

# EXTRACTION OF TEMPLATE FROM DIFFERENT WEB PAGES

**Thota Srikeerthi<sup>1\*</sup>, Ch. Srinivasarao<sup>2\*</sup>, Vennakula I s Saikumar<sup>3\*</sup>**

1. M.Tech (CSE) Student, Dept of CSE, Pydah College of Engg & Tech, Vishakapatnam.

2. Assoc Professor, Dept of CSE, Pydah College of Engg & Tech, Vishakapatnam.

3. M.Tech (CSE) Student, Dept of CSE, Pydah College of Engg & Tech, Vishakapatnam.

## Keywords:

Surface Template,  
XSLT, Wrapper  
Development

**Abstract:** Information in these days becomes the most wanted in the era of globalization. People for the use mankind make information as the most complex component in their life. According to the need technology changes from time to time, especially considering to Computer Science, it's a changing technology. These days compared to early ages of Computer Science there is a lot of changes occurred. Considering the 'WWW' in other terms the internet age, we come across a terminology like web site which we call a page as web page. In order to attract the user to we use template for the sake of ease and a layman can understand. Web page consists of some framework or template likely to look like as of business requirement; may be of scientific or banking etc. In order to achieve high productivity in terms Interface which we likely to tell as the end user page and designed in such way that a layman can use it easily and effectively. Hence of this paper tries to concentrate on the concept of extraction of template from various web pages and tries to populate a general or common template which can be easy to end user. Keeping in mind there are many research is going on how effective is user interface and the point come into existence that is it useful to satisfy the accuracy and performance that depends upon intelligence how effective we use it which uses the clustering mechanism to effort the web page effectively. Hence we likely to concentrate on the clustered mechanism concept where a large volume of template having high positivity towards the usability.

## 1. INTRODUCTION

In the information technology world, in order to publish information we need the help of www, in technical side we need a page which likely to written in HTML, XML or some other language by using some protocol to govern it. Hyper Text Markup Language is not designed for structured data extraction at the first place, but mainly for data presentation. HTML pages are usually 'dirty', may have improper closed tags, improper nested tags, and incorrect parameter value of a tag are examples of such problems. In order to reduce parsing error over ill formed web pages caused by loose HTML standard, World Wide Web Consortium recommend stricter standard Markup

**\* Thota Srikeerthi**

M.Tech (CSE) Student, Dept of CSE, Pydah College of Engg & Tech, Vishakapatnam

Languages, such as XHTML and XML. However, parser of search engine still needs to cope with existing web pages that do not conform to the new standards.

Text extraction from HTML page is a critical preprocess for information retrieval. In this phase, page content must be extracted and irrelevant template data removed. Hyper Text Markup Language is not designed for structured data extraction at the first place, but mainly for data presentation. HTML pages are usually "dirty", i.e., improper closed tags, improper nested tags, and incorrect parameter value of a tag are examples of such problems. In order to reduce parsing error over ill formed web pages caused by loose HTML standard, World Wide Web Consortium recommend stricter standard Markup Languages, such as XHTML and XML. However, parser of

search engine still needs to cope with existing web pages that do not conform to the new standards.

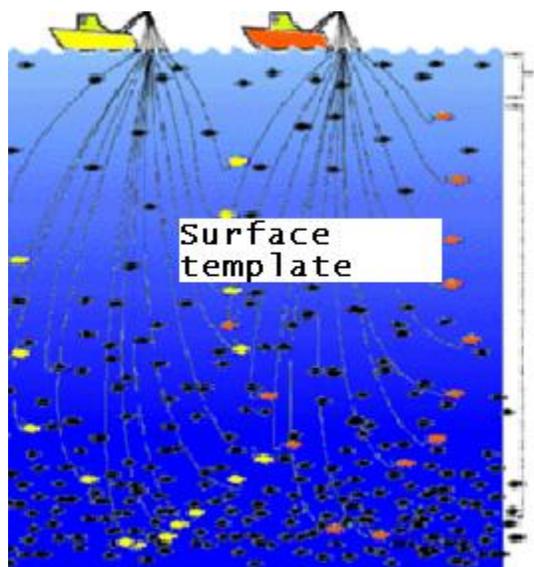


Fig 1.1 Showing Surface template to the ocean of web pages.

As shown next in figure 1, we must extract unstructured data from the web page; meanwhile remove irrelevant template navigation bars and advertisements. Based on our observations of html page, we propose a novel approach to achieve this goal. We do not retrieve the multimedia content at the current stage, thus no advertisements are included in extracted text.

## 2. RELATED WORK

Local algorithms based on machine learning have been proposed to remove certain types of template material. It uses decision tree learning to remove “nepotistic” links, which are not present for a valid navigational purpose, and Kushmerick introduces AdEater, a browsing assistant that learns to automatically remove banner advertisements from pages. There are various approaches to understanding the relative merits of different parts of web pages that address problems raised by the presence of templates and related phenomena. Kao et al. Propose a scheme based on information entropy to focus on the links and pages that are most information-rich, hopefully downgrading template material in the process. carefully decompose web pages using features of the layout into

blocks, and then judge the quality and salience of the blocks in order to rate their importance.

The Web pages consist of template which may in turn differ from context to context, Hence of keeping in mind such concepts that We briefly describe how we have organized the experimental results. For each input collection of web pages that we used we present the following information as part of the experimental results.

```
<template schema="1">
  <start-string context="2">
    <![CDATA[<html> <body> Book:]]>
  </start-string>
  <start-string context="5">
    <![CDATA[Author:]]>
  </start-string>
  <end-string context="2">
    <![CDATA[</body> </html>]]>
  </end-string>
</template>
```

The encoding of the value above using the template above results in the following page:

```
<html>
  <body>
    Book: Computer Science
    Author: Mr AP
    Author: Mr Wolber
  </body>
</html>
```

A template is just set of optional start-string and end-strings associated with each type in the schema. The context attribute in the <start-string> and <end-string> elements identifies the type in the schema that the element is associated with. In an encoded page, the "start-string" occurs before the encoding a sub-value of the type that it is associated with, and the "end-string" after. The above representation of the template is equivalent to our definition of a template in the paper.

In our research on Web data extraction, we emphasize the importance of middleware solutions that extract entire databases from target Web sites and make these datasets available for data mining and other analysis (similar to the Junglee system). This involves crawling target Web sites periodically, extracting structured data, and performing domain-specific feature extraction.

Hence, our work encompasses a middle ground between database systems, AI, and Web systems

research in general. Our ideas have been implemented in ANDES, a software framework that merges crawler technology with XML-based data extraction technology. ANDES is similar to other extraction systems in that it defines a wrapper for each Web site of interest. The underlying data extraction method uses XSLT, which combines templates, path expressions, and regular expressions into a concise package. Templates can be used to decompose the data extraction process hierarchically, as is done in. An XML path expression can traverse an HTML document recursively and express predicates (Weblog, context and delimiter patterns, and token features). Finally, regular expressions, which are an extension to XSLT, permit decomposition of plain-text fields (leaf nodes) of an HTML tree. Note that many other data extraction systems process documents linearly; consider the forward and backward token rules in STALKER, the head-left-right-tail delimiters in the HLRT wrapper class, and the prefix/ infix/postfix expressions in. In contrast, XML path expressions take full advantage of the tree structure of HTML documents, making it easy to visit ancestors, siblings, and children before and after the current position in the document. For instance, finding the *n*th top-level table element in a document is trivial using an XPath expression, but may be impossible using linear expressions due to the possibility of table elements containing an unknown number of nested tables.

### 3. METHOD

Large-scale learning of scripts and narrative schemas also captures template-like knowledge from unlabeled text. Scripts are sets of related event words and semantic roles learned by linking syntactic functions with coreferring arguments. While they learn interesting event structure, the structures are limited to frequent topics in a large corpus. We borrow ideas from this work as well, but our goal is to instead characterize a specific domain with limited data. Further, we are the first to apply to the concept of mining in extraction the relevant template.

The proposed method outlined comprises the encoding and tagged mechanism, where schema object describes the context of web pages which is most component of template. The use of templates as a Meta programming technique requires two distinct operations: a template must be defined, and a defined template must be instantiated. The template definition describes the generic form of the generated source code, and the instantiation causes a specific set of source code to be generated from the generic form in the template.

A document is labeled for a template if two different conditions are met: (1) it contains at least one trigger phrase, and (2) its average per-token conditional probability meets a strict threshold. Both conditions require a definition of the conditional probability of a template given a token. The conditional is defined as the token's importance relative to its uniqueness across all templates. This is not the usual conditional probability definition as

Kidnap Bomb Attack Arson Precision Trigger phrases are thus template-specific patterns that are highly indicative of that template. After identifying triggers, we use the above definition to score a document with a template. A document is labeled with a template if it contains at least one trigger, and its average word probability is greater than a parameter optimized on the training set. A document can be (and often is) labeled with multiple templates. We next extract entities into the template slots. Extraction occurs in the trigger sentences from the previous section. The extraction process is two-fold: 1. Extract all NPs that are arguments of patterns in the template's induced roles.

Kidnap Bomb Attack Arson Precision Trigger phrases are thus template-specific patterns that are highly indicative of that template. After identifying triggers, we use the above definition to score a document with a template. A document is labeled with a template if it contains at least one trigger, and its average word probability is greater than a parameter optimized on the training set. A document can be (and often is) labeled with multiple templates. We next extract entities into the template slots. Extraction occurs in the trigger sentences from the previous section. The extraction process is two-fold: 1. Extract all NPs that are arguments of patterns in the template's induced roles.

2. Extract NPs whose heads are observed frequently with one of the roles

Algorithm for surface template extraction:

Use Template::ExtractIn;

Use Data::Dumperout;

My obj1 = Template::ExtractIN->new;

My \$template = << ' .';

  If(found)

    <ul>[% FOREACH record %]

      <li><A HREF="[% url %]">[% title %]</A>: [% rate %] – [% comment %].[% ... %]

  [% END %]</ul>

Else

```
My $document = << ';
```

```
<html><head><title>an link to
establish</title></head><body>
```

```
<ul><li><A HREF="xyz url">details.</A>: A+ - nice.
```

```
this text is ignored.</li>
```

```
<li><A HREF="some url to match">position.</A>: Z! -
```

```
Text to be ignored, too.</li></ul>
```

```
Print Data::Dumperout::Dumperout(obj-> extractIn
($template, $document))
```

In the above algorithm, Templates are different from macros. A macro, which is also a compile-time language feature, generates code in-line using text manipulation and substitution. Macro systems often have limited compile-time process flow abilities and usually lack awareness of the semantics and type system of their companion language (an exception should be made with Lisp's macros, which are written in Lisp itself and involve manipulation and substitution of Lisp code represented as data structures as opposed to text).

Template meta programs have no mutable variables—that is, no variable can change value once it has been initialized, therefore template meta programming can be seen as a form of functional programming.

URL and web page for commercial sites backed by dynamic commodity database are usually generated dynamically. Corresponding to different category of products, slightly different template is applied. Thus there can be more than one template from a web site. Assuming pages that look similar belong to the same template; we cluster html pages based on Crezenti et al's algorithm, and then identify elements of web pages that were generated by a common template. Crezenti et al analyzed tag and link property of html pages. Features such as distance from the home page, tag periodicity, URL similarity, and tag probability were applied as measures of HTML page similarity. Their method produces high accuracy in clustering pages generated from the same template. Clustering and mining as of both are dynamic template extraction mechanism is the further technology and research paper. Hence In this paper we try to give to solution of template of surfacing mechanism.

#### 4. CONCLUSION AND FUTURE ENHANCEMENT

Time and technology change with respect to human need. As of technology is concerned, the trend goes on internet world of which websites in turn gives rise to concept of web pages. It also gives rise to Many HTML web pages incorporate dynamic technique such as JavaScript and PHP script. As template mining is a typical issue to be under similar category of template. For instance, sub menu of navigation bars are rendered by "document. write" method of JavaScript on the event of mouse over. Image map are also implemented with JavaScript code. We certainly don't want to include any template navigational bar in the indexing phase. However, there are cases when page content are rendered with JavaScript methods. For example, the link the bottom of google home page "Make google your home page" appears only if you use other site as your home for Internet Explorer. In our future work to avoid the risk of missing critical information we need to handle those scripts.

As of that this research paper will give dynamic solution to the web world having high density of data,

#### REFERENCES

1. Jussi Myllymaki, Effective Web Data Extraction with Standard XML Technologies, World Wide Web Conference (WWW10), 2001.
2. W3C HyperText Markup language Home page, <http://www.w3.org/MarkUp/>
3. Valter Crescenzi, Giansalvatore Mecca, Paolo Merialdo, RoadRunner: Automatic Data Extraction from DataIntensive Web Sites, International Conference on Management of Data and Symposium on Principles of Database Systems (SIGMOD02), 2002.
4. Document Object Model (dom) Level 1 Specification Version 1.0, <http://www.w3.org/TR/REC-DOM-Level-1>, 2010.
5. <http://www.wikipedia.com>
6. [patent/www.google.com](http://patent/www.google.com)
7. <http://infolab.stanford.edu/~arvind/extract>

8. J. Cho and U. Schonfeld, "Rankmass Crawler: A Crawler with High Personalized Pagerank Coverage Guarantee," Proc. Int'l Conf. Very Large Data Bases (VLDB), 2007.
9. T.M. Cover and J.A. Thomas, Elements of Information Theory. Wiley Interscience, 1991.
10. V. Crescenzi, G. Mecca, and P. Merialdo, "Roadrunner: Towards Automatic Data Extraction from Large Web Sites," Proc. 27th Int'l Conf. Very Large Data Bases (VLDB), 2001.
11. V. Crescenzi, P. Merialdo, and P. Missier, "Clustering Web Pages Based on Their Structure," Data and Knowledge Eng., vol. 54, pp. 279- 299, 2005.
12. M. de Castro Reis, P.B. Golgher, A.S. da Silva, and A.H.F. Laender, "Automatic Web News Extraction Using Tree Edit Distance," Proc. 13th Int'l Conf. World Wide Web (WWW), 2004.
13. I.S. Dhillon, S. Mallela, and D.S. Modha, "Information-Theoretic Co-Clustering," Proc. ACM SIGKDD, 2003.
14. M.N. Garofalakis, A. Gionis, R. Rastogi, S. Seshadri, and K. Shim, "Xtract: A System for Extracting Document Type Descriptors from Xml Documents," Proc. ACM SIGMOD, 2000.
15. D. Gibson, K. Punera, and A. Tomkins, "The Volume and Evolution of Web Page Templates," Proc. 14th Int'l Conf. World Wide Web (WWW), 2005.
16. K. Lerman, L. Getoor, S. Minton, and C. Knoblock, "Using the Structure of Web Sites for Automatic Segmentation of Tables," Proc. ACM SIGMOD, 2004.
17. B. Long, Z. Zhang, and P.S. Yu, "Co-Clustering by Block Value Decomposition," Proc. ACM SIGKDD, 2005.
18. F. Pan, X. Zhang, and W. Wang, "Crd: Fast Co-Clustering on Large Data Sets Utilizing Sampling-Based Matrix Decomposition," Proc. ACM SIGMOD, 2008.
19. M.D. Plumbley, "Clustering of Sparse Binary Data Using a Minimum Description Length Approach," <http://www.elec.qmul.ac.uk/staffinfo/markp/>, 2002.
20. J. Rissanen, "Modeling by Shortest Data Description," Automatica, vol. 14, pp. 465-471, 1978.
21. J. Rissanen, Stochastic Complexity in Statistical Inquiry. World Scientific, 1989.