

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/261848753>

# Web page analysis based on HTML DOM and its usage for forum statistics and alerts

Conference Paper · April 2010

CITATIONS

0

READS

181

4 authors, including:



**Robert Gyorodi**

University of Oradea

54 PUBLICATIONS 265 CITATIONS

[SEE PROFILE](#)



**Cornelia Győrödi**

University of Oradea

55 PUBLICATIONS 274 CITATIONS

[SEE PROFILE](#)



**George Pecherle**

University of Oradea

23 PUBLICATIONS 107 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Convergence of university practical training for integration with success in the labor market [View project](#)



Applications of Data Mining in the Decision Making Process [View project](#)

# Web Page Analysis Based on HTML DOM and Its Usage for Forum Statistics and Alerts

ROBERT GYÖRÖDI, CORNELIA GYÖRÖDI, GEORGE PECHERLE, GEORGE MIHAI CORNEA

Department of Computer Science

Faculty of Electrical Engineering and Information Technology, University of Oradea

Str. Universitatii 1, 410087, Oradea

ROMANIA

rgyorodi@uoradea.ro, cgyorodi@uoradea.ro, gpecherle@uoradea.ro, generalmip@yahoo.com

**Abstract:** - Message boards are part of the Internet known as the 'Invisible Web' and pose many problems to traditional search engine spiders. The dynamic content is usually very deep and difficult to search. In addition, many of these sites change their locations, servers, or URLs almost daily creating problems with the indexing process. However, during the growth of the World Wide Web and with the help of search engines, they represent an important source of information to solve different problems. Another interesting feature of this type of webpages is that a big community has been developed, expressing different opinions and discussing various topics. Using special retrieval and indexing algorithms, mostly based on the HTML DOM tree, we have developed an algorithm to obtain detailed and accurate trend statistics that can be used for different marketing solutions and analysis tools.

**Key-Words:** - data analysis, data models, HTML DOM, information extraction, text recognition

## 1 Introduction

During the last decade, most websites have been providing the information generated from their structured data in an underlying database through certain predefined templates or layouts. The message boards are also part of these categories.

Following the great number of message boards available on the Internet, these semi-structured web sources contain rich and unlimited valuable data for a variety of purposes. Extracting this data and then rebuilding them into a structured database represents a challenge to perform automatic data mining from web sources.

Several approaches have been reported in the literature for the purpose of building and maintaining mining scripts for semi-structured web sources. Some of them can be classified as the so-called wrappers [1]. The wrapper technique allows automatic data extraction through a predefined wrapper created for each target data source. The wrappers then accept a query against the data source and returns a set of structured results to the calling application. This method is easy to implement but is hard to maintain and extend.

On the other hand, there are several automatic methods without requiring an initial manual learning process. For example, some methods are based on the automatic generation of a template from the first multiple pages, before extracting the rest of the data based on the template [2], [3], [4]. Another method to generate a template automatically is based on finding certain repetitive patterns contained in a single page

based on the structure of the DOM tree. The method is extensively described in [5].

Taking into account the demand for such services, we have developed an algorithm that can accurately analyze web forums and message boards and provide accurate statistics. This algorithm has been developed using special retrieval and indexing methods, that will be described later in this paper.

Even if the existing methods described above can be successfully applied in our system to detect different repetitive structures like forums, threads and posts, it alone can't gather all the information required for the further processing of the message board and for an effective update process.

Until we further develop the fully automatic template generation system, we are using an alternative method that combines a reverse learning process with a system to test the templates we have already used in our system. Because of the special needs of our crawling process, we will always look for certain types of information that have a certain link between them.

Besides offering accurate trend statistics, the unique way in which our spider works allows us to give almost real time alerts when someone posts a message which contains a certain phrase in it. This way, we can offer services to our users, regarding a wide range of topics discussed in the forums: buying and selling of different products, buying and selling of services, offer help on different topics, etc.

Alexa.com [9] is a web information service that offers free traffic metrics, search analytics and demographics data. Coupled with the services offered by

Alexa.com, we can also offer geo-targeting features for the statistics and the alert system that can greatly improve the quality of the services. This way, we will only provide graph results for a required region. Being able to analyze, for example, the impact of a TV advertising campaign launched only in a certain region for testing purposes, the company doing this analysis can combine the internal intelligence (consisting in sales trends, for example) with the external intelligence provided by our system.

Without the use of an advanced system to index and process the message boards, all the above mentioned features would not be possible. The old methods of indexing would require huge levels of bandwidth and computation power to be able to provide daily alerts of new messages that are updated in the forums. However, using our implementation, the bandwidth usage required by the indexing of a new message will be of an average of 70% the page that contains the message. This rate is even better for frequently updated forums in which the number of new messages between two consecutive scans is higher. The worst case appears when we perform a search for a single message contained in a topic located within the oldest topics. However, the majority of the new messages will be located within the first topics, as they are the most active ones.

The introduction of our service will also bring great benefits to the forum owners, who will get more visitors and more potential members to their communities through the alert system.

Our statistics system will be a useful tool for all the marketing companies offering different marketing solutions to their clients, by providing external intelligence. It will also be useful for managers, political analysts, etc.

## 2 Our Implementation of the Reverse Process to Determine the Templates

The reverse process to determine the templates [1] has to take place if all the previously determined templates failed to return an acceptable dataset, after they have been applied over the desired page. In this case, human intervention is required to generate the template. We will illustrate the process only for message boards homepages.

This process is divided in two steps:

1. After it is determined that human intervention is required, a separate service will be used, that will provide support for region-of-interest (RoI) input. At the beginning of this step (because of the constant updates that are made in a forum's homepage, such as the last user who posted, number of messages, etc.), we have to save the HTML source of the page so that the reverse

processing application can work on the same dataset as the human did.

2. When the database is populated with the targeted regions-of-interest, the reverse process can start to analyse the available data. It will search for all the nodes that contain the regions-of-interest provided by the human at step one above [1]. If more than one node is found to contain a targeted region-of-interest, the process will try to group them into the smallest subgraph containing all of the determined paths, taken once. This case is valid, for example, when the system has to process a page containing posts that contain other quoted posts.

For the first step, we have developed a web application to handle the interaction with the human, because it facilitates a faster and easier to implement access to the database. It consists in three major parts: a frame to browse the targeted message board page, a textbox containing the page's HTML source in case it is needed and the RoI input boxes. Using this interface an average user can insert or update around 30 message boards per hour. This results in more than 200 new templates added daily. Together with the increasing number of templates stored inside the database, this also generates an exponential capacity to process new message boards using the templates already stored in the database.

The second step is the one when the actual template consisting in the absolute un-indexed path to every RoI is determined, by searching for the node containing the targeted information.

In order to eliminate all possible errors due to repetitive information (for example, two forums with the same number of messages), all the matching nodes have to be selected. At the end of the process all these nodes will be grouped, one of each category, in the minimum tree spanning all the required vertices once [6]. All of these operations are done in the previously saved page's HTML source, to eliminate any discordances.

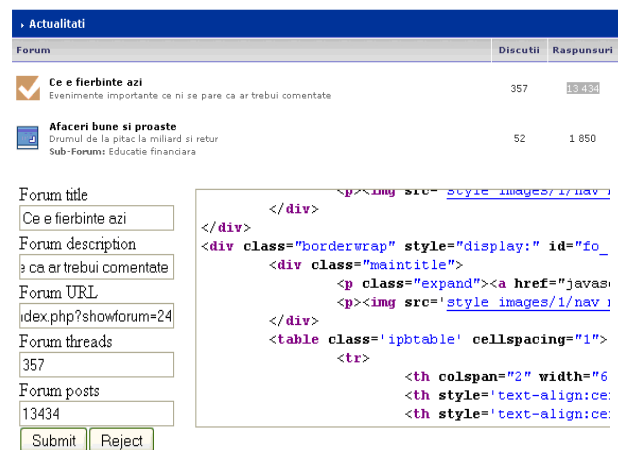


Fig. 1. Interface of our web application to provide RoI for forums

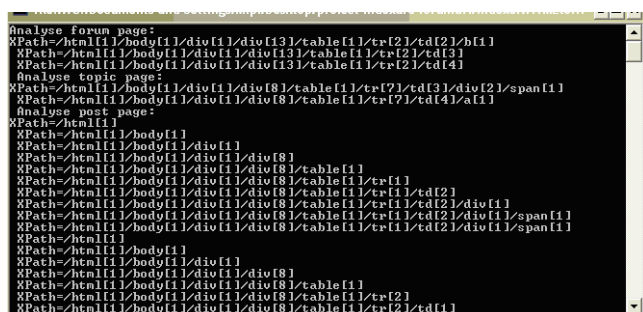


Fig. 2. Example of resulting paths from the DOM tree, after applying our algorithm

After the paths for a certain forum have been determined, the template path is considered to be the deepest sub-tree containing one of the predetermined RoIs, followed upwards until we reach a common root for all the selected RoI.

During this process, the following errors could appear resulting in a request to re-enter another set of RoI:

1. The most common one: a certain RoI was not found – it is usually caused by human mistake or by charset problems.
2. After the template was generated and a run test has been done, the template returns only one row – this is more a kind of a warning, because the template can be good but the page may consist of only one forum, thread or post.

### 3 Our Implementation of the Process to Determine the Match of a Template with a New Forum

This is the first step taken every time a new forum has been added in the indexing list. This was not implemented only for this reason, because it is also used to determine if there are changes in the structure of a forum that already has a template. With some new methods added to the first class (that only determines if a template returns a plausible result set), it is also possible to determine if a template in the database can be used over the new forum. A clustering based on the partial roots gives us an algorithm with  $O(\log(n))$ , where  $n$  is the number of templates stored in the database.

The result check is done based on a few simple presumptions (this example is for forums):

1. The messages and the threads in a forum will always be numbers containing spaces, dots, or commas among numeric characters.
2. The number of messages will always be higher or equal to the number of threads.
3. The average length of the forum title will always be smaller than the average length of the description.

## 4 Our Indexing and Alert System

As it was told before, the main feature on which we are focusing is to provide a cost effective and fast way to crawl the deep invisible Internet. Other methods that have been implemented so far use an algorithm to determine the rate at which a page is updated and approximate the next date at which the new scan should occur.

This method doesn't provide either a cost effective way to crawl the Internet or an accurate one. The first inconvenience comes when a new page is crawled for the first time. In order to calibrate and obtain the most convenient update time, the page will be crawled several times more fervently. After the statistics algorithm has enough data to return a smaller error, it can be subject to an error increase if the way in which the page changes is unpredictable. Another disadvantage is that, this kind of algorithm will never be able to provide real time access to page changes without heavily overloading both the crawled server and the crawling server, in order to minimize the time between two scans.

The way in which our system works allows us to determine if there was an update on a certain forum. Tracking down its structure, we can determine exactly what page was updated and take the required actions.

The typical formulas to calculate the costs are as follows:

1. Forum update check:

$$\text{Cost/day [KB/day]} = \text{mainPage [KB]} * 6 * 24$$

mainPage = the size of the main (home) page in KB

This is the cost representing the amount of information exchange between the crawling and the crawled server to determine if a new message has been added on a certain forum. The update is done once, every 10 minutes. This way, the size of the main page is multiplied by 6 (in one hour, there are 6 groups of 10 minutes) and then by 24 (the number of hours in one day).

2. Topic update check:

$$\text{Cost/topic [KB/topic]} = \text{mainPage [KB]} + \text{topicOffset [number]} * \text{forumSize [KB]}$$

mainPage = the size of the main (home) page in KB

topicOffset = the nearest integer greater than or equal to (the number of topics, between the first topic and the topic that was updated / number of topics per page)

forumSize = the size of the forum in KB

This is the cost representing the amount of information exchange between the crawling and the crawled server to determine which topic has been updated. Since we already have the information contained in mainPage, we don't have to download it again, in case of a periodic update.

As it can be seen from the cost formula, the worst case occurs when we have an update in the last few pages of the forum (topicOffset takes the maximum value). However, this case does not occur too often, because the last topics are the older ones, abandoned or closed.

### 3. Message update check

**Cost/message [KB/message] = mainPage [KB] + topicOffset [number] \* forumSize [KB] + k [number] \* topicSize [KB]**

mainPage = the size of the main (home) page in KB

topicOffset = the nearest integer greater than or equal to (the number of topics, between the first topic and the topic that was updated / number of topics per page)

forumSize = the size of the forum in KB

k = the nearest integer greater than or equal to (msgNb / msgPerPage)

msgNb = the number of messages that have been updated

msgPerPage = the number of messages in a page

topicSize = the size of the topic, in KB

This is the cost representing the amount of information exchange between the crawling and the crawled server to determine and to process the new messages. As it can be seen from the formula, because the new messages are well ordered there is no need to check all the messages, only the last msgNb of them determined at the previous stage. The first two arguments are already obtained from the previous step so they will be shared every time there is more than one message to update.

Considering the fact that at every new update that we perform, we have more than one topic to process and more than one message to add to the database and to process, we won't have the entire cost calculated using the cost/message formula. Instead, parts of that cost will be shared by multiple updated messages, resulting in an even lower total cost per number of messages indexed. As an example from a medium Romanian forum, we have the following average values during 24 hours of testing:

1. 25 of 144 main page checks returned a positive message to follow through. This represents a percentage of 17.36% .

2. 3 follows were performed through sub-forums adding an extra cost required to further follow the forum tree downwards.

3. All the messages were in the first page of the topics, resulting in an offset cost of 1 (the best case because all the messages were in new threads).

4. The messages were clustered through the positive follow through with an average of 2.32 further decreasing the average cost per message.

At the end of the test period we had:

- 8485.7 KB downloaded for the forum homepage (consisting in the root of the forum). From this bandwidth usage, 17.36% was used to further detect threads.

- 100% of the threads were found in the first page of the forum. With an offset cost of 1, the total bandwidth used here was 1726.6 KB.

- there were 92 messages to process. The total bandwidth used to process those messages was 10952 KB. 74 page loads were required for this task.

The total bandwidth used for the scanning of the forum through the test period of 24 hours was of 21671.8 KB (21.2 MB). The extra bandwidth of 435.1 KB comes from the messages to which the tree has been followed through sub-forums.

The share rate resulting in decreased bandwidth usage and greater performance was as follows:

- The main page was loaded 144 times during the test period of 24 hours. Only the 25 follow through are taken into account further.

- 31 threads were marked as updated using the 25 scans of the main page. This results in a share rate of 80.64%.

- furthermore, 74 message pages were determined by analyzing the thread information, resulting in a share rate of 41.89%. From those message pages, we were able to extract 92 messages.

The approximate bandwidth that could be used to periodically update this site in the worst possible case is 35MB, calculated using only the positive presumption that the forum offset is still equal to 1 because we have updates in the most recent threads. The increased size appears because of zero clustering. So for every new message the tree has to be followed downwards without sharing nodes.

## 5 Complete Process Diagram for Our System

In the next diagram we will try to illustrate the entire process that takes place every time a new forum is added to the database to be scanned and every time a forum is updated to scan for new messages.

The diagram has two main parts:

- the reverse process to determine the templates and the first indexing process
- the periodical update indexing, when data can be either added to the statistics database or to the alert system database

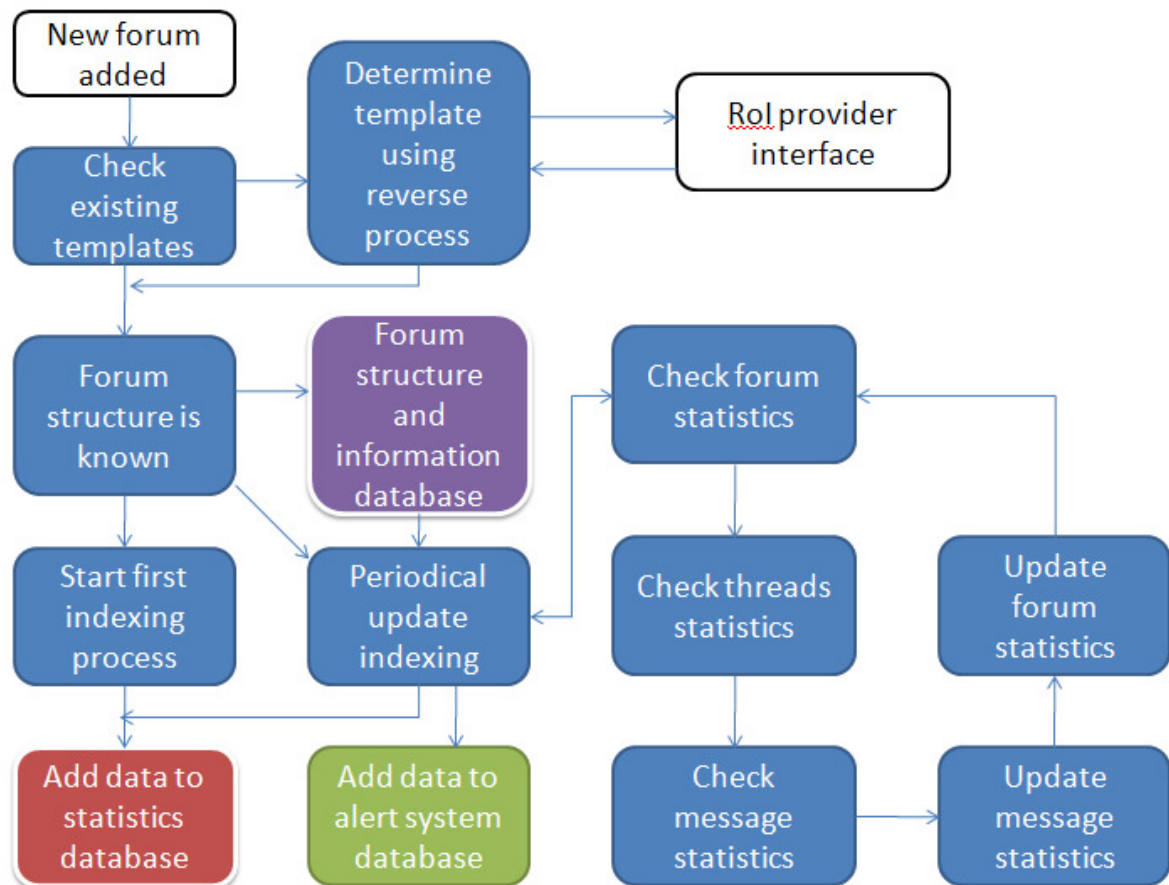


Fig. 3. Complete process diagram for our system

## 6 Large Database Support

Over time, the database where we store all the information may become very large, due to the amount of data gathered from multiple forums. Also, our system should support expanding to an unlimited number of forums to be analyzed, therefore large database support is needed.

For practical reasons, we use MySQL, however any other database management systems that support large databases can be used. MySQL has a feature called partitioning, that allows you to distribute portions of individual tables across a file system according to rules which you can set largely as needed. In effect, different portions of a table are stored as separate tables in different locations [7].

Some of the advantages of using partitioning include:

- Being able to store more data in one table than can be held on a single disk or file system partition.
- Data that loses its usefulness can often be easily removed from the table by dropping the partition containing only that data. Conversely, the process of adding new data can in some cases be greatly facilitated by adding a new partition specifically for that data.

- Some queries can be greatly optimized in virtue of the fact that data satisfying a given WHERE clause can be stored only on one or more partitions, thereby excluding any remaining partitions from the search. Because partitions can be altered after a partitioned table has been created, you can reorganize your data to enhance frequent queries that may not have been so when the partitioning scheme was first set up. This capability, sometimes referred to as partition pruning, was implemented in MySQL 5.1.6 [7].

One of the main reasons we need to use partitioning is that usually data collections are far beyond any reasonable amount of RAM that you can install. In most database systems, the indexes are cached in RAM, to allow fast retrieval of records. That's why we need to use a different approach, because data can reach sizes that can't be stored in RAM [8].

One useful way to apply partitioning in our system is to make partitions based on the year, because each message in a forum has a timestamp attached. This way, we will have all our forum messages organized by year, allowing an easier management and faster data retrieval.



Since the native date type is not supported, we must convert the date into an integer. In addition to the list of allowed functions, we must take into account the fact that only two date functions can trigger the partition pruning. Thus, if we have to deal with a date column, we need to use one of them (YEAR or TO\_DAYS).

When using the YEAR() function, partitioning is easy, readable, and straightforward [8].

```
CREATE TABLE by_year (d DATE)
PARTITION BY RANGE (YEAR(d))
(PARTITION P1 VALUES LESS THAN (2007),
PARTITION P2 VALUES LESS THAN (2008),
PARTITION P3 VALUES LESS THAN (2009),
PARTITION P4 VALUES LESS THAN
(MAXVALUE))
```

The basic concept behind partition pruning is that of not scanning partitions where there can be no matching values. For example, when we need to retrieve messages older than 2007, we will only scan partition p1 from the example above and do not scan partitions p2, p3 or p4. This greatly improves performance and the time needed to retrieve data from our tables.

## 7 Conclusions and Future Work

The algorithm we propose is a new concept in the process of indexing structured web pages with a lower cost (lower bandwidth and increased performance). The method we described in this article can be applied to other types of structured web pages, not only for web forums.

Also, it can be a useful tool for gathering marketing statistics, based on the discussions on a certain product or service, identified by one or more keywords.

There are several issues that require further modifications and implementation in the algorithm. The first one represents the poor reachability of a message that is placed among the last threads. In the case in which we have an extremely active forum with large numbers of updates, the time between a scan of the main page and until all the required messages were found is big enough to manage new messages. The scan for new messages will be stopped as soon as we reach the pre-determined number of messages, so that it is possible to miss the last messages. Another issue appears when there are deleted messages or threads, or moved threads. This case is completely unhandled, human intervention being required to correct the data and continue the scan.

In the near future, we want to implement handling routines for these special cases. Another future implementation we want to do is to switch from the test server to a live one and optimize it so that it can work on multiple machines simultaneously.

## References:

- [1] Z. Akbar and L.T. Handoko, "Reverse method for labeling the information from semi-structured web pages", proceeding of the 2009 International Conference on Signal Processing Systems pp. 551-555
- [2] L. Arlota, V. Crescenzi, G. Mecca, P. Merialdo, Automatic annotation of data extracted from large websites, in: Proceedings of the WebDB Workshop, 2003, pp. 7-12.
- [3] N. Kushmerick, Regression testing for wrapper maintenance, Proceedings of the sixteenth national conference on Artificial intelligence and the eleventh Innovative applications of artificial intelligence conference innovative applications of artificial intelligence, July 18-22, 1999, pp. 74-79.
- [4] J. Wang, F.H. Lochovsky, Data extraction and label assignment for web databases, Proceedings of the 12th international conference on World Wide Web, May 20-24, 2003, pp. 187-196.
- [5] Li WeiDong, Dong Yibing, Wang RuiJiang, Tian HongXia, "Information Extraction from Semi-Structured WEB Page Based on DOM Tree and Its Application in Scientific Literature Statistical Analysis System", 2009 IITA International Conference on Services Science, Management and Engineering
- [6] Jie Zou Le, D. Thoma, G.R., "Combining DOM tree and geometric layout analysis for online medical journal article segmentation", Digital Libraries, 2006. JCDL '06. Proceedings of the 6th ACM/IEEE-CS Joint Conference, pp. 119 - 128
- [7] MySQL Developer Zone (<http://dev.mysql.com>) – "Chapter 18. Partitioning"
- [8] MySQL Developer Zone (<http://dev.mysql.com>) – "MySQL partitions in practice – Giuseppe Maxia"
- [9] Alexa.com – The Web Information Company - Free traffic metrics, search analytics, demographics