# Xpath에 의한 인터넷 문서의 레이아웃 추출 방법에 관한 연구
## A Study on Layout Extraction from Internet Documents Through Xpath

선복근, 한광록

호서대학교 컴퓨터공학부

Bok-Keun Sun(bksun@office.hoseo.ac.kr), Kwang-Rok Han(krhan@office.hoseo.ac.kr)

## 요약

현재 뉴스 데이터 등 대부분의 인터넷 문서는 일정한 템플릿을 기반으로 작성되고 있으며 템플릿은 메인 데이터 이외에 인덱스, 광고, 헤더데이터 등 정보검색에 도움이 되지 않는 형태로 구성되어 있다. 이는 인터넷 문서를 정보검색의 데이터로서 사용하려고 할 때 적합한 형태가 아니다. 그러므로 다양한 정보검색 분야에서 인터넷 문서를 처리하기 위해선 광고, 페이지 인덱스 등의 부가정보를 분별해야 한다.

따라서 본 논문에서는 웹페이지의 레이아웃에 영향을 미치는 블럭 태그의 특징과 구조를 파악하고 웹페이지간의 거리를 계산하여, 웹페이지의 레이아웃을 검출하는 방법을 제안한다. 실험결과 1000개의 문서 중 640개를 분류했으며, 평균 64%의 recall 수치를 얻을 수 있었다. 이 방법을 데이터 추출, 문서요약 등의 정보검색 분야의 전처리 과정에 적용할 경우 문서의 자동화 처리 시간을 감소시키고 처리의 효율성을 높일 수 있을 것으로 기대된다.

■ 중심어 : | 정보검색 | 데이터 추출 | 레이아웃 | HTML | XML 기술 |

## Abstract

Currently most Internet documents including news data are made based on predefined templates, but templates are usually formed only for main data and are not helpful for information retrieval against indexes, advertisements, header data etc. Templates in such forms are not appropriate when Internet documents are used as data for information retrieval. In order to process Internet documents in various areas of information retrieval, it is necessary to detect additional information such as advertisements and page indexes.

Thus this study proposes a method of detecting the layout of web pages by identifying the characteristics and structure of block tags that affect the layout of web pages and calculating distances between web pages. As a result of experiment, we can successfully extract 640 documents from 1000 samples and obtain 64% recall rate. This method is purposed to reduce the cost of web document automatic processing and improve its efficiency through applying the method to document preprocessing of information retrieval such as data extraction and document summarization.

■ keyword : | Information Retrieval | Data Extraction | Layout | HTML | XML Technologies |

# I. Introduction

Because most web contents are user-friendly, users who use information regard the user-friendliness as natural. With the explosive growth of web contents and data, however, documents must be converted into forms that can be processed automatically for structural data process [1]. Accordingly, academic and industrial areas have been studying on methods of abstracting web data to use in information retrieval. HTML has not been fully standardized and data are user-friendly not machine-friendly, however, there is established tool, technology or standard [2].

Researchers are actively studying methods of processing web contents mechanically through XML and expect even more active research in the future. Although a relatively late technology compared to HTML, XML is evolving into a language fit for data exchange between heterogeneous machines and automatic data process [3]. However, the number of XML documents accessible on the web is still insignificant compared to that of HTML documents. Although XML technology will become common on the web in the future, HTML will be used continuously until XML technology reaches the level of design convenience and user-friendliness of HTML. Thus there are various researches on methods of abstracting data by analyzing HTML documents or converting HTML documents into XML documents, and these researches are considered useful in the areas of information retrieval and data processing. Using these methods, data on websites can be processed automatically by users and applications.

HTML-based web data are usually developed focused on design and layout rather than on data processing. This makes it difficult to use web data in data-based information retrieval or calculation or as source data for applications.

Most sites including news sites that contain data to be used in information retrieval use data presentation structure called template in arranging information to enhance legibility [4-6]. As shown in [Fig. 1], the main content is positioned at the center and headers, footnotes, advertisements, etc. are placed around. In the template, information is divided visually and semantically.
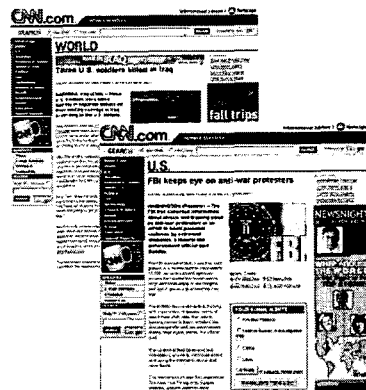


Fig. 1. Example page using a template.

Pages using the same template have similar layouts and the roles and positions of semantic groups are also similar. Thus, it may be useful for automatic data processing to identify pages of the same layout.

To develop a method of identification, this study examines how to abstract the structure of documents from web pages. The research is focused on the structure and characteristics of specific HTML tags that affect the layout of web pages. It calculates distances between web pages and, if the distance is over a specific threshold, the pages are judged to use the same layout. The

developed method is supposed to be applied to information retrieval such as automatic data extraction and provide information about the structure of web pages using a template, and ultimately, to reduce the cost of web document automatic processing and enhance processing efficiency.

## II. Related works

Recently there are increasing interests in researches to analyze the internal structure of web pages and abstract data. Such researches usually purpose to abstract data through analyzing HTML documents and converting their structure into document forms.

Artificial intelligence areas such as machine learning are focused on formulating rules that can abstract sharable general rules through analyzing sample pages and using them for abstracting data from other similar pages [2].

WYSIWYG Web Wrapper Factory(W4F) is a tool kit to generate web wrappers [7]. It includes a unique language to formulate rules for recognizing web sites and for abstracting data from web pages. In addition, it provides a mechanism of converting abstracted data into a structure for a specific purpose. The language used here, however, is system-specific one not one compatible with XML.

Web Language(WebL)[8][9] of Compaq is a procedural language for web wrappers. It provides a powerful abstraction language, mixing regular expressions and reflexive pass expressions, but not compatible with XML technology.

ANDES of IBM uses XML core technologies such as XHTML[10], XPATH[11] and XSLT[12] and provides data abstraction framework [13]. Its approach to XHTML document tree processing uses Xpath(XML Path) and, if necessary, can be combined with regular expressions. The mechanism of XSLT style sheet templates is useful in identifying data extraction patterns represented by the elements of XHTML documents and Xpath, and defines style sheets using a simple description method. However, there are many difficulties in automatic conversion of HTML into XHTML using XSLT including the irregularity of HTML, client scripts and dynamic HTML, so many heuristic methods should be employed. If it is possible to convert all HTML into XHTML correctly using XSLT, any XML tool can be applicable in handling the XHTML pages because XHTML is based on XML [13].

Sometimes annotation-based methods are used in data conversion and extraction. Annotation is a kind of meta data and additional information of web pages [1][4][6][14]. Using annotation, it is possible to abstract and convert information. In general, annotation is related to URI(Uniform Resource Identifiers), Xpath and Xpointer[15], which indicate the position of annotation in documents. The author classifies web pages into a number of semantic groups and save information about each semantic group such as its role and position as an annotation in XML form. Using annotation, it is possible to change the structure of web pages without changing semantic groups or losing information. Adding annotation into contents, it is possible to convert and abstract contents accurately and precisely. The biggest problem of annotation, however, is cost. It is almost impossible to add annotation to a large volume of frequently updated pages like news. However, most of such web pages usually use

templates. Pages using the same template have similar layouts and the paths and roles of their semantic groups are almost identical.

In order to research web document data extraction and solve problems in it, this study examines a method of abstracting layout information from web documents. This study is focused on the structure and characteristics of specific HTML tags that affect web page layout. It calculates distances between web pages and if the distance is below a certain threshold it judges that the pages use the same layout. Through this, this study purposes to provide a method that can be applied to all pages of the same layout in information retrieval areas such as data extraction and annotation making.

## III. Layout group detection

In general the main content of a web document is positioned at the center and additional information such as indexes and advertisements are around the main content [Fig. 1]. Most documents of this type are created and updated using templates. Some of HTML tags are related to template and layout and information about the tags can be used in calculating distances between pages. First of all, let us discuss the characteristics of HTML tags and then how to calculate distances between pages.

## 1. Block level tag

Block level HTML Tag is a key element in the layout structure of web pages. [Table 1] shows an example of Block Level Tag.

Such a tag is called a layout tag in this paper. After analyzing HTML and abstracting layout

tags, we structuralize each tag as an Xpath expression and refer it as data source for calculating distances between pages.

### Table 1. An example of block level tag.

| Feature | Block level tag |
|---------|-----------------|
| Table | TABLE, THEAD, TBODY, TR, TH, TD |
| Form | BUTTON, FORM, TEXTAREA |
| others | HR, OL, UL, LI, DIV, SPAN, P |

[Fig. 2] shows a HTML document and an example of Xpath abstracted from the document.

## 2. Detection algorithm

To determine if HTML documents use the same template and, as a result, have the same layout, we analyze layout tags as presented above. For this, we convert layout tags into Xpath and calculate distances, and by doing so compare the layouts of web documents.
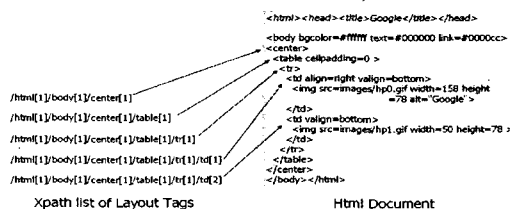


Fig. 2. Xpath expression of a layout tag

When there are two HTML document HDA, HDB, let us say the layout tag lists of the two document HDA, HDB layout as LA and LB. Then the distance (D) between the two documents is calculated using Equation (1).

$$D = \sum_i f_i(L_i) \qquad (1)$$

where Li , LA, LB and fi is a function to calculate the distance between LA and Li, which

is the ith layout of LB. The distance calculating function fi is as Equation (2).

$$f_i(L_i) = \begin{cases} 0 & L_i \in (L_A \cap L_B) \\ L_i & L_i \notin (L_A \cap L_B) \end{cases} \qquad (2)$$

In Equation (2) above, if the layout tags expressed in document HDA and HDB, namely, Xpaths are identical with each other 0 is returned, and if not a positive number is returned. When Xpaths are identical with each other the values of attributes such as bgcolor and align in the corresponding tags may be compared, but in this study if Xpaths are identical with each other the distance calculating function simply returns 0 [16]. Li is a parameter value. Using the equation above, this study calculates the distance between web document HDA and HDB. If the result is not larger than T, a specific threshold, the two documents are judged to have the same layout.

## 3. Parameter setting

The values of layout tags were set using a heuristic method as shown in [Table 2].

Table 2. The values of L parameters.

| Layout Tag | L Value |
|---|---|
| THEAD,TBODY,HR | 1 |
| other block tags | infinite |

Parameter values were set so that if Xpaths are not identical with each other the structures of the pages are judged to be different from each other except the corresponding tag is THEAD, TBODY or HR.

## IV. Implementation

This study implemented a HTML analysis too and HALD(HTML Analyzer for Layout Detection)

system based on the algorithm and parameter values described in Chapter 3 to evaluate the performance of the algorithm and the system. [Fig. 3] and [Fig. 4] show the outline of the HTML analysis tool system and a screen executing the system.
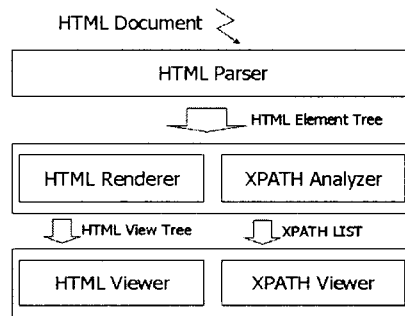


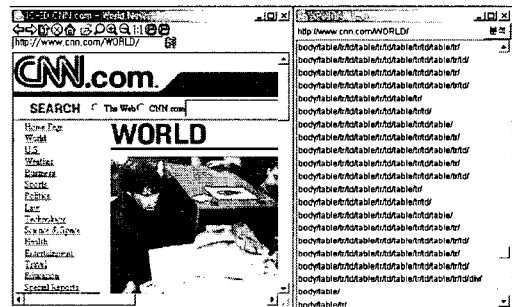Fig. 3. Outline of HTML analysis tool system.



Fig. 4. Screen executing HTML analysis tool.

HTML parser in [Fig. 3] parses documents and creates a HTML element tree similar to DOM tree. Using the element tree, the tree generator creates view trees and Xpath analyzer creates Xpath lists. The left of [Fig. 4] shows a screen executing the viewer and the right shows a screen executing Xpath analyzer.
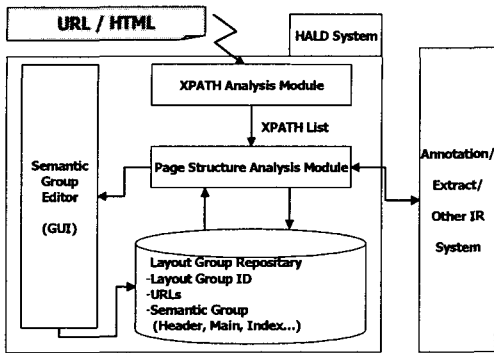
Fig. 5. Outline of HALD system

Xpath analyzer is included in Xpath analysis module within HALD system in [Fig. 5]. The structural diagram of HALD system to be implemented for data extraction system is as in [Fig. 5].

Xpath analysis module in [Fig. 5] abstracts Xpath lists from input HTML documents and page structure analysis module uses the lists as input and calculates distance to Xpath in storage. If the module judges that a page of the same structure is in the storage it assigns the corresponding semantic group to the HTML document and save it. If the module fails to find a page of the same structure, it assigns a new semantic group through semantic group editor and registers the page in the storage as a document of a new structure.

Systems in information retrieval areas such as data extraction can perform information retrieval effectively using HALD through API provided by the page structure analysis module.

## V. Performance evaluation

This chapter analyzes the structure of web pages using the algorithm and equations presented in Chapter 3 and evaluate the results. First, [Table

3] is the result of a test using 5 and 10 as threshold T.

Table 3. Result of document structure analysis(recall).

| Web Site | T = 5 | | T = 10 | |
|---|---|---|---|---|
| | classifying | recall | classifying | recall |
| www.cnn.com | 70 | 35 | 76 | 38 |
| news.yahoo.com | 30 | 15 | 44 | 22 |
| www.khan.co.kr | 147 | 73.5 | 148 | 74 |
| www.ytn.co.kr | 99 | 49.5 | 122 | 61 |
| www.nytimes.com | 110 | 55 | 137 | 68.5 |

According to the result, World News of CNN, Latest News of Yahoo News and Politics News of Kyunghyang Shinmun were all created using templates and each site used the same template to create its documents.

Two hundred articles in each site all used the same template, and the test was carried out by measuring recall values. According to the result of the test, when documents contain relevant article boxes inserted using <DIV>, <TABLE>, etc. the system often judged that their structures are different.
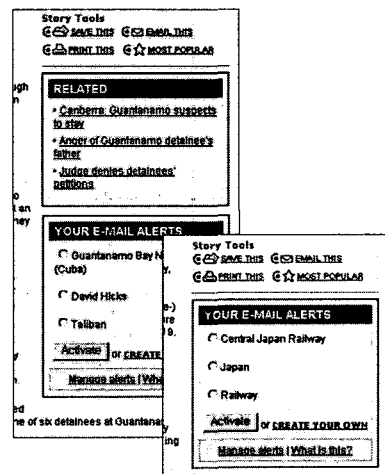


Fig. 6. Examples of relevant article boxes in CNN News site.

[Fig. 6] shows examples of pages that used the same template but were judged to have different structure as they had relevant article boxes. When a heuristic method of excluding relevant article boxes tagged with <DIV>, <TABLE>, etc. was applied and structure analysis was carried out again on the pages, the result was as shown in [Table 4], which shows that the application of the heuristic method improved the result to some degree.

### Table 4. Result of analysis applying a heuristic method.

| Web Site | T = 5 | | T = 10 | |
|---|---|---|---|---|
| | classifying | recall | classifying | recall |
| www.cnn.com | 96 | 48 | 104 | 52 |
| news.yahoo.com | 75 | 37.5 | 88 | 44 |
| www.khan.co.kr | 180 | 90 | 182 | 91 |
| www.ytn.co.kr | 114 | 57 | 131 | 60.5 |
| www.nytimes.com | 121 | 60.5 | 135 | 67.5 |

We extracted 200 articles from each of 5 sites to find out  satisfactory value of threshold without regard to the use of template and computed the precision value. The result was as shown in [Table 5].

### Table 5. Result of document structure analysis(precision).

| Web Site | T = 5 | | T = 10 | |
|---|---|---|---|---|
| | classifying | precision | classifying | precision |
| www.cnn.com | 80 | 78.7 | 110 | 72.7 |
| news.yahoo.com | 65 | 95.3 | 98 | 86.7 |
| www.khan.co.kr | 110 | 91.8 | 137 | 91.1 |
| www.ytn.co.kr | 88 | 86.3 | 114 | 80.7 |
| www.nytimes.com | 106 | 89.6 | 128 | 82.0 |

Based on [Table 4] and [Table 5],  [Table 6] shows  the  result  of  performance  evaluation according  to  the  threshold  value  using  dice coefficient formula[17].

### Table 6. Result of dice coefficient measures

| Web Site | T = 5 | T = 10 |
|---|---|---|
| www.cnn.com | 59.6 | 60.6 |
| news.yahoo.com | 53.8 | 58.3 |
| www.khan.co.kr | 90.8 | 91 |
| www.ytn.co.kr | 68.6 | 69.1 |
| www.nytimes.com | 72.2 | 74.0 |

According to the result of [Table 6], we can say that the performance of this method is most satisfactory at T = 10. However, we need to test and evaluate  the HALD system with more web documents in order to improve the performance and get the reliability of the system.

## VI. Conclusions and future works

Web services and designs are growing more complicated  for instance, a system showing user's opinion on articles and they are making automatic data processing more difficult. The development of many heuristic methods, however, is expected to resolve such difficulties.

Many web pages including news are created using  data  presentation  structures  called templates. Thus this study developed and tested a system that analyzes the structure of web pages and applies the results in various areas including information retrieval. This study developed first HTML analysis tool and Xpath analyzer and tested their performance. This study will be continued to develop Xpath analyzer and HALD system  of  upgraded  performance  through improving  the  algorithm,  applying  various parameter  values  and  discovering  additional heuristic methods and apply them to information extraction and processing in information retrieval areas.

## References

[1] K. Nagao, Y. Shirai and K. Squire, "Semantic Annotation and Transcoding : Making Web Content More Accessible," IEEE MultiMedia, Vol.8, pp.69-81, 2001.

[2] J. Myllymaki and J. Jackson, "Robust Web Data Extraction with XML Path Expressions," IBM Research Report, 2002.

[3] http://www.w3c.org/XML.

[4] C. Asakawa and H. Takagi, "Annotation-Based Transcoding for Nonvisual Web Access," Proceedings of ACM ASSETS 2000, pp.172-179, 2000.

[5] T. Sullivan and R. Matson, "Barriers to Use: Usability and Content Accessibility on the Web's Most Popular Sites," Proceedings of ACM CUU 2000, pp.139-144, 2000.

[6] A. W. Huang and N. Sundaresan, "Semantic Transcoding System to Adapt Web Services for Users with Disabilities," Proceedings of ACM ASSETS 2000, pp.156-163, 2000.

[7] A. Sahuguet and F. Azavant, "Building Light-Weight Wrappers for Legacy Web Data-Sources Using W4F," International Conference on Very Large Data Bases(VLDB), Edinburgh, Scotland, 1999.

[8] C. Allen, "WIDL:Application Integration with XML," WWW Journal2(4), 1997.

[9] http://www.w3c.org/TR/NOTE-WIDL.

[10] http://www.w3c.org/TR/xhtml2.

[11] http://www.w3c.org/TR/xpath.

[12] http://www.w3c.org/TR/xslt.

[13] J. Myllymaki, "Effective Web data extraction with standard XML technologies," The International Journal of Computer & Tele-communications, Vol.39, No.5, pp.635-644, 2002.

[14] M. Hori, G. Kondoh, K. Ono, S. Hirose, and S. Singhal, "Annotation-Based Web Content transcoding," Proceedings of the 9th International WWW Conference, pp.197-211, 2000.

[15] http://www.w3c.org/XML/Linking

[16] K. Fukuda, H. Takagi, J. Maeda, and C. Asakawa "Layout Group Extraction from Web Content for Effective Adaptation," IBM Research Report, 2002.

[17] Ian H. Witten and Eibe Framk, Data Mining, Morgan Kaufmann Publishers, 1999.

## 저 자 소 개

**선 복 근(Bok-Keun Sun)**　　　　　　정회원

- 1999년 호서대학교 컴퓨터공학과 졸업(학사)
- 2001년 호서대학교 벤처전문대학원 컴퓨터응용기술학과 졸업(공학석사)

- 2003~현재 : 호서대학교 반도체제조장비국산화연구센터 전임연구원, 현재 호서대학교 대학원 컴퓨터공학과 박사과정 재학

<관심분야> : 정보검색, 에이전트시스템, HCI

**한 광 록(Kwang-Rok Han)**　　　　　정회원

- 1984년 : 인하대학교 전자공학과 졸업(공학사)
- 1986년 : 인하대학교 대학원 정보공학 전공(공학석사)
- 1989년 : 인하대학교 대학원 정보공학 전공(공학박사)
- 1989~1991년 : 한국체육과학원 선임연구원
- 1991년~현재 : 호서대학교 컴퓨터공학부 교수
- 2001년~2002년 : ISI University of South California 방문연구원

<관심분야> : 정보검색, 자연어처리, 기계번역,