



Constraint Model for libwww Webbot

[Manolis Kamvysselis](#) - [Henryk Frystyk Nielsen](#)
[World Wide Web Consortium](#) - Summer 96

Jump to:

- [Design](#) issues and purpose
- [Rules](#) description and issuing of commands
- [URLs](#) and Web Elements
- [Function](#) description for extending the robot
- [Extensions](#) to the robot, and what's next

I. Design issues and Purpose of the Robot

The W3C robot's purpose in crawling the web is not to accumulate data from web sites, to generate a database, or an index, or to update documents, or verify links, even though all of the functionality described above can easily be added, without changing the basic skeleton of the robot as it is today.

The primary purpose of the robot is to understand the relations between web pages, and not the content of the pages themselves. How many paths can lead to a given page, what tree structures can be generated by following links under a certain constraint, how many degrees of separation lie between any two pages, what is the shortest path to reach a site, and how can intelligent navigation be automated.

The robot will also be used as a test tool interacting with the w3c-libwww-99. The robot code is written in Tcl, the graphic interface to the user is written in Tk, and the library code used to provide all the network functionality is written in C. All functionality is therefore portable in almost any system with network access.

The most elementary web robot can fetch a web document, extract the links from that document, follow the urls, fetch the documents, and restart for every document. This leads to an exponentially growing number of documents to follow every time you reach a new depth. Therefore, for a robot to be of any use, we need a way of specifying a **Constraint Model** that will allow the robot to choose which documents to follow and which not to.

The Constraint Model is basically a set of rules, that the user specifies, or that the program itself can formulate, to follow a crawling model specified by the user at the beginning, or during the crawling. The next section describes the constituent blocks of the rules (or Rule Atoms) and the ways to combine those rules.

II. Rule Description

A rule can be simple or complex. Complex rules can be expressed as logical combinations of simpler rules, called Rule Atoms.

A. Rule Atoms

Definition

A Rule Atom is the simplest form of a rule, and the building block of every rule. Anything simpler can not be called a rule, and anything more complicated cannot be called a rule atom.

Examples

```
+url:*.mit.edu/*
  keep documents whose url matches *.mit.edu/*
  matches: http://web.mit.edu/manoli/www/ and http://manoland.mit.edu/ but not http://web.mit.edu.ca/
-body:*badword*
  reject documents whose body contains badword
```

```
+head:*meta*
  keep documents whose head contains the meta keyword
-date:*december*
  reject documents whose date contains the word december
```

Understanding the Rule Atoms

As you can see from those examples, each **rule atom** has three elements:

- a. The **category** of the rule can be one of the following four, in order of priority:
 1. **Imperative rules (*)**: those documents will be followed first and always
 - Direction rules (+)**: Keep: the documents will be followed next if they don't match -
 - Rejection rules (-)**: Reject: document will not be followed, unless it matches an imperative rule
 - LastResort rules (?)**: Default: if the robot has no more documents matching * or +, it follows ?
- b. The **field** of the document to search
 1. Sample fields are: url, head, body, date, etc...
The program calls `Url_Geturl`, `Url_Getbody`, etc, depending on the field specification
To add a new field, all you need to do is provide a function
- c. The **pattern** to match
 1. For the moment, it is a glob-type pattern: * matches any combination of characters, ? matches one character. Later on, functionality might be added, to handle regular expression matching.
Please Note: If you want to match the pattern in the middle of a string, you have to add an initial and final * Here are some examples, that you should take a look at:
 - *.mit.edu does not match `http://web.mit.edu/` since the final * is not present
 - *.mit.edu/* matches `http://web.mit.edu/people/`, `ftp://www.mit.edu/pub/`
 - *.edu* matches `http://www.sfu.edu.co/`, `ftp://www.mit.edu.ca/pub/bin/`

B. Combining Rules to form Rules

Definition of a Rule

The definition is recursive: a rule can be either a rule atom, or a logical combination of rules.

Undersating Rule Combination

Here are the forms a rule can take and the result of applying the rule to a list of documents

Applying ... to a list of documents	Returns the list of documents that...
RuleAtom	match an imperative rule, or that match a directional rule and don't match a rejection rule
{OR rule1 rule2 ...}	either returned by applying rule1, or returned by applying rule2, or ..
{AND rule1 rule2 ...}	urls returned by rule1, and by rule2, and by...
{NOT rule1 rule2}	urls returned by rule1 and not by rule2.

Example

With those simple rules, we can express very precise constraints

```
{NOT {AND +url:*.mit.edu/* -url:*.mit.edu/personal/*
      {OR -body:*badword* +head:*goodcontext*}}
 +date:*october*}
```

C. Applying Rules

The notion of truth

Applying a rule such as `+url:*.mit.edu/*` to a list of elements

more...

Please note that the logical NOT operator is not defined in its common interpretation of a unary operator, negating something that is false, to give something true. This is because testing a rule against a set of web elements, does not return true or false, but the list notion of true and false, which is sublists. Therefore, negating a list list2 cannot return the contrary of list2, unless we know what represents the "truth" which in this case is list1. Therefore, one must understand that we can only negate a list relative to a larger list. {NOT 1} is always 0 and {NOT 0} is always 1, but the result of evaluating a rule against a list of documents (that we will assimilate to URLs in this paragraph) is

not a boolean, but a sublist. Therefore, if `http://web.mit.edu/` is the result of applying rule2, the `{NOT rule2}` will be `{NOT http://web.mit.edu/}`, which would represent the entire web, except `http://web.mit.edu/`. One must therefore provide the part of the web against which `http://web.mit.edu/` will be negated, and that in form of a list.

Applying a rule to an Element list

- a. Reminder
 1. As you saw earlier, a rule atom has the form `+url:*.mit.edu/*`, and is formed of the category specification (`*`, `+`, `-`, `?`), the field to check (`url`, `head`, `body`, `date`, etc) and the pattern to match (`*.mit.edu/*` in this example).
- b. Applying a Rule Atom to a Web Element
 1. Depending on the field specification of the rule, a corresponding `Get$field` function is called on the Web Element (where `$field` can be any of `url`, `head`, `body`, `date`, etc...), and the pattern is matched against that result. To be able to match against more fields, all one has to do is provide the code for getting the corresponding field from the Web Element, and the rule will be used with no difference. For example, to be able to use a rule `+hat:top` or `+hat-color:red`, all one has to do is provide the functions `URL_Gethat` and `URL_Gethat-color`. The name can be anything, but cannot contain a column (`:`), since it is used to separate the field from the pattern in a rule.
- c. Applying a Rule Atom to a List of Web Elements
 1. When the Web Atom is applied to a List of Elements, two local lists are created, named `matched`, and `unmatched`, and each Element checked goes to one list or the other. After every element in the list has been checked, the result (the return value) of applying the Rule to the Element List is the list of matched elements in the case of an imperative (`*`) or directional (`+`) rule atom, or the list of unmatched elements, in the case of a negational (`-`) rule.
- d. Applying Rule Combinations to Element Lists
 1. When applying `{AND rule1 rule2 {OR rule3 rule4}}` to a list of elements, `rule1`, `rule2`, `rule3`, and `rule4` are applied separately to the list of elements, yielding `list1`, `list2`, `list3`, `list4`. Then is evaluated the expression `[List_and list1 list2 [List_or list3 list4]]`, which returns the desired list value, the result of applying the rule combination to the element list.
- e. Accepting or rejecting an element

D. Sets of Rules

Global Rules and Local Rules

- a. Design issues and purpose
 - o The Constraint Model of the Robot was built on the idea that the links a user will choose to follow, do not only depend on the purpose of his navigation and his choices of sites, but also, on the history of his navigation, the documents accessed previous to the current page. Therefore, a constraint model would be more efficient by having not only a set of static, global rules, that are set by the user when the navigation begins, but also a set of local rules for each page, that are set automatically, as the navigation progresses, and that are specific to a page, or a family of pages.
- b. Global Rules and Constraints
 - o Global rules are set by the user before the navigation begins. They represent the rules that the program has to follow throughout his navigation, and cannot be overwritten by local rules. For example, if a user doesn't want to access personal pages, a global rule would be `-url:*/Personal/*`, or if he doesn't want to access only educational pages, a rule might be `+url:*.edu*`
- c. Local Rules and Navigation method
 - o The local rules are set every time a new element is created. They are set based on a global variable, `nav`, that the user sets at the beginning of the navigation, and that he can change at any time. This variable can now have four values: `dumb`, `deep`, `broad`, `strategic`.

III. Representing Documents and Urls

A. A Web Element

Definition

A Web Element is the sum of all the information gathered on a certain web document reached following a certain path, with a known set of rules.

Representation of a Web Element

A Web Element consists of:

Family

The documents reached from parent 0302 will have family indexes 030201, 030202, 030203. The family index is unique, and fully identifies a given Element.

Url

The same document can be reached from two different paths. The same url therefore can be shared by 030202 and 0401.

Priority

It shows how important a certain element is for the navigation of the robot. Here's a table of values:

Value	Category	Corresponding rule	Description
0	Untested	none	Element has not been tested yet
1	First	imperative	Will be followed first
2	Next	directional	Followed next
3	Last	lastresort	Followed last
5	Rejected	rejection	Never followed

Rules

The local rules of the url. Updated automatically.

State

Has the url been followed yet. 1 for yes, 0 for no.

B. The Web

As was said earlier, the main interest of the robot, is not how the web is constructed, how a site is managed, or how the web should be ordered, but the relations that exist among documents,

IV. Program Functions

A. Rule Functions

function	input	outputdescription
Rule_Init		Initializes rules
Rule_ShowAll		
Rule_ApplyAtom	elementlist ruleatom	list of matching elements
Rule_Apply	elementlist rule	
Rule_and		
Rule_or		

B. Handling Urls

function	input	outputdescription
Url_Init		Initializes urls
Url_ShowAll		
Url_ApplyAtom	elementlist ruleatom	list of matching elements
Url_Apply	elementlist rule	

C. Handling Lists

function	input	outputdescription
List_and	list1 list2 ...	list of common elements to list1 list2 ...
List_or	list1 list2 ...	list of elements either in list1 or list2 ...
List_not	list1 list2	list of elements in list1 and not in list2

V. Extending the Robot

Gathering more information

What is really easy to do with the current design of the robot, is to make it gather more information from the documents that it accesses. All one has to do is change the `Url_Load` function.

Graphical Representation

Once elements are gathered by the robot, their family indexes describe their relations very precisely. One can then choose a representation for trees, and provide a graphical representation of the portion of the web accessed.

Distributing Web indexing and crawling

When a stable design of storing results is reached and proven general and extensible, then different parties can run their robot, and create tree structures inside their own sites, and store these results in an agreed-upon location, that will be searched by the robot, when accessing that site, and the results will be reused.

Strategic navigation

The present robot doesn't know how to handle strategic navigation, which would be trying to reach a url. A special way of calculating rules will have to be formulated.

More

And much much more...

[Manolis Kamvysselis](mailto:manoli@mit.edu) - manoli@mit.edu - <http://mit.edu/manoli>

Summer 96