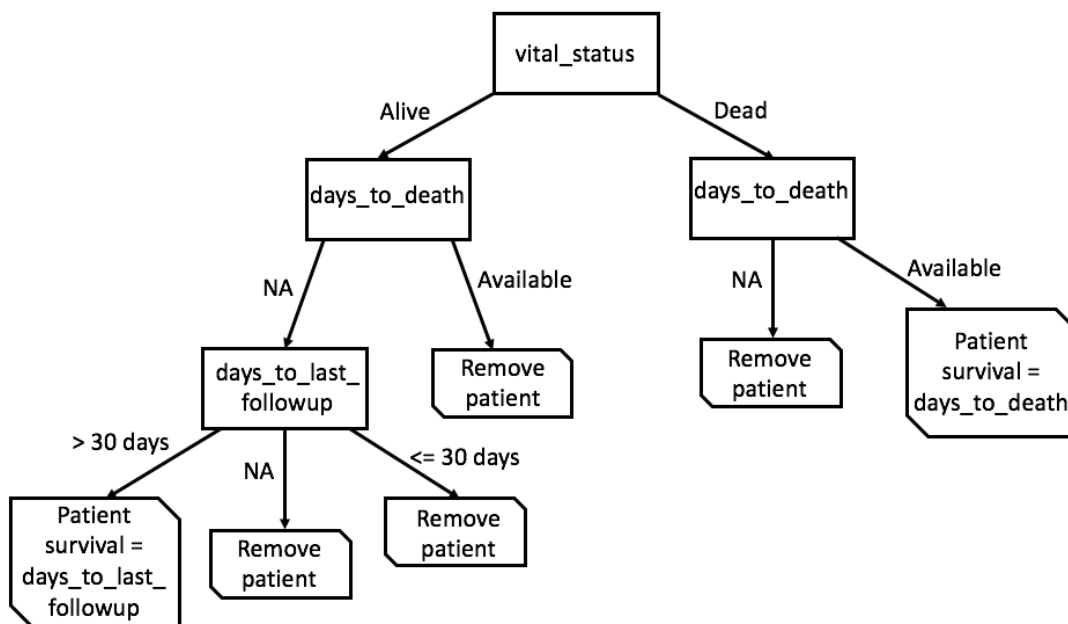


Figure 2-1: Procedure for obtaining patient survival time

Chapter 3

Gene Set Selection Methods Background and Implementation

We applied different statistical and machine learning techniques to whole transcriptome data sets with the goal of identifying genes whose expression levels are most likely to act as relevant biomarkers in the context of lung adenocarcinoma and metastatic skin cutaneous melanoma. In order to select the final gene set to proceed with, we compared the performance of different work-flows that utilize those techniques by testing the predictive power of the gene-set selected by each work-flow. Figure 3-1 shows the different work-flows compared. The next few sections describe the background and implementation details of the statistical and machine learning methods utilized. Those sections are followed by a section that describes the predictive modelling procedure and evaluation metrics used. The final part of this chapter describes how those methods were put together into work-flows, and how the work-flows were evaluated and the final gene set selected.

3.1 Feature Selection Methods

3.1.1 Tumor-Normal Differential Expression Analysis

Background

In differential expression analysis, normalized gene expression levels are analyzed to characterize the difference in expression of a particular gene between different cohorts. The goal in tumor-normal differential expression analysis is to identify genes for which the expression levels in the disease cohort differs significantly from the expression level in the same tissue within a healthy cohort. A statistically significant difference in expression levels, means that the difference is too big to have been caused by random chance, and might point to an underlying genotypic difference between the groups. In such comparisons, gene expression levels must be normalized before comparison to control for differences in library preparation and complexity.

Implementation Details

To perform tumor-normal differential expression analysis, we obtained normalized FPKM data for tumor and normal samples from the MSKCC RNAseqDB, which includes TCGA and GTEx gene expression data (see Data section for more details). We compared each tumor type, against the normal expression levels of the tissue associated with the disease. After obtaining the data, we performed a Mann-Whitney test to compare the tumor and normal expression levels for each gene. We used the Python package SciPy to run the Mann-Whitney test [18]. If the p-value of the Mann-Whitney test was a certain threshold, we considered the differences between the two cohorts to be statistically significant. Genes with a statistically significant difference between the two cohorts were considered differentially expressed, and passed onto the next step of analysis. To account for the multiple testing problem, we applied the Bonferroni correction to our p-value threshold, setting it to $0.05/(\# \text{ of genes tested})$.

3.1.2 Minimum Redundancy Maximum Relevance

Background

Many feature selection methods focus on 'maximum relevance', meaning they attempt to identify variables which are most likely associated with the classification variable of interest. Those methods often don't take into account the relationship between the selected features themselves, leading to some redundancy in the selected feature subset. For example if gene A was selected as relevant to the classification variable of interest, another gene B which is highly correlated with A is likely to be selected as well. However, in this scenario, gene B may not add new information to the selected feature subset, making it a redundant feature. Minimum Redundancy Maximum Relevance (mRMR) is a feature selection algorithm which attempts to address this problem by taking into account correlation between features being selected [9]. In this thesis, we follow the mRMR procedure presented in [9] with some changes.

Ding et. al use the F -statistic to measure the relevance of each gene to the desired classification variable. Let Ω be the original gene set, S is the set of genes selected so far, and $\Omega_S = \Omega - S$ is the remaining set of genes being selected from. Also, let h be the classification variable, $F(g_i, h)$ be the F -statistic of gene g_i with respect to h , and $c(g_i, g_j)$ be the correlation coefficient of gene g_i with gene g_j . For a feature set S maximum relevance is calculated by:

$$\begin{aligned} & \max V_F, \\ V_F &= \frac{1}{|S|} \sum_{i \in S} F(i, h) \end{aligned} \tag{3.1}$$

And minimum redundancy is calculated by:

$$\begin{aligned} & \min W_c \\ W_c &= \frac{1}{|S|^2} \sum_{i, j \in S} |c(i, j)| \end{aligned} \tag{3.2}$$

[9]

The mRMR feature set is selected by optimizing both Eqs. 3.1 and 3.2 simultaneously by combining into the F -test correlation quotient (FCQ):

$$\begin{aligned} & \max(V_F/W_c) \\ & = \max_{i \in \Omega_S} \left[F(i, h) / \left(\frac{1}{|S|} \sum_{j \in S} |c(i, j)| \right) \right] \end{aligned} \quad (3.3)$$

[9] In the overall algorithm, the first feature is selected by identifying the gene with the maximum F -statistic with respect to the desired classification variable, and the remainder of the features are selected by optimizing eq. 3.3 [9].

Implementation Details

In our implementation of mRMR, our classification variable of interest was the survival cutoff used for each cancer type. For example, for Lung Adenocarcinoma, our classification variable was 1-year survival status. We used one-way Analysis of Variance (ANOVA) to calculate the F -statistic and Spearman-rank correlation to compute the correlation coefficients between pairs of genes. ANOVA is a statistical hypothesis test for which the null hypothesis is that the groups being tested are all random samples from the same population. Our target selected feature set size was 50 genes, and we used the Python library SciPy [18] to implement both ANOVA and Spearman-rank correlation tests.

3.2 Dimensionality Reduction Methods

3.2.1 Principal Components Analysis

Background

Principal components analysis (PCA) is a dimensionality reduction procedure, in which a feature matrix is transformed into a lower dimension matrix with new features that are linearly uncorrelated with one another. These features are called principal

components (PCs) and are computed by singular value decomposition of the original data matrix. The principal components in the transformed matrix are ordered by the amount of variance they explain of the original data, such that the first principal component explains the most variance of the original data.

Implementation Details

To perform dimensionality reduction of gene expression data, we used the Python library Scikit-learn [26] with default parameters. PCA was applied only to differentially expressed genes in datasets where tumor-normal differential expression analysis was performed. After transforming the data, we selected the principal components that explained at least 70% of the variance in the original data, and discarded the remaining PCs. Using the transformed matrix, we proceeded in three different directions, shown in Figure 3-1. 1) We ran the remaining feature selection and survival prediction procedure directly on the transformed matrix, and then mapped the PC's selected by that workflow back to genes, 2) we first mapped the top PC's back to genes, and then proceeded with the feature selection and survival prediction procedures only on those genes, and 3) We identified the most highly correlated gene with each PC and built a predictive model only on those genes. In order to map PC's to genes we used Scipy's Spearman-rank correlation test implementation [18] to identify genes which have a high correlation coefficient with each PC. We considered Spearman-rank correlation to be an appropriate measure because it is a non-parametric test and does not assume that the datasets being compared are normally distributed. We used a multiple-hypothesis corrected p-value cutoff to identify highly correlated genes, such that $\alpha = \frac{0.05}{(\#PCs)*(\#genes)}$, where α is the p-value cutoff used for significance. We noted that only the top 8-11 PC's tend to be significantly correlated with genes, even if many more PCs (40 - 45) were needed to explain at least 70% of the variance of the original data.

3.3 Prediction Methods

3.3.1 Logistic Regression

Background

Logistic regression is a classification algorithm that models the probability that a sample belongs to each class by a logistic function. When applying logistic regression it is often advisable to use regularization to control the complexity of the model and avoid over-fitting to the training data. Two popular regularization methods are l_1 and l_2 regularization. l_1 regularization attempts to minimize the l_1 norm of the weights assigned to the feature space, which leads to assigning zero weights to many of the features. l_2 regularization minimizes the sum of squares of the weights, forcing some of the weights to be small but doesn't force them to zero.

To evaluate each model's performance, we computed the accuracy and the Area Under the Receiver Operating Characteristic Curve (ROC AUC). The AUC computes the probability that the model ranks a randomly chosen positive sample higher than a randomly chosen negative sample. In our case, a 'positive sample' is one where survival is at least the threshold used for that cancer type (for example 1 year for lung adenocarcinoma), and a 'negative sample' is one for which survival is below the threshold. The AUC takes into account both the true positive rate and the false positive rate of the classifier. Accuracy compares the number of correctly classified positives and correctly classified negatives to the total number of actual positives and negatives in the population.

Implementation

We implemented logistic regression using the Python library Scikit-learn [26]. We implemented two logistic regression models for each proposed gene-set, an l_1 -regularized model and an l_2 -regularized model. The Scikit-learn function used implemented 10-fold cross-validation for parameter selection. The evaluation parameters were also computed using the Scikit-learn library [26].

3.4 Work-flow Procedure and Evaluation

In this section we describe how the methods described were put together into work-flows for gene-set recommendation, and how those work-flows were evaluated. To identify the gene-set for each cancer type, we started by splitting the patients into a 'modelling' set and 'validation set'. The modelling set comprised 80% of the data for lung adenocarcinoma and 90% for skin cutaneous melanoma, and it was used for feature selection and predictive modelling. The validation set was used for final model evaluation and gene-set identification. In creating the modelling-validation split, we used scikit-learn's [26] StratifiedShuffleSplit function to ensure that both the modelling and validation sets have equal proportions of patients above and below the specified survival cutoff. The survival cutoff for each cohort was chosen to minimize class-imbalance.

We then proceeded to perform tumor-normal differential expression analysis as described above on the modelling set. Once a matrix of survival and gene-expression data with only differentially expressed genes was obtained, we split the modelling samples further into a training set and a testing set using stratified split. The training set comprised 80% of the modelling data. We proceeded to perform feature selection on the training set, and went in two different directions. 1) We ran the mRMR algorithm directly on that matrix and built a predictive model based on the 50 selected genes. 2) We performed principal component analysis to further reduce the dimensionality of that matrix and created a transformed matrix with the principal components that explained at least 70% of the variance in the data. We then proceeded in three different directions from that matrix. 2a) We ran the mRMR algorithm directly on the PCs, then identified the genes correlated with the mRMR-selected PCs and built a predictive model based on those genes. 2b) We identified the genes which correlate with the PCs, ran the mRMR algorithm on all the correlated genes, and built a predictive model based on the selected genes. 2c) We built a predictive model directly on the set created by identifying the most correlated gene to each PC. All predictive models built were evaluated on the testing set.

Next, we validated each gene-set's predictive potential using the following steps:

- 1) We created a sub-matrix of the validation matrix with only the genes in that gene-set.

- 2) We used the predictive mode built for that gene-set on the validation sub-matrix, and computed the prediction accuracy and AUC.

We then selected the best performing gene-sets based on the results of the validation evaluation metrics and combined those gene-sets into our proposed set of biomarker genes for that cancer type. Figure 3-1 shows a schematic summarizing the steps described here.

Finally, we proceeded to perform further analysis to evaluate the proposed gene-set's potential as a prognostic biomarker and immunotherapy-response predictor. This analysis is further discussed in the Results chapter. In the downstream gene analysis, all Kaplan-Meier plots were created using the Python library Scikit-survival [28, 29, 27], and all logrank tests for comparing survival functions were performed using the Python Lifelines library [8].

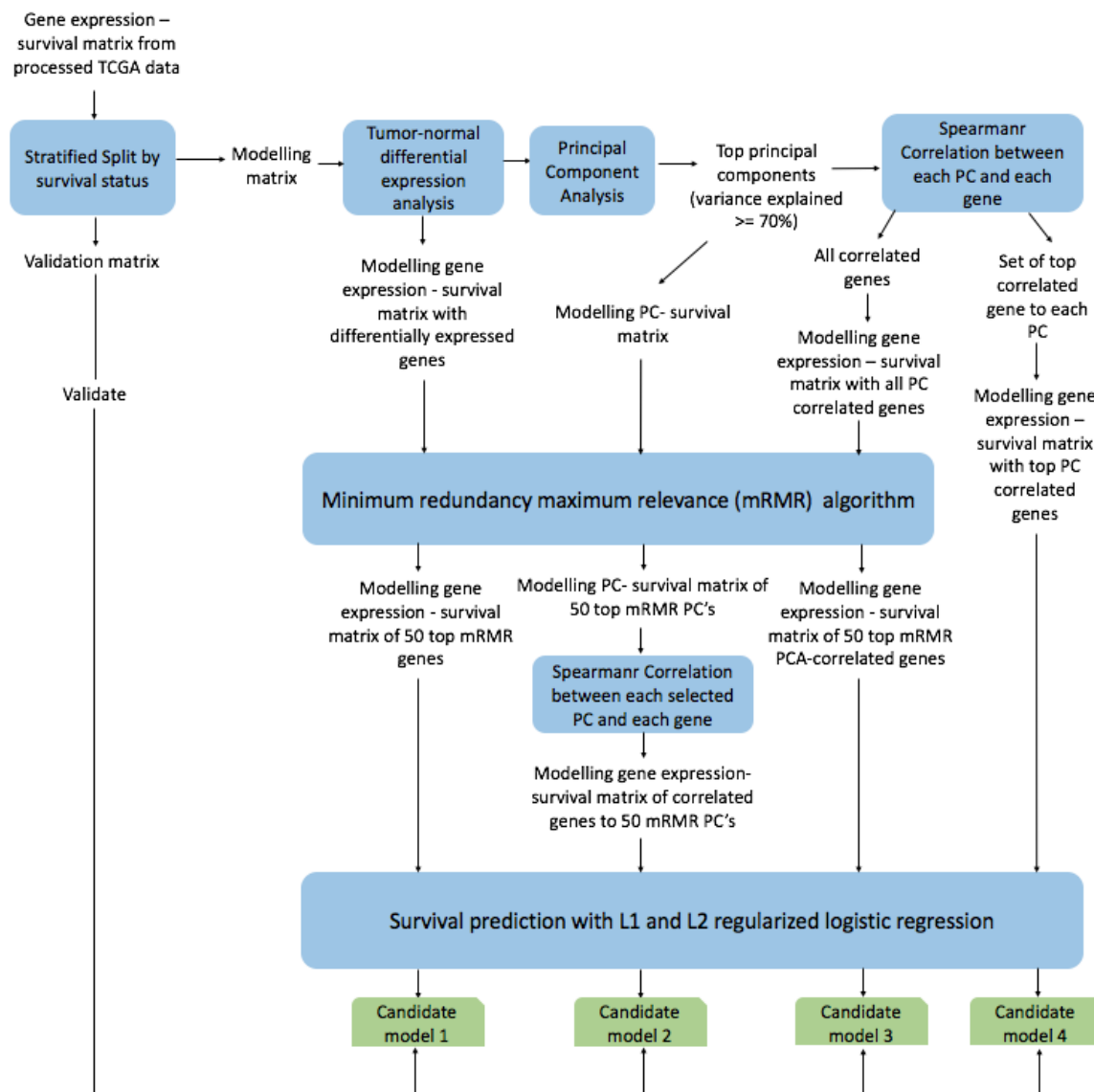


Figure 3-1: Work-flow comparison procedure

Chapter 4

Results

4.1 Lung Adenocarcinoma (LUAD)

Lung adenocarcinoma is a type of non-small cell lung cancer (NSCLC), which constitutes the biggest percentage of lung cancer diagnoses. Lung cancer in general has poor prognosis, so the identification of prognostic or therapy-response predictive biomarkers can be of clinical benefit [16].

In the lung adenocarcinoma cohort, we had 412 patients with RNA-seq data from primary tumors and valid clinical data. After filtering out genes with missing gene-expression data in some of those patients, we had 13,431 genes remaining for analysis. We ran tumor-normal differential expression analysis for this cohort against the GTEx lung cohort, and found 7,859 differentially expressed genes which we proceeded with in the remainder of the analysis. We used 1-year survival as a cutoff for low survival for this cohort.

4.1.1 LUAD Gene Set Selection

Table 4.1 shows the prediction metrics of the best performing feature-selection protocols for this cohort. The two methods we found most effective were: 1) performing mRMR analysis on the genes correlated with the top PCs, and 2) identifying the genes most highly correlated with each of the top PCs, one gene for each PC. When

analyzing the PCA results, we found that the top 40 principal components were required to explain 70% of the data. However, only the top 8 PCs as well as the 11th PC were found to be highly correlated with some genes. Those PCs explain about 47% of the variance in the data, and significantly correlate with 2081 genes.

We proceeded with a gene-set that combines the genes identified by both protocols. This was a total of 58 genes as there was one gene (ANKRD36) in common. We ran gene-set enrichment analysis on this gene-set using the Gene Set Enrichment Analysis in Python library (GSEAPy) [24, 32]. Figure 4-1 panel (b) was produced from GSEAPy library and shows that enrichment of the selected gene set in the tumor is correlated with patient survival status. Specifically, there are some genes in the gene-set for which enrichment correlates with low survival, and another set of genes for which enrichment correlates with high survival. Also, panel (a) of figure 4-1 shows how 1-year survival is affected by the expression value of each gene. The figure shows some genes for which expression values in the lower quartile is associated with better survival than expression values in the upper quartile and vice-versa. Quartiles were calculated separately for each gene based on its expression distribution in this cohort. Both of these figures support our belief that the selected gene-set is associated with patient survival and includes useful prognostic biomarkers for patients with lung adenocarcinoma.

Figure 4-1 panel (c) is a Kaplan-Meier plot that shows no significant difference in survival function between patients with high enrichment and low enrichment of our gene-set. We can make sense of this result in context of the other two panels in figure 4-1. As discussed above, both of those panels demonstrate that there are some genes in the gene-set for which enrichment is associated with better survival, and others for which enrichment is associated with less-favorable survival. So enrichment of the gene-set as a whole is not expected to differentiate patient survival. Further, it is possible that low or high enrichment of certain genes correlates with survival only for a subset of patients who satisfy other criteria, and in that case we would not necessarily see a trend in the overall population.

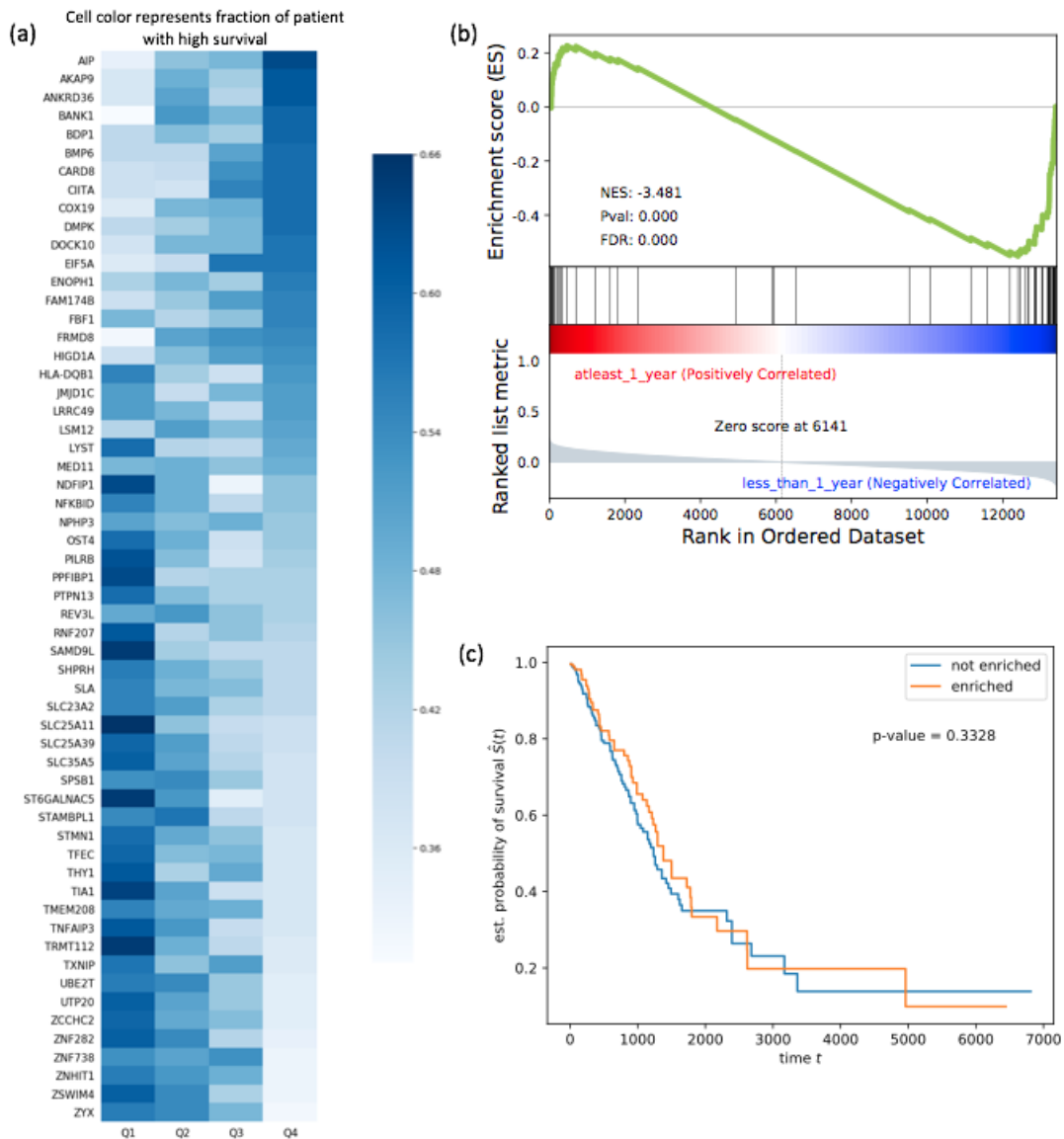


Figure 4-1: Expression levels of the identified genes are associated with survival. (a) Heatmap showing 1-year survival rate for each gene divided by expression quartiles. The depth of color of each cell represents the percentage of patients with at least 1 year survival. There are some genes in the gene-set for which expression in the lower quartile is associated with high survival and in the upper quartile is associated with low survival, and others for which the trend is flipped. (b) GSEA plot showing that high enrichment of some of the genes in the gene-set is associated with high survival, and high enrichment for others is associated with low survival. (c) Kaplan-Meier plot comparing survival function for patients with above-median enrichment for the gene-set, and patients with up-to-median enrichment for the gene-set.

Table 4.1: Validation Metrics for Best Performing Models for LUAD

Gene Selection Protocol	# selected genes	Model	Validation AUC	Validation Accuracy
mRMR on PC correlated genes	50	L1 logistic regression	0.707	0.651
		L2 logistic regression	0.707	0.706
genes correlating with top PCs	9	L1 logistic regression	0.688	0.614
		L2 logistic regression	0.689	0.614

4.1.2 Genes in Cancer Hallmark Pathways

We used the Broad Institute’s MutSigDB [32, 24] to identify cancer hallmarks represented in our gene-set. Cancer hallmarks are a set of biological processes affected during tumor development [14]. We found that the gene-set includes genes from 24 different cancer hallmark pathways and that 19 out of our 58-gene set are part of at least one of those pathways. Table 4.2 shows the list of cancer hallmark pathways represented and the genes in our gene-set which are part of those pathways.

The Interferon-gamma Response pathway is significantly represented in our gene set. The statistical significance of hallmark representation was determined by MutSigDB [32, 24] based on the hypergeometric distribution, with a multiple hypothesis testing corrected p-value cutoff. Interferon-gamma (IFN- γ) is a cytokine involved in both innate and adaptive immunity. It has been shown to be associated with tumor-elimination [5]. One of the IFN- γ associated genes in our geneset is CIITA, which is a class II major histocompatibility complex (MHC class II) regulator [19]. MHC proteins are cell-surface proteins which allow the immune system to identify foreign particles [19]. IFN- γ has been show to upregulate CIITA expression, leading to MHC class II expression, increasing the immune surveillance of the tumor [5]. Figure 4-2 shows a Kaplan-Meier plot comparing the survival function of patients with high tumor CIITA expression to that of patients with low tumor CIITA expression in our cohort. Consistent with previous findings, the plot shows that high expression of CI-

Table 4.2: List of Cancer Hallmark Pathways Represented in LUAD Gene Set

Hallmark	Genes
TNFA SIGNALING VIA NFKB	SPSB1, TNFAIP3
HYPOXIA	TNFAIP3
DNAREPAIR	REV3L
G2M CHECKPOINT	STMN1
APOPTOSIS	TXNIP
MYOGENESIS	DMPK
INTERFERON ALPHA RESPONSE	SAMD9L, TXNIP
INTERFERON GAMMA RESPONSE	BANK1, CIITA, SAMD9L, TNFAIP3, TXNIP
APICAL JUNCTION	THY1, ZYX
APICAL SURFACE	THY1
HEDGEHOG SIGNALING	THY1
COMPLEMENT	DOCK10, TNFAIP3
PI3K AKT MTOR SIGNALING	SLA
MTORC1 SIGNALING	SLA
E2F TARGETS	STMN1, UBE2T
MYC TARGETS V2	UTP20
EPITHELIAL MESENCHYMAL TRANSITION	THY1, TNFAIP3
OXIDATIVE PHOSPHORYLATION	SLC25A11
GLYCOLYSIS	STMN1
P53 PATHWAY	TXNIP
BILE ACID METABOLISM	BMP6, SLC23A2
PEROXISOME	SLC23A2
ALLOGRAFT REJECTION	EIF5A, THY1
KRAS SIGNALING UP	TNFAIP3

ITA is associated with better survival. The difference in survival was determined to be statistically significant via the logrank test, with a p-value of 0.0394. The plot also includes a sub-plot which serves to cross-validate the significance of the difference between survival functions in the original plot. The subplot was created by making 1000 label permutations and computing the survival function for above-median expression and below-median expression for each permutation, and then computing the p-value of the difference between the two survival functions. These permutations allows us to observe the background p-value distribution and determine whether our data is biased towards a lower p-value. We plot the histogram of p-values across permutations and we put a red line on the histogram indicating where the p-value of our original result lies on that distribution. We observe a uniform distribution of p-values, as expected of p-values under the null-hypothesis. This suggests that the data is not inherently biased towards lower p-values, supporting the significance of our result.

Another hallmark represented in our gene-set is the E2F targets hallmark. E2F genes encode transcription factors with targets involved in the cell cycle pathway [4]. One of the E2F targets in our gene-set is UBE2T, which has been shown to be involved in the DNA repair pathway [34]. UBE2T has been shown to have oncogenic properties in prostate cancer [34]. In our cohort, we find that its over-expression is associated with poor prognosis as shown in figure 4-3. The difference in survival function between patients with low versus high expression of UBE2T was determined to be statistically significant by the logrank test, with a p-value of 0.0111. Again, we validate this result with 1000 permutations of survival labels, and we find the p-values obtained from those permutations follow a uniform distribution as expected under the null hypothesis. The result of these permutations suggest that the data is not biased towards lower p-values and it supports the significance of our result.

Allograft rejection is a third hallmark represented in our gene-set, and it includes genes that are up-regulated during transplant rejection. One of the allograft rejection hallmark genes in our gene-set is EIF5A. EIF5A is a translation factor that has been shown to correlate with the development of several cancer types [23]. EIF5A has been shown to correlate with poor prognosis in NSCLC, among other cancer types

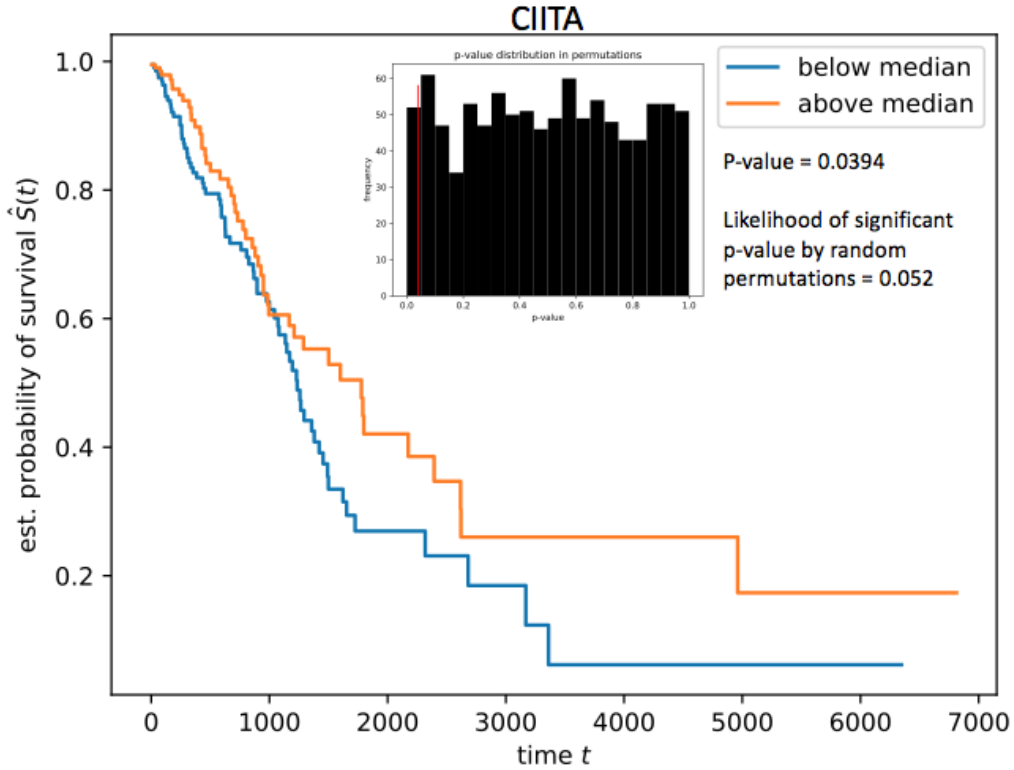


Figure 4-2: Kaplan-Meier plot comparing high vs. low expression of CIITA. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations.

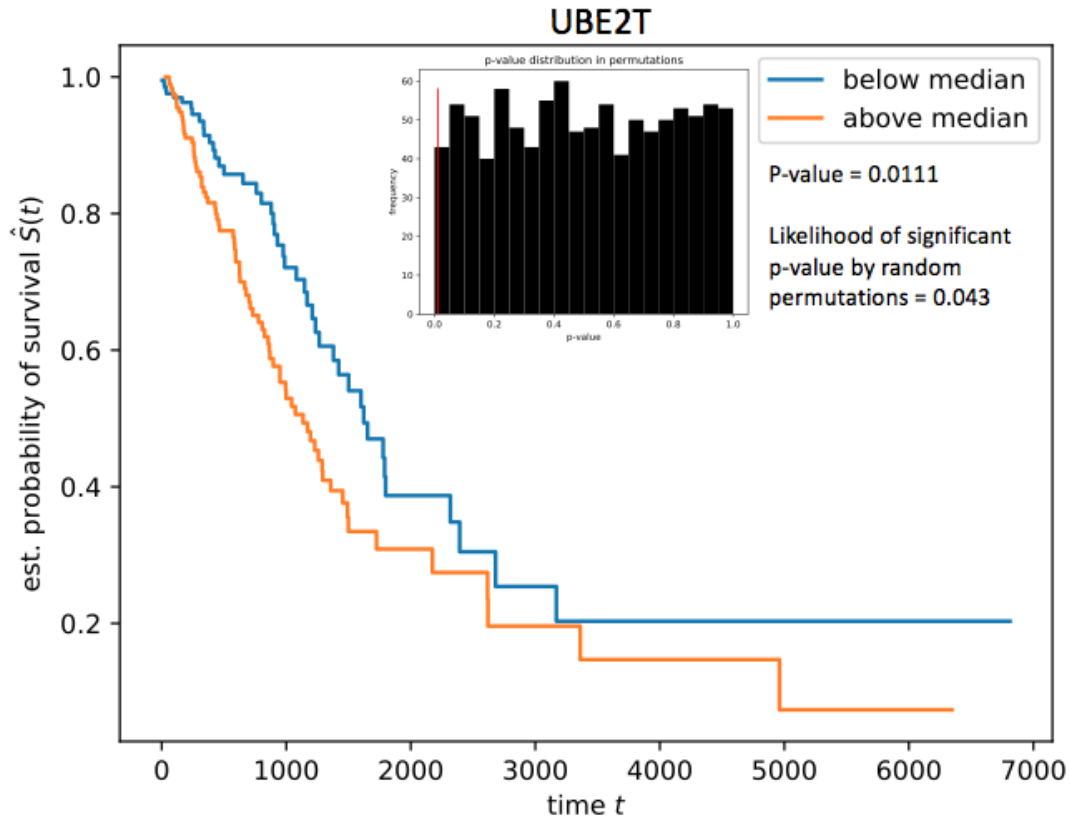


Figure 4-3: Kaplan-Meier plot comparing high vs. low expression of UBE2T. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations.

[23]. Figure 4-4 includes a Kaplan-Meier plot comparing survival in patients with high versus low expression of EIF5A in our cohort. Consistent with previous findings summarized in [23], we find that high expression is associated with poorer survival. The difference was found to be statistically significant by a logrank test, with a p-value of 0.0277. Again we found that only a small fraction of the 1000 random label permutations lead to a statistically significant difference, and that the p-values obtained from those permutations are uniformly distributed as is expected under the null-hypothesis. This supports the significance of our result, as it suggests that the data is not biased towards lower p-values.

4.1.3 Association with Immune Cell Signatures

We used CIBERSORT [25] to estimate immune cell abundance of tumors in our cohort. CIBERSORT makes this estimation by deconvolving the bulk gene expression data of the tumor samples using immune cell gene expression signatures [25]. We then computed the correlation between each immune cell type’s abundance and the enrichment score of our gene-set to determine whether our gene-set associates with immune-cell signatures. Association with immune-cell signatures would suggest that our gene-set can be used to infer immune cell content in the tumor environment, as well as immunotherapy response potential. Figure 4-5 describes the procedure and results of this analysis. The correlation coefficients in this figure are Spearman’s rank correlation coefficients, and significance was determined based on a multiple-hypothesis-corrected p-value of $\frac{0.05}{22} = 0.0023$. The rows with a blue border in the table in figure 4-5 are the ones where the correlation was deemed statistically significant based on this cutoff.

Among the immune cell signatures that correlate with our gene set are CD8 T-cells and activated memory CD4 T-cells. Both CD8 and CD4 T-cells have been shown to fight cancer [10]. Further, enhanced T-cell response can improve immunotherapy response and reduce tumor immune evasion mechanisms, leading to better clinical outcomes [10]. Our gene-set’s correlation with CD8 and CD4 T-cell signatures supports the gene-set’s potential in forecasting anti-tumor immune response and predicting

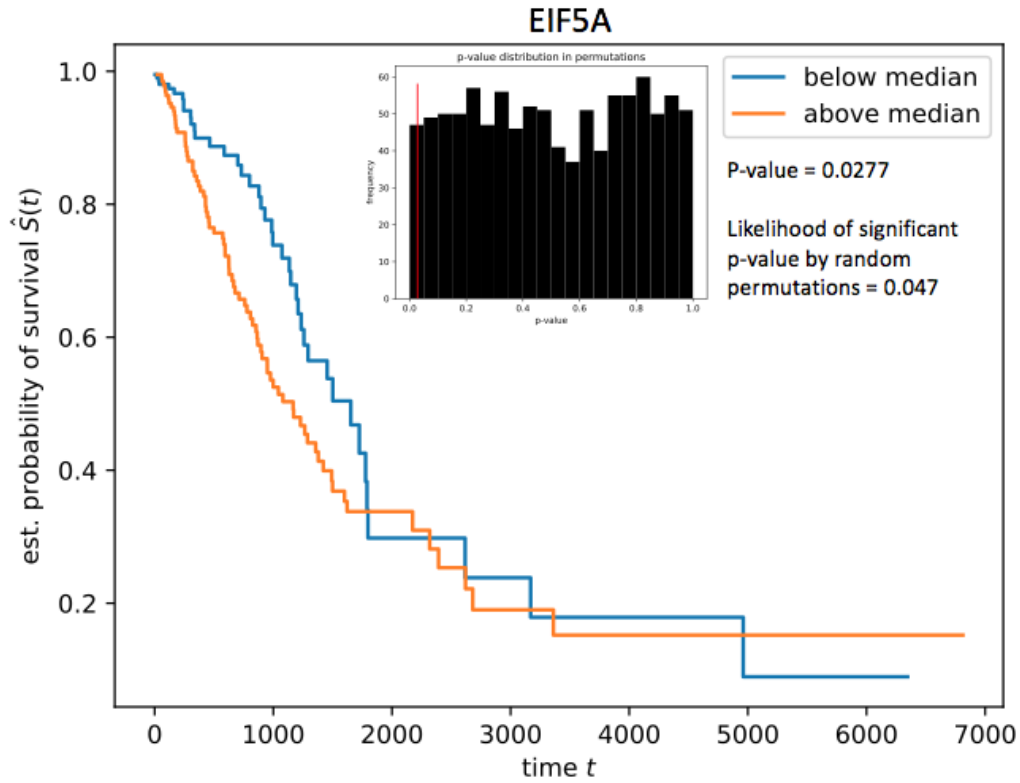


Figure 4-4: Kaplan-Meier plots comparing high vs. low expression of EIF5A. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations.

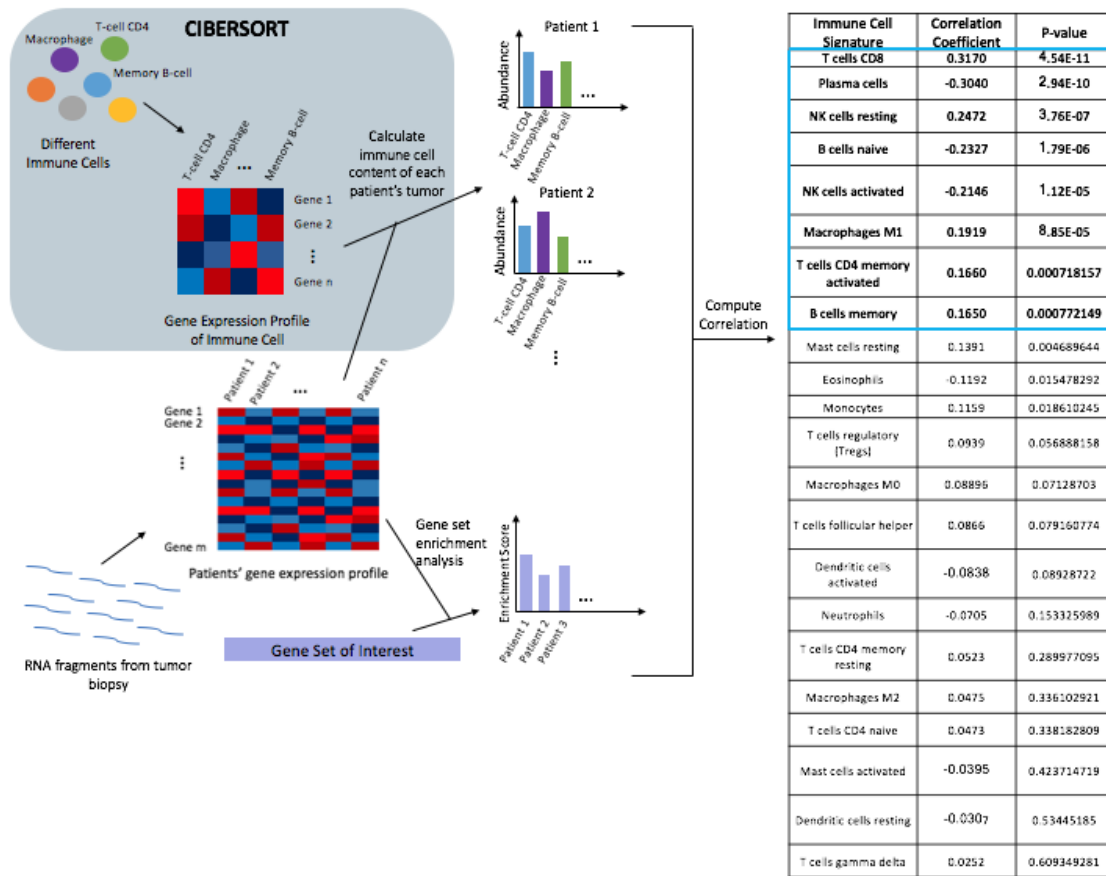


Figure 4-5: Analysis of LUAD gene-set correlation with estimated immune cell abundance in the tumor environment

immunotherapy response potential.

Our gene-set also correlates with macrophages M1 cell signature. Macrophages are a type of white blood cells of the immune system that engulf and digest unrecognized particles and cells. Macrophages are very diverse and can have both pro-tumor and anti-tumor effects when present in the tumor micro-environment [2]. M1 macrophages are believed to have anti-tumor activity [2] and a recent study showed that high levels of M1 macrophages in the tumor micro-environment is associated with improved survival in non-small cell lung cancer (NSCLC) [17]. Thus, this result again demonstrates the potential of our gene-set as a prognostic biomarker in LUAD.

Table 4.3: SKCM Sample Breakdown

Disease Status	Tissue Site	Number of samples
Primary Tumor	cutaneous or subcutaneous	36
Metastatic Tumor	distant metastasis	59
Metastatic Tumor	cutaneous or subcutaneous	71
Metastatic Tumor	regional lymph node	212

4.2 Skin Cutaneous Melanoma (SKCM)

Most of the TCGA SKCM samples were metastatic samples, as such, the overall cohort was quite heterogeneous. Table 4.3 shows the breakdown of sample types we had after processing and merging gene expression and clinical data. We performed initial model analysis on all four sub-cohorts shown in table 4.3, but had very poor performance on all but the samples obtained from regional lymph nodes. This was not surprising because the other sub-cohorts are relatively small and, in the case of samples obtained from distant metastases, can be very heterogeneous. As such we proceeded with only the samples obtained from regional lymph nodes of patients with metastatic disease.

4.2.1 SKCM Gene Set Selection

Upon filtering out genes with missing gene expression values, we had 12,208 remaining genes. We did not perform tumor-normal differential expression analysis on this cohort. We used a survival cutoff of three years for this cohort, as this led to the most balanced class partition.

We found that the best performing models for SKCM used genes selected by running the mRMR algorithm on genes highly correlated with the principal components explaining $\geq 70\%$ of the data. The top 35 PCs were required to explain at least 70% of the data, however, we found that only the top 8 PCs, the 10th PC and the 14th PC were highly correlated with any particular gene. Those components explain

Table 4.4: Best Performing Predictive Model for SKCM

Gene Selection Protocol	# selected genes	se-Model	Validation AUC	Validation Accuracy
mRMR on PC correlated genes	50	L1 logistic regression	0.667	0.636
		L2 logistic regression	0.717	0.682

about 48.7% of the variance in the data, and are highly correlated with 2446 genes.

Table 4.4 shows the evaluation metrics of the best performing model on the validation set. Figure 4-6 panel (a) shows how patient 3-year survival rate changes with each gene’s expression. Specifically, we can see that for most genes, patients with an expression value in the upper quartile for that gene have a much higher three-year survival rate than patients whose expression value falls in the lower quartile, and there are a few genes for which the trend is reversed. Quartiles were calculated separately for each gene, based on its expression distribution in the cohort. Similarly, panel (b) of figure 4-6 demonstrates that 3 year survival is associated with high enrichment of most genes in our selected gene-set, and that there are a few genes for which enrichment is associated with lower survival. Finally, panel (c) of figure 4-6 is a Kaplan-Meier plot, which shows that patients with tumors enriched for our gene set have a more favorable survival function than patients whose tumors are not enriched for the gene set. The difference between the two survival functions was found to be statistically significant using the logrank test.

4.2.2 Genes in Cancer Hallmark Pathways

We again used the Broad Institute’s MutSigDB [32, 24] to identify cancer hallmarks represented in our gene-set. Cancer hallmarks are a set of biological processes affected during tumor development [14]. 14 out of our 50-gene set are members of a total of 19 cancer hallmarks. While no pathways representation in our gene set was deemed sta-

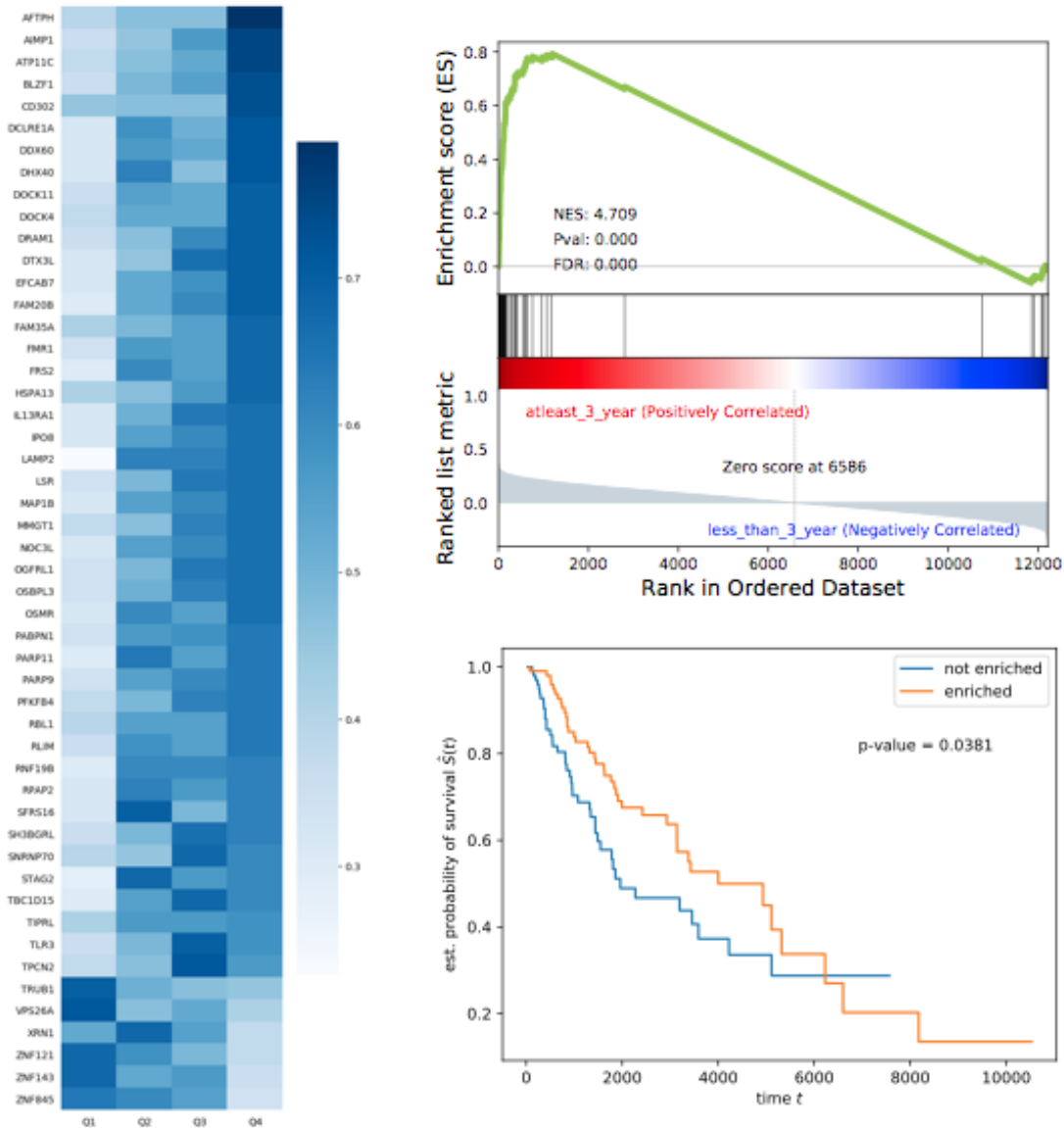


Figure 4-6: Expression levels of the identified genes are associated with survival. (a) Heatmap showing 3-year survival rate for each gene divided by expression quartiles. The depth of color of each cell represents the percentage of patients with at least 3 year survival. For most genes in the gene-set expression in the upper quartile is associated with high survival and expression in the lower quartile is associated with low survival. There are a few genes for which the trend is reversed. (b) GSEA plot showing that high enrichment of most of the genes in the gene-set is associated with high survival, and high enrichment for a few others is associated with low survival. (c) Kaplan-Meier plot comparing survival function for patients with above-median enrichment for the gene-set, and patients with up-to-median enrichment for the gene-set.

Table 4.5: List of Cancer Hallmark Pathways Represented in SKCM Gene Set

Hallmark	Genes
TNFA SIGNALING VIA NFKB	DRAM1, RNF19B
MITOTIC SPINDLE	DOCK4
IL6 JAK STAT3 SIGNALING	IL13RA1, OSMR
G2M CHECKPOINT	RBL1
ADIPOGENESIS	CD302
ESTROGEN RESPONSE LATE	LSR
PROTEIN SECRETION	LAMP2
INTERFERON ALPHA RESPONSE	DDX60, PARP9
INTERFERON GAMMA RESPONSE	DDX60
COMPLEMENT	DOCK4, LAMP2
INFLAMMATORY RESPONSE	OSMR, TLR3
GLYCOLYSIS	IL13RA1
P53 PATHWAY	DRAM1, RNF19B
UV RESPONSE UP	HSPA13
UV RESPONSE DN	MAP1B
HEME METABOLISM	LAMP2
COAGULATION	LAMP2
ALLOGRAFT REJECTION	TLR3

tistically significant by the hypergeometric distribution, the presence of genes which are part of those pathways in our gene-set points to the genes' potential as cancer biomarkers. Table 4.5 shows the hallmarks represented and the genes representing them in our gene-set.

Two of the genes are part of the Interferon-alpha response hallmark (DDX60 and PARP9), and DDX60 is also part of the Interferon-gamma response hallmark. As discussed in the previous section, interferon-gamma is a cytokine involved in immune response, and has been shown to be associated with both tumor elimination and tumor immune-system evasion [5]. Similarly interferon-alpha is another cytokine involved in immune response. Figure 4-7 shows Kaplan-Meier curves for both PARP9 and DDX60, comparing survival for patients with low versus high expression of those genes. The difference in survival was found to be statistically significant using the logrank test, with a p-value of 0.0003 for DDX60 and 0.0 for PARP9. Similar to the Kaplan-Meier plots in the previous section, the plots also include validation subplots

showing distribution of p-values obtained by comparing high-expression and low-expression survival functions based on 1000 random label permutations. For both genes we find that only a small fraction of the permutations lead to a significant p-value, and that the p-value distribution across the random permutations is uniform, which is expected under the null hypothesis. This suggests that the data is not biased and supports the significance of our result.

Another pathway represented in our gene-set is the P53 pathway. The p53 pathway activates cell-cycle arrest and apoptosis in response to a stress signal, and is known to have a tumor-suppressor phenotype [15]. DRAM1 (DNA damage regulated autophagy modulator 1) has been shown to regulate cell proliferation and is thought to have tumor-suppressing effects [21]. Indeed, in our cohort we found that high expression of DRAM1 is associated with improved survival. Figure 4-8 demonstrates this observation with a Kaplan-Meier plot showing the survival function for patients with above-median expression of DRAM1 and patients with below-median expression of DRAM1. The difference between the two groups was found to be statistically significant via a logrank test, with a p-value of 0.0097. Again, we include in figure 4-8 a subplot which shows the distribution of p-values computed from 1000 label permutations. We observe a uniform distribution, suggesting that the data is not biased and supporting the significance of our result.

4.2.3 Association with Immune Cell Signature

As with the lung adenocarcinoma gene-set, we used CIBERSORT [25] to estimate the immune-cell abundance of samples in our cohort, which is done by deconvolving the bulk gene-expression data of our samples with known immune-cell gene signatures [25]. We then computed the correlation between enrichment of our SKCM gene-set and immune-cell abundance in our cohort. Figure 4-9 demonstrates the analysis procedure and results.

Again, we find that enrichment of our gene-set correlates with higher abundance of T-cells in our samples, specifically activated memory CD4 T-cells. As discussed in the previous section, CD4 T-cells have been shown to have an anti-tumor function

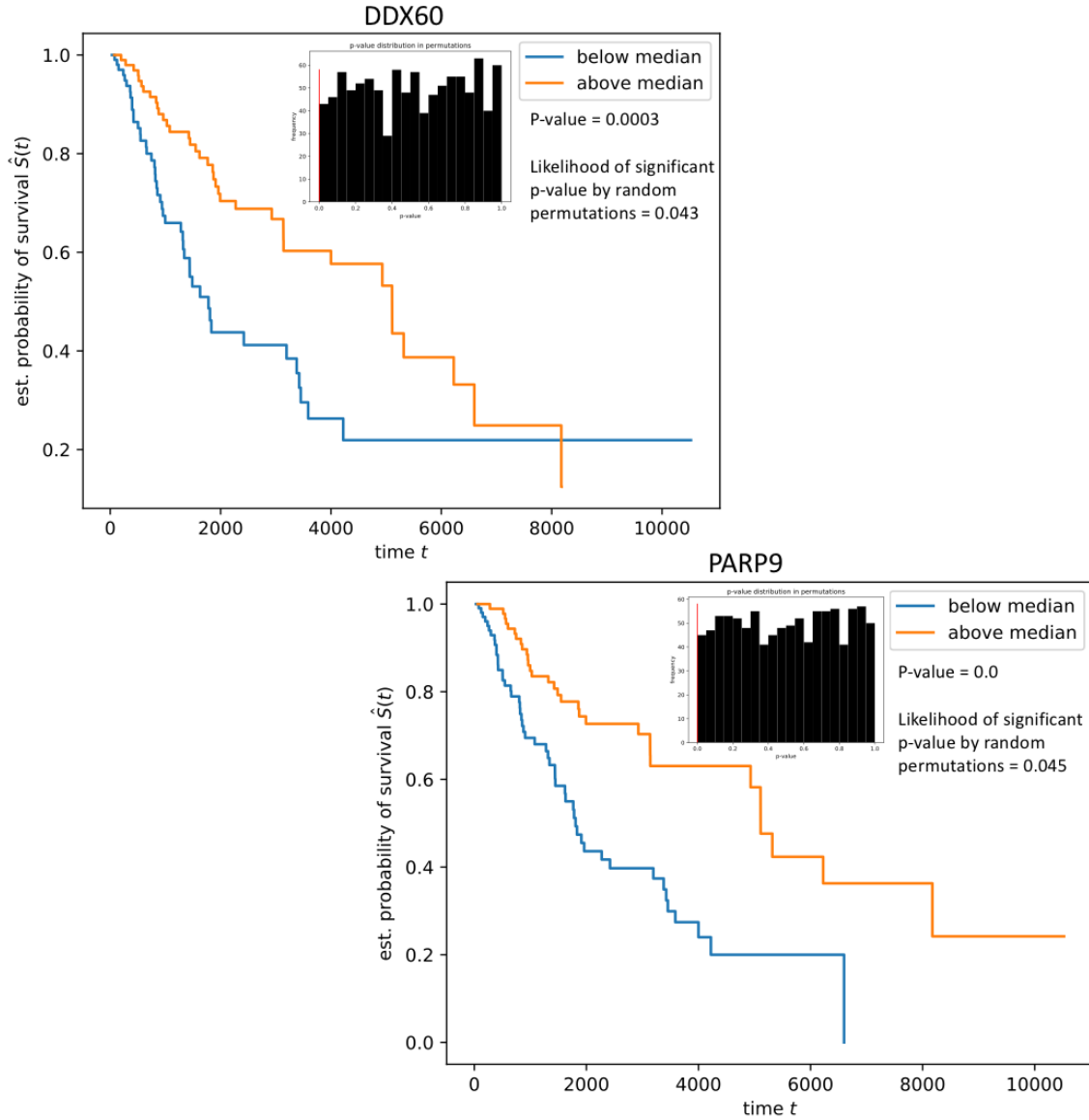


Figure 4-7: Kaplan-Meier plots comparing high vs. low expression of DDX60 and PARP9. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations.

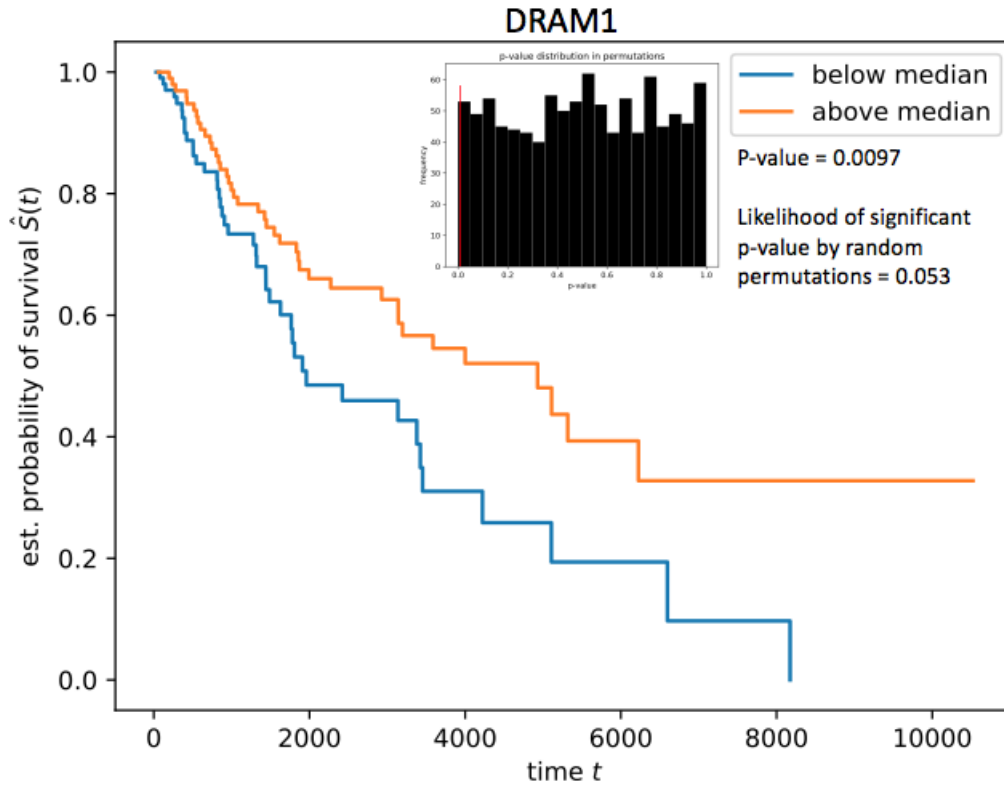


Figure 4-8: Kaplan-Meier plots comparing high vs. low expression of DRAM1. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations.

and can serve as prognostic biomarkers or therapeutic targets. [10]. Our gene-set's correlation with CD4 T-cell signatures demonstrates the gene-set's potential as an SKCM biomarker.

We also find a correlation with abundance of M1 macrophages. While macrophages can have either pro-tumor and anti-tumor effects when present in the tumor micro-environment, M1 macrophages are believed to have anti-tumor activity as discussed in the previous section [2]. Again, correlation between enrichment of our gene-set and M1 macrophages abundance in the tumor demonstrates the gene-set's potential as a prognostic biomarker.

4.2.4 Differentiating Immunotherapy Responders and Non-Responders

Riaz et al. [31] evaluated the genomic profile of 68 melanoma patients undergoing immunotherapy. They found that patients' pre-treatment samples had differentially expressed genes between responders and non-responders, most of which were immune-related genes, supporting our claim that high immune-cell abundance in the tumor micro-environment can predict immunotherapy response. We compared the genes they found to be significantly differentially expressed between responders and non-responders with the genes in our gene-set, and found 5 over-lapping genes: DRAM1, AFTPH, PARP11, LSR, and DTX3L. This overlap was determined to be statistically significant with a p-value of 0.0156 using a hypergeometric distribution. As discussed previously, DRAM1 has been shown to regular cell-proliferation and to have tumor-suppressing effects [21], and figure 4-8 shows that the survival function of patients in our cohort with high expression of DRAM1 is more favorable. Futhermore, figure 4-10 shows that above-median expression of each of PARP11 and AFTPH is associated with a more favorable survival function in our cohort, and that those differences are statistically significant. Statistical significance was determined using the logrank test. The figures also show that the p-values obtained from randomly permuting the data follow a uniform distribution for both genes. These findings further support our claim

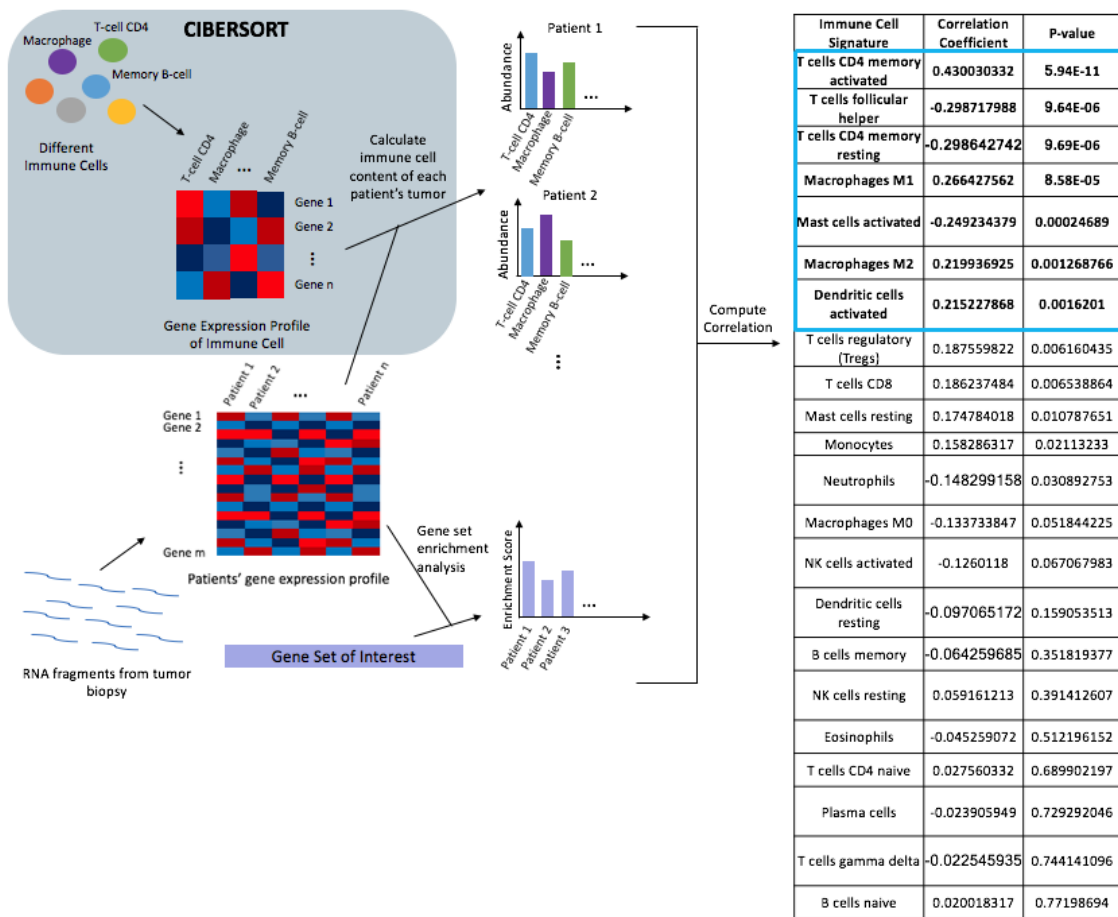


Figure 4-9: Analysis of SKCM gene-set correlation with estimated immune cell abundance in the tumor environment

that the genes we identified have the potential of predicting immunotherapy response potential of SKCM patients.

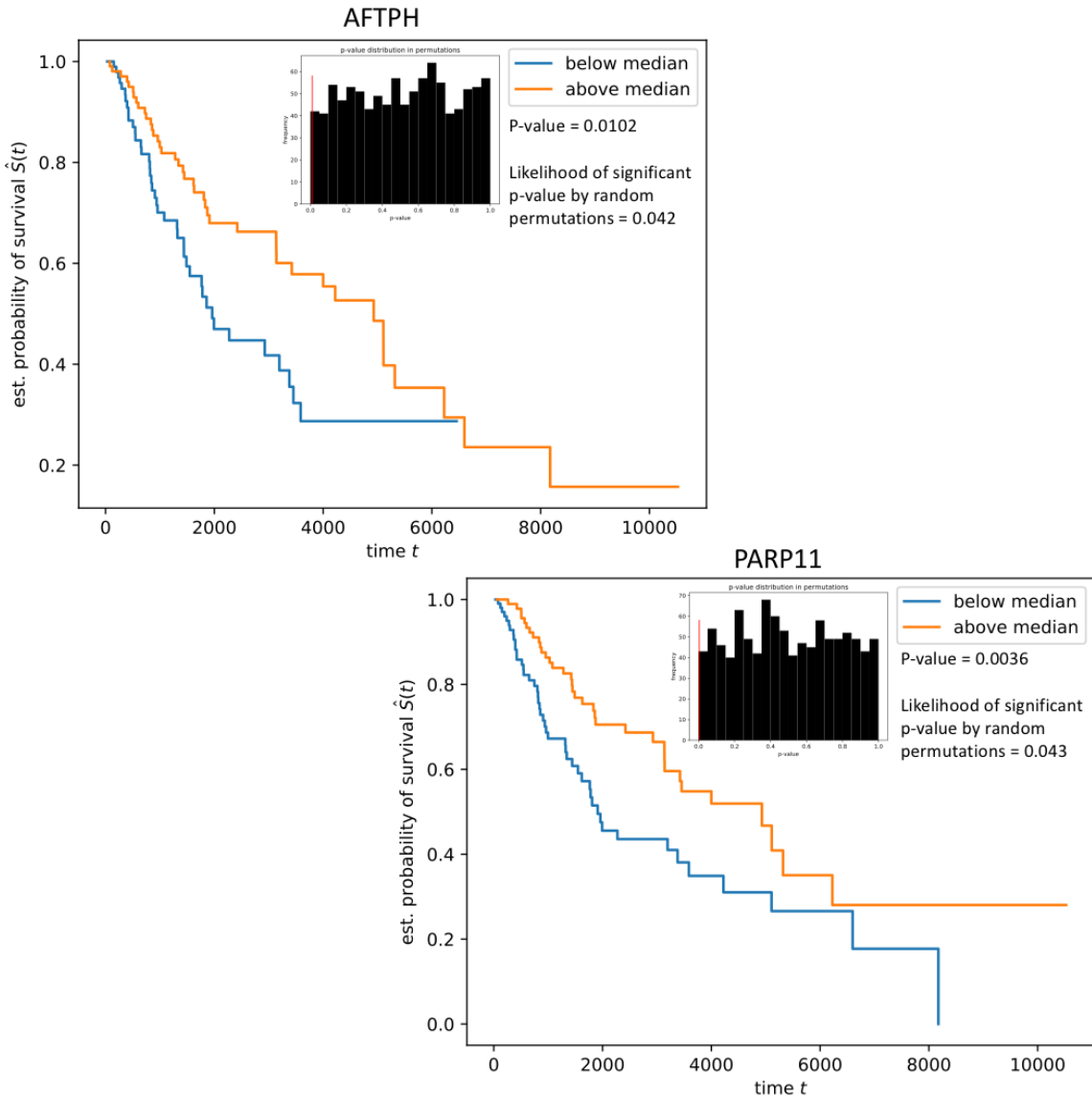


Figure 4-10: Kaplan-Meier plots comparing high vs. low expression of AFTPH and PARAP11. The subplot in the top center shows the p-value distribution obtained from randomly permuting survival labels 1000 times and comparing the difference between the survival function associated with high expression and the survival function associated with low expression under the random permutations. The red line on the histogram indicates the p-value of that difference in the original data. The uniform histogram serves to validate the significance of the difference observed between the two original curves. The figure includes the p-value of the difference between the two survival function, and the proportion of significant p-values observed in the cross-validation permutations.

Chapter 5

Discussion

5.1 Motivation and Contributions

Cancer is a leading cause of death world-wide. A large body of cancer research aims to identify biomarkers that correlate with patient prognosis and response to treatment [6]. The advancement of next-generation sequencing technologies has led to the availability of large amounts of whole-genomic and whole-transcriptomic tumor data-sets that can be used for such research. Given the heterogeneity of cancer, personalized treatment decisions enabled by such research has the potential to improve patient outcomes.

This project utilizes whole-transcriptome data from the Cancer Genome Atlas to identify biomarkers in lung adenocarcinoma and skin cutaneous melanoma. We implemented and evaluated various statistical and machine learning techniques for the identification of such biomarkers, and present here a work-flow that utilizes those techniques. We used this work-flow for the unbiased identification of gene-sets that can act as prognostic biomarkers for LUAD and metastatic SKCM. We note that the work-flow can be adjusted for use in other cancer types, as well as other diseases, due to the fact that no disease-specific knowledge was used in the work-flow implementation. For both disease types we showed that expression levels of genes in our gene-set correlates with survival and we found various cancer hallmarks that were represented in our gene-sets demonstrating the genes' potential as useful clinical biomarkers. Fi-

nally, we found that enrichment of our gene-set in the tumor transcriptome correlated with immune cell abundance, specifically CD4 and CD8 T-cells, and macrophages. These results further support our belief that the gene-set can serve as a prognostic biomarker and inform the extent of the patient’s anti-tumor immune response and potential response to immunotherapy. We were also able to validate the utility of some of the genes in our SKCM gene-set in differentiating responders and non-responders to immunotherapy using data from a separate study on immunotherapy in melanoma patients by Riaz et al.

5.2 Limitations and Future Work

We would like to discuss some limitations of our project that can be addressed in future studies. One of the main limitations of our work lies in the sparsity of the available clinical data. Since TCGA combines data from different institutions and different studies, we could not reliably obtain sufficient data pertaining to treatments, response status, and disease-free survival. As such, we relied on overall survival as a proxy for prognosis and treatment-response potential. This can be limiting as many confounding factors, such as other medical conditions, or death by other causes, can affect patients’ overall survival. Furthermore, patients may have received and responded to therapy, changing their tumor profile and extending their survival time. Our procedure did not account for such events. To address these limitations, we recommend the following future steps:

- 1) Controlling for patients with known actionable mutations in the data-set. These patients are likely to have received targeted therapy, giving them a survival benefit over patients with no such known mutations.

- 2) Validating the disease-free survival predictive potential of the identified genes in data-sets where such information is reliably available.

- 3) Further validating the immunotherapy-response predictive potential of the identified genes in data-sets with patients who have undergone immunotherapy.

In addition to the above, we also recommend including additional DNA focused

data modalities in the feature selection. Mutational information, for example, can be very valuable in characterizing the tumor as it would provide insight regarding potential functional alterations in the translated proteins. Also, copy-number variation and DNA methylation data can allow us to indirectly infer up-regulation or down-regulation of transcription. Inferring this information using DNA-based biomarkers will allow the measurement of the identified biomarkers in liquid biopsies using cell-free DNA data. Cell-free DNA can be obtained from a blood sample, providing a minimally invasive procedure to characterize a patient's disease. This would allow for disease characterization in patients where a tissue biopsy is not available, and it would also allow for the collection of longitudinal data and tracking of patients' tumor progression by sampling at multiple time-points.

Bibliography

- [1] The double edge of cancer immunotherapy, February 2017.
- [2] Shweta Aras and M Raza Zaidi. Tameless traitors: macrophages in cancer progression and metastasis. *British Journal Of Cancer*, 117:1583+, October 2017.
- [3] Colin N Dewey Bo Li. Rsem: accurate transcript quantification from rna-seq data with or without a reference genome. *BMC bioinformatics*, 12:323, August 2011.
- [4] Adrian P. Bracken, Marco Ciro, Andrea Cocito, and Kristian Helin. E2f target genes: unraveling the biology. *Trends in Biochemical Sciences*, 29(8):409–417, 2004.
- [5] Flavia Castro, Ana Patricia Cardoso, Raquel Madeira Goncalves, Karine Serre, and Maria Jose Oliveira. Interferon-gamma at the crossroads of tumor immune surveillance or evasion. *Frontiers in Immunology*, 9:847+, May 2018.
- [6] Eytan Ruppin Ze’ev A. Ronai Daniela Senft, Mark D.M. Leiserson. Precision oncology: The road ahead. *Cell*, 23:874–898, September 2017.
- [7] Jose Baselga David M. Hyman, Barry S. Taylor. Implementing genome-driven oncology. *Cell*, 168:584–599, February 2017.
- [8] Cameron Davidson-Pilon, Jonas Kalderstam, Paul Zivich, Ben Kuhn, Andrew Fiore-Gartland, Luis Moneda, Gabriel, Daniel Wllson, Alex Parij, Kyle Stark, Steven Anton, Lilian Besson, Jona, Harsh Gadgil, Dave Golland, Sean Hussey, Ravin Kumar, Javad Noorbakhsh, Andreas Klintberg, Jakub Kaluzka, Isaac Slavitt, Eric Martin, Eduardo Ochoa, Dylan Albrecht, dhuynh, Denis Zgonjanin, Daniel Chen, Chris Fournier, Arturo, and Andr   F. Rendeiro. Camdavidson-pilon/lifelines: v0.21.1, April 2019.
- [9] C. Ding and H. Peng. Minimum redundancy feature selection from microarray gene expression data. In *Computational Systems Bioinformatics. CSB2003. Proceedings of the 2003 IEEE Bioinformatics Conference. CSB2003*, pages 523–528, Aug 2003.
- [10] Aurelie Durgeau, Yasemin Virk, Stephanie Corgnac, and Fathia Mami-Chouaib. Recent advances in targeting cd8 t-cell immunity for more effective cancer immunotherapy. *Frontiers in immunology*, 9:14+, January 2018.

- [11] Hawazin Faruki, Gregory M. Mayhew, Jonathan S. Serody, D. Neil Hayes, Charles M. Perou, and Myla Lai-Goldman. Lung adenocarcinoma and squamous cell carcinoma gene expression subtypes demonstrate significant differences in tumor immune landscape. *Journal of Thoracic Oncology*, 12(6):943–953, 2017.
- [12] Geoffrey T Gibney, Louis M Weiner, and Michael B Atkins. Predictive biomarkers for checkpoint inhibitor-based immunotherapy. *The Lancet. Oncology*, 17(12):542+, December 2016.
- [13] Salvatore Siena Alberto Bardelli Giulia Siravegna, Silvia Marsoni. Integrating liquid biopsies into the management of cancer. *Nature Reviews Clinical Oncology*, 14:531–548, March 2017.
- [14] Douglas Hanahan and Robert A. Weinberg. Hallmarks of cancer: The next generation. *Cell*, 144(5):646–674, 2011.
- [15] Sandra L Harris and Arnold J Levine. The p53 pathway: positive and negative feedback loops. *Oncogene*, 24(17):2899–2908, 2005.
- [16] Roy S. Herbst, Daniel Morgensztern, and Chris Boshoff. The biology and management of non-small cell lung cancer. *Nature*, 553(7689):446–454, Jan 25 2018. Name - Food Drug Administration–FDA; Copyright - Copyright Nature Publishing Group Jan 25, 2018; Last updated - 2018-01-30; SubjectsTermNotLit-GenreText - United States–US.
- [17] Jurgita Jackute, Marius Zemaitis, Darius Pranys, Brigita Sitkauskiene, Skaidrius Miliauskas, Simona Vaitkiene, and Raimundas Sakalauskas. Distribution of m1 and m2 macrophages in tumor islets and stroma in relation to prognosis of non-small cell lung cancer. *BMC immunology*, 19(1):3+, January 2018.
- [18] Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source scientific tools for Python, 2001–.
- [19] Salome LeibundGut-Landmann, Jean-Marc Waldburger, Michal Krawczyk, Luc A. Otten, Tobias Suter, Adriano Fontana, Hans Acha-Orbea, and Walter Reith. Mini-review: Specificity and expression of ciita, the master regulator of mhc class ii genes. *European Journal of Immunology*, 34(6):1513–1525, 2004.
- [20] Giulia C Leonardi, Luca Falzone, Rossella Salemi, Antonino Zanghi, Demetrios A Spandidos, James A Mccubrey, Saverio Candido, and Massimo Libra. Cutaneous melanoma: From pathogenesis to therapy. *International journal of oncology*, 52(4):1071–1080, February 2018.
- [21] Ting Lu, Zhou Zhu, Junchao Wu, Hua She, Rong Han, Haidong Xu, and Zheng-Hong Qin. Dram1 regulates autophagy and cell proliferation via inhibition of the phosphoinositide 3-kinase-akt-mtor-ribosomal protein s6 pathway. *Cell Communication and Signaling*, 17(1):28+, 2019.

- [22] Guan-Yi Lyu, Yu-Hsuan Yeh, Yi-Chen Yeh, and Yu-Chao Wang. Mutation load estimation model as a predictor of the response to cancer immunotherapy. *npj Genomic Medicine*, 3(1):12+, April 2018.
- [23] Michael B Mathews and John W B Hershey. The translation factor eif5a and human cancer. *Biochimica et biophysica acta*, 1849(7):836–844, 2015.
- [24] Vamsi K Mootha, Cecilia M Lindgren, Karl-Fredrik Eriksson, Aravind Subramanian, Smita Sihag, Joseph Lehar, Pere Puigserver, Emma Carlsson, Martin Ridderstrale, Esa Laurila, Nicholas Houstis, Mark J Daly, Nick Patterson, Jill P Mesirov, Todd R Golub, Pablo Tamayo, Bruce Spiegelman, Eric S Lander, Joel N Hirschhorn, David Altshuler, and Leif C Groop. Pgc-1 α -responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes. *Nature Genetics*, 34(3):267–273, 2003.
- [25] Aaron M Newman, Chih Long Liu, Michael R Green, Andrew J Gentles, Weiguo Feng, Yue Xu, Chuong D Hoang, Maximilian Diehn, and Ash A Alizadeh. Robust enumeration of cell subsets from tissue expression profiles. *Nature Methods*, 12:453+, March 2015.
- [26] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [27] S Pölsterl, P Gupta, L Wang, S Conjeti, A Katouzian, and N Navab. Heterogeneous ensembles for predicting survival of metastatic, castrate-resistant prostate cancer patients [version 1; peer review: 3 approved with reservations]. *F1000Research*, 5(2676), 2016.
- [28] Sebastian Pölsterl, Nassir Navab, and Amin Katouzian. *Machine Learning and Knowledge Discovery in Databases*. Springer International Publishing, 2015.
- [29] Sebastian Pölsterl, Nassir Navab, and Amin Katouzian. An efficient training algorithm for kernel survival support vector machines. *CoRR*, abs/1611.07054, 2016.
- [30] C. Zhang A.V. Penson E. Reznik L. Zhang T. Minet-A. Ochoa B.E. Gross C. A. Iacobuzio-Donahue D. Betel B.S. Taylor J. Gao N. Schultz Q. Wang, J Armenia. Unifying cancer and normal rna sequencing data from different sources. *Scientific Data*, 5:180061, April 2018.
- [31] Nadeem Riaz, Jonathan J Havel, Vladimir Makarov, Alexis Desrichard, Walter J Urba, Jennifer S Sims, F Stephen Hodi, Salvador Martin-Algarra, Rajarsi Mandal, William H Sharfman, Shailender Bhatia, Wen-Jen Hwu, Thomas F Gajewski, Jr Slingluff, Craig L, Diego Chowell, Sviatoslav M Kendall, Han Chang, Rachna Shah, Fengshen Kuo, Luc G T Morris, John-William Sidhom, Jonathan P Schneck, Christine E Horak, Nils Weinhold, and Timothy A Chan.

Tumor and microenvironment evolution during immunotherapy with nivolumab. *Cell*, 171(4):934–949, November 2017.

- [32] Aravind Subramanian, Pablo Tamayo, Vamsi K. Mootha, Sayan Mukherjee, Benjamin L. Ebert, Michael A. Gillette, Amanda Paulovich, Scott L. Pomeroy, Todd R. Golub, Eric S. Lander, and Jill P. Mesirov. Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences*, 102(43):15545–15550, 2005.
- [33] Suzanne L. Topalian, F. Stephen Hodi, Julie R. Brahmer, Scott N. Gettinger, David C. Smith, David F. McDermott, John D. Powderly, Richard D. Carvajal, Jeffrey A. Sosman, Michael B. Atkins, Philip D. Leming, David R. Spigel, Scott J. Antonia, Leora Horn, Charles G. Drake, Drew M. Pardoll, Lieping Chen, William H. Sharfman, Robert A. Anders, Janis M. Taube, Tracee L. McMiller, Haiying Xu, Alan J. Korman, Maria Jure-Kunkel, Shruti Agrawal, Daniel McDonald, Georgia D. Kollia, Ashok Gupta, Jon M. Wigginton, and Mario Sznol. Safety, activity, and immune correlates of anti-pd-1 antibody in cancer. *New England Journal of Medicine*, 366(26), 2012.
- [34] Mingxin Wen, Yongwon Kwon, Yongsheng Wang, Jian-Hua Mao, and Guangwei Wei. Elevated expression of ube2t exhibits oncogenic properties in human prostate cancer. *Oncotarget*, 6(18):25226–25239, July 2015.
- [35] Lu Sun Chunying Song Blanca Homet Moreno Siwen Hu-Lieskovan Beata Berent-Maoz Jia Pang Bartosz Chmielowski Grace Cherry Elizabeth Seja Shirley Lomeli Xiangju Kong Mark C. Kelley Jeffrey A. Sosman Douglas B. Johnson Antoni Ribas Roger S. Lo1 Willy Hugo, Jesse M. Zaretsky. Genomic and transcriptomic features of response to anti-pd-1 therapy in metastatic melanoma. *Cell*, 165(1):35–44, March 2017.
- [36] Jin Zhang, Stefan Endres, and Sebastian Kobold. Enhancing tumor t cell infiltration to enable cancer immunotherapy. *Immunotherapy*, 11(3):201–213, 2019. PMID: 30730277.