

International Journal of Computer Science and Mobile Computing

A Monthly Journal of Computer Science and Information Technology



ISSN 2320-088X
IMPACT FACTOR: 6.017

IJCSMC, Vol. 6, Issue. 3, March 2017, pg.14 – 19

INFORMATION EXTRACTION FROM WEB PAGES USING PATTERN MATCHING

Mrs. Gitanjali Shirsath **Dr. Rekha Rathore**

D.C. E., RKDF, Indore

D.C. E., RKDF, Indore

git08shirsath@gmail.com

rekharathore23@gmail.com

Abstract - Web content nowadays is designed in such a way that it becomes arduous for software tools to access them readily. Even web content that is automatically generated from back-end databases are usually presented with lot of extra add-on styles. Information Extraction from Web Pages Using Pattern Matching involves direct extraction of data from various web pages, where they mostly formed in an unstructured HTML format. In this paper we review to develop such systems which will automatically extract the user required data and deliver it to user through various means. In this paper, we present an automated information extraction system that can extract the relevant data from product descriptions across different sites. Extracting structured data from deep Web pages is a challenging problem due to the underlying entangled structures of such pages. Problems arise in querying data sources due to unstructured contents of web pages (HTML) we cannot directly extract data into a new structured form. To address this problem, we propose a system which will automatically extract the user required data and deliver it to user through various means (Email, Sms, and Webpage). We provide a fully visual and interactive user interface with new technique and approach.

Keywords— *Pattern Matching, Web Data Extraction, HTML, Relevant data*

I. Introduction

Information Extraction from WebPages system is a system which extracts data from web page and delivers it to user through various means. The system will repeatedly extract the data through web pages using Template Matching. Data on the deep web pages are usually stored in the form on templates. Extraction in Proposed system is concerned on data which is stored in a form of template. The Report provides complete description of all the function and specification of Project Information Extraction from Web Pages System. Proposed system will automatically and repeatedly extracts data from web pages with changing content and presents the extracted data in various forms, such as e-mail, sms updates, and as a web page.

World Wide Web provides huge amount of data which is a good source of valuable information. This information is obviously very useful for a variety of purposes. Thus we need to automate the translation of web pages into structured data and make software tools gain database system functionality over this structured data. This enables easier comparison between products from

different online stores and can support a variety of advanced applications such as product recommender, online false advertisement detection, and demand forecasting systems. The main obstacle for providing better support to these applications is that, at present, the most of web content is not machine-accessible. Problems arise in querying data sources due to unstructured contents of web pages (HTML), we cannot directly extract data into a new structured form. To address this problem, we propose a system that can perform web data extraction in a Mash up format. Proposed system has a fully visual and interactive user interface having new technique and approach. The availability of information on the Internet has grown exponentially; the Internet users now face problems of extracting information relevant to their needs. Proposed application provides users means to directly extract data from various web pages. These data are mostly found in an unstructured HTML format. Proposed application has the capability to convert this unstructured data into fully structured forms so that it can be easily sent out as emails and SMS updates.

II. Related Work

Web page extraction approaches are mostly based on either generating wrappers, or tag based or tree based approaches. These approaches are either domain- dependent, or need human intervention, or require large amount of dataset for training purpose. Our focus is on generating wrappers using DOM tree and than using Simple Tree Pattern Matching Algorithm automatically load the extractor so that it will extract the data from various web pages, where they are mostly formed in an unstructured HTML format, then transformed the data into a structured format. It also focused on delivering the extracted data through various means, such as *delivering* the resulting structured data to external applications such as database management systems, data warehouses, business software systems, content management systems, decision support systems, RSS publishers, email servers, or SMS servers. Alternatively, the output can be used to generate new web Services out of existing and continuous changing web sources. Many techniques use Tag-based approach that uses HTML parser to convert each web page into DOM tree. It then applies traditional pattern reduction techniques are template dependent and include substantial features of the DOM tree.

III. Proposed System

In the proposed system is development of application that automatically and repeatedly extracts data from web pages with changing content and it delivers the extracted data to a database or some other application. It should include all the facilities that a user expects generally. The project objectives are:

1. Extract data from various web pages, where they mostly formed in an unstructured HTML format, into a new structured format such as XML or XHTML.
2. Implementation of web data extraction and stages in making a Mash up.
3. Implement web data extraction by visually extract targeted data from data sources (web pages). Afterward, we combined web data extraction with the stages of making a Mash up, e.g. data retrieval, data source modelling, data cleaning/ filtering, data integration and data visualization.

The task of web data extraction performed by such a system is usually divided into five different functions:

1. Web interaction, which comprises mainly the navigation to usually pre-determined target web pages containing the desired information.
2. Support for wrapper generation and execution, where a wrapper is a program that identifies the desired data on target pages, extracts the data and transforms it into a structured format.
3. Scheduling, this allows repeated application of previously generated wrappers to their respective target pages.
4. Data transformation, which includes filtering, transforming, refining, and integrating data extracted from one or more sources and structuring the result according to a desired output format (usually XML or relational tables).
5. Delivering the resulting structured data to external applications such as database management systems, data warehouses, business software systems, content management systems, decision support systems, RSS publishers, email servers, or SMS servers. Alternatively, the output can be used to generate new web services out of existing and continually changing web sources.

IV. System Design And Architecture

Figure shows the Architecture of System. It is a high-level view of a typical fully-fledged semi-automatic interactive web data extraction system. This system comprises several tightly connected components and interfaces three external entities.

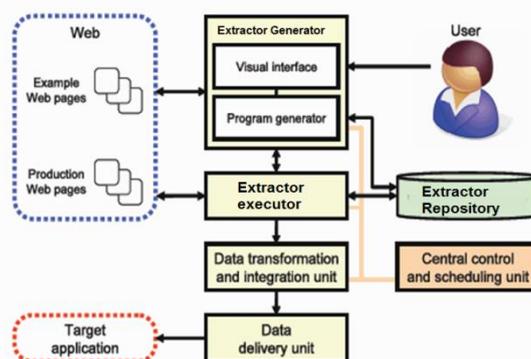


Fig. 1 Extension system

- **Web Interaction:** specification of URL and marking of data.
- **Extractor Generation:** a program that we shall design that will download the required page, parse it and generate tree from parsed output.
- **Extraction:** comparing generated tree with saved one to extract updated data from current web page available.
- **Scheduling :** repeated execution of above
- **Data Transformation:** E-Mail Updates, SMS Updates , Data Warehousing, XML, Tabular View
- **Generation of Mash-ups:** Deliverable APIs or libraries that can be used in other applications to extract required information from various web sources.

The web contains pages with information of interest. In the target application, extracted and refined data will be ultimately delivered to the user, who interactively designs the extractor. The extractor generator supports the user during the extractor design phase. It commonly has a visual interface that allows the user to define which data should be extracted from web pages and how this data should be mapped into a structured format such as xml.

V. Algorithm Used For the System

HTML parser (DOM parser) is used for parsing the web page. DFS algorithm is for searching appropriate data on web page.

1. DOM PARSER:

The Document Object Model is a programming API for HTML and XML documents. It defines the logical structure of documents and the way a document is accessed and manipulated. However, XML presents this data as documents, and the DOM may be used to manage this data. With the Document Object Model, programmers can create and build documents, navigate their structure, and add, modify, or delete elements and content. Anything found in an HTML or XML document can be accessed, changed, deleted, or added using the Document Object Model, with a few exceptions - in particular, the DOM interfaces for the internal subset and external subset have not yet been specified. The Document Object Model (DOM) is a cross-platform and language-independent convention for representing and interacting with objects in HTML, XHTML and XML documents. Objects in the DOM tree may be addressed and manipulated by using methods on the objects. The public interface of a DOM is specified in its application programming interface (API). To render a document such as an HTML page, most web browsers use an internal model similar to the DOM. The nodes of every document are organized in a tree structure, called the DOM tree, whose uppermost node is the *Document* object. When an HTML page is rendered in a browser, the browser downloads the HTML into local memory and parses it, using the DOM to construct the internal data structures employed to display the page in the browser window. The DOM is

also the way in which JavaScript sees the state of the browser and the current HTML page. JavaScript code uses the DOM APIs to inspect or improve a web page. Because DOM supports navigation in any direction (e.g. parent and previous sibling) and allows for arbitrary modifications, an implementation must at least buffer the document that has been read so far (or some parsed form of it.).

2. DFS Algorithm:

Depth-first search (DFS) is an algorithm for traversing or searching a tree, tree structure, or graph. Depth-first search is an organized way to find all the vertices reachable from a source vertex. We first select the root node of a tree or any random node (selecting some node as the root in the graph case) and explore as much as possible in a branch and then return to a fixed point (backtracking). DFS is an uninformed search that progresses by expanding the first child node of the search tree that appears and thus going deeper and deeper until a goal node is found, or until it hits a node that has no children. Then the search backtracks, returning to the most recent node it hasn't finished exploring. Here, the word backtrack means that when you are moving forward and there are no more nodes along the current path, you move backwards on the same path to find nodes to traverse. All the nodes will be visited on the current path till all the unvisited nodes have been traversed after which the next path will be selected.

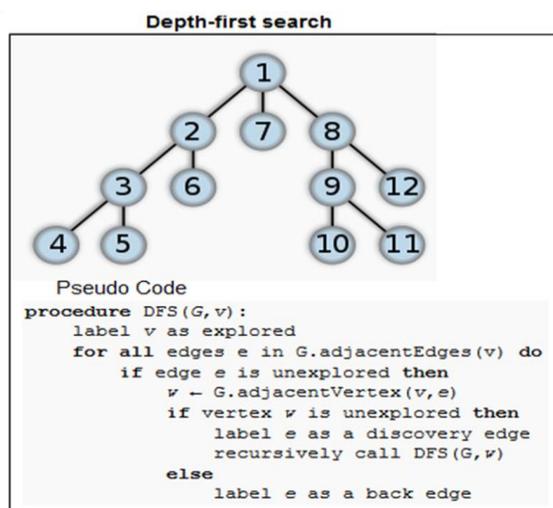


Fig. 2 DFS algorithm

A convenient explanation of a depth first search of a graph is in terms of a spanning tree of the vertices reached during the search. The advantage of depth-first Search is that memory requirement is only linear with respect to the search graph.

3. HTML Parser

In the HTML Parser a Java library used to parse HTML in either a linear or nested fashion. The mainly used for transformation or extraction. The features of HTML parser are filters, visitors, custom tags and easy to use JavaBeans. It is a very fast, robust and also well tested package. HTML Parser is a super. It is fast to real-time parser for real-world HTML. What has attracted most developers to HTML Parser has been its simplicity in design, speed and ability to handle streaming real-world html. The two fundamental use-cases are used in the HTML Parser. That are handled by the parser are extraction and transformation (the syntheses use-case, where HTML pages are created from scratch, is better handled by other tools closer to the source of data). While earlier versions concentrated on data extraction from web pages. The version 1.4 of the HTML Parser has substantial improvements in the area of transforming web pages. It can be simplified tag creation and editing, and HTML method output. In general, to use the HTML Parser you will need to be able to write code in the Java programming language. Although some example programs are provided that may be useful as

they stand, it's more than likely you will need (or want) to create your own programs or modify the ones provided to match your intended application

4. Simple tree pattern matching algorithm

Simple tree pattern matching is the act of checking a derived sequence of tokens for the presence of the constituents of some pattern. In contrast to pattern recognition, the match usually has to be exact. The patterns generally have form of either sequences or tree structures. Uses of pattern matching include outputting the locations of a pattern within a token sequence. The output is component of the matched pattern. Then it substitutes the matching pattern with some other token sequence for example search and replace. The sequence patterns (e.g., a text string) are frequently described using regular expressions and it will matched using techniques such as backtracking.

The Tree Match algorithm, can find all different matchings of a tree pattern in the data source directly. It is no longer a decomposition-matching as well as merging process. Tree patterns are used in some programming languages as a general tool to process data based on its structure, e.g., Haskell, ML. The symbolic mathematics language is Mathematica. The special syntax for this language is expressing tree patterns. Mathematica language construct for conditional execution and value retrieval based on it. For simplicity and efficiency reasons, these tree patterns defect some features that are available in regular expressions. Often it is possible to give alternative patterns that are tried one by one, which yields a powerful conditional programming construct. Pattern matching usually include support for guards.

VI. Result And Discussion

Following Fig shows a better example of working of our system; consider a Horoscope website where the data gets updated every day. In this scenario user will specify some data about he wants to regularly (in the given fig the dark circled area of birthday horoscope).



Fig. 3 Result of System

With the help of our proposed system whenever the data Specified will get change or updated user will get notification through various means.

VII. Conclusion

In this paper we proposed an Information Extraction from Web Pages system which will able to extract web data through various regularly changing websites. This paper mainly concentrates on deep websites as data on deep websites is in the form of a particular template, Template matching approach will be used for extraction.

By using DOM parser, DFS algorithm and Simple tree pattern Matching Algorithm system will is able to extract proper data from unstructured website. In this paper, we tried to design a system that gives complete solution for Data Extraction from web pages. The system automatically and repeatedly extracts data from the web pages with changing contents and delivers the extracted data to the users through various means like SMS, Email, and web page.

Acknowledgement

It is my privilege to acknowledge with deep sense of gratitude to my research work guide Dr. Rekha Rathore for her valuable suggestions and expert guidance throughout my course of study and timely help given to me in the completion of my paper work.

References

- [1] Craig. A. Knoblock, Kristina Lerman, Steven Min ton, Ion Muslea, “Automated Web Data Extraction Using Template Matching”, IEEE International Advance Computing Conference, 2013.
- [2] R. A. Gultom, R. F. Sari, and B. Budiardjo, “Implementing web data extraction and making mashup with xtractorz,” IEEE 2nd International Advance Computing Conference, pp. 385393, 2010.
- [3] X. Zheng, Y. Gu, and Y. Li, “Data extraction from web pages based on structural-semantic entropy,” pp. 93102, April 2012.
- [4] Sandeep Sirsat¹, Dr. Vinay Chavan, “Pattern Matching for Extraction of Core Contents from News Web Pages” Second International Conference on Web Research (ICWR) 978-1-5090-2166-6/16/\$31.00 ©2016 IEEE
- [5] Robert Baumgartner Institute for Information Systems Vienna University of Technology and Lixto Software GmbH, Austria, “Web Data Extraction”, 2009.