

# Analysis of HTML Structure in E-commerce Websites Using Tree Representation

Jose E. Ventura\* · Jeong-Sun Park\*

\*Dept of Industrial and Management Engineering, Myongji University

## 트리 표현을 사용하는 전자상거래 웹 사이트의 HTML 구조 분석

호세 벤토라\* · 박 정 선\*

\*명지대학교 산업경영공학과

### Abstract

개인화된 제품과 서비스에 대한 소비자의 요구는 성공적인 전자상거래 플랫폼을 기반으로 하고 있다. 성공적인 전자상거래 플랫폼을 개발함에 있어 자주 간과되고 있는 중요한 요소는 바로 웹 페이지의 HTML 구조이다. HTML 구조는 전자상거래 웹사이트의 속도와 랭킹을 결정짓는 기본적인 요소이다. 본 논문은 HTML 구조를 분석하기 위한 효율적이고 다소 생소한 시각화 기법을 제안하는데, 이러한 기법을 사용하여 개발자는 잠재적인 프로그래밍 오류와 개선 사항을 발견할 수 있다. 본 논문은 하나의 사례를 이용하여 제안된 기법을 더욱 구체화 시킨다.

**Keywords :** HTML Structure, E commerce, Tree Representation.

### 1. Introduction

HTML documents consist of elements which are constructed with tags that define the body of the page, paragraphs, links, tables, images, forms and others. The HTML hierarchy (DOM: Document Object Model) is a parent child sibling relationship.

An element that is directly above another element in the hierarchy is called parent of the element below it. The element below the parent is called the child. When two elements are equal in the hierarchy, they are known as siblings. The Figure 1 shows the basic hierarchy of a basic HTML document.

The goal of the analysis of HTML (DOM) structure is to get a semantic structure, a structure

that uses HTML tags for the purpose what they are designed to be used to, for example, <table> only for tabular data like a calendar, <div> for divisions of a document only for generic block level element you want to style. Having a clean semantic document has a lot of benefits. It's more logical for machines to understand. That means more compatibility as old browsers and other devices like mobiles phones or PDA's will understand it. Also, it will be more logical for other machine type things to read like search engine robots giving higher search ranks. For the analysis, an open source tool (Graphviz) is used to show good or bad structure and a bad example is modified as a good one.

† 교신저자: 박정선, 경기도 용인시 처인구 명지로 116 제1공학관 516호, 명지대학교 산업공학과

M · P: 031-330-6453, E-mail: jspark@mju.ac.kr

2011년 7월 18일 접수; 2011년 9월 24일 수정본 접수; 2011년 12월 7일 게재확정

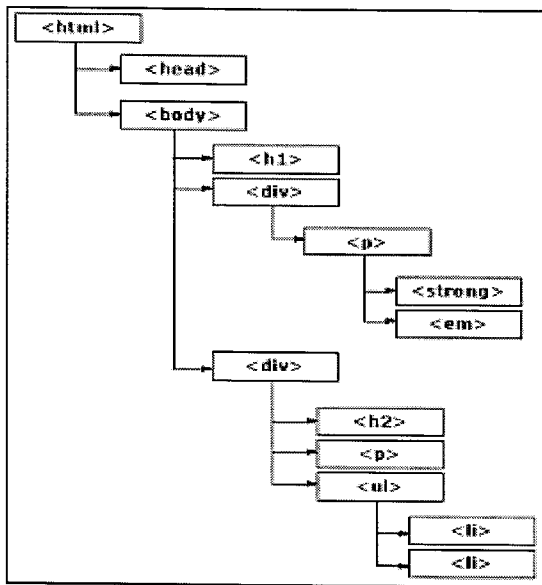


Figure 1. Hierarchy of a HTML document

## 2. Visualization of HTML Structure Using Tree Representation

One novel technique to analyze the HTML structure of an E commerce website is by visualizing it using tree representation where each HTML tag is represented by a node connected according to his HTML hierarchical structure.

The HTML structure can be converted in a tree graph using the following tools:

a) Perl script (html2gld.pl): This Perl script extracts the tags of a HTML documents turning in a GV file (ASCII text representation of a graph. This file describes a graph in terms of Nodes, Edges, Subgraphs, Attributes. Graphs, nodes and edges may have attributes that specify details of their appearance on the screen such as colors, sizes, shapes, etc.).

b) Graphviz (version 2.28, an open source graph visualization software): This program takes descriptions of graphs in a simple text language (GV file), and makes diagrams in useful formats, such as images and SVG for web pages, PDF or Postscript for inclusion in other documents.

Graphviz has many useful features for concrete diagrams, such as options for colors, fonts, tabular node layouts, line styles, hyperlinks.

The basic layout representations are:

- dot : "hierarchical" or layered drawings of directed graphs. This is the default tool to use if edges have directionality.
- neato and fdp : "spring model" layouts.
- Twopi : radial layouts.
- Circo : circular layout.

A basic HTML document is represented by Figures 2 and 3:

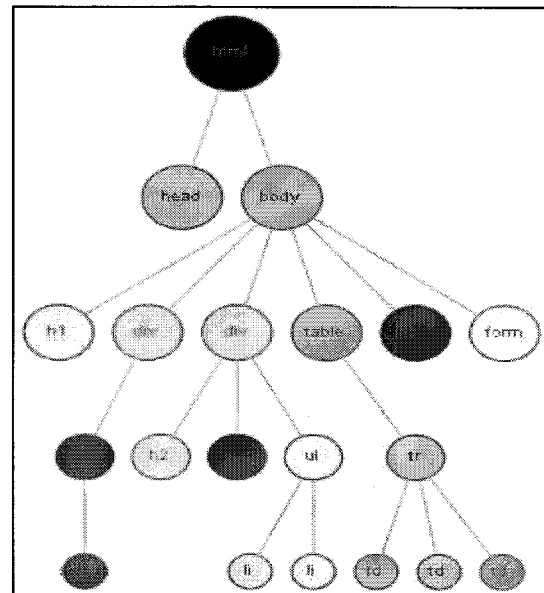


Figure 2. Visual representation of a basic HTML document in (hierarchical layout)

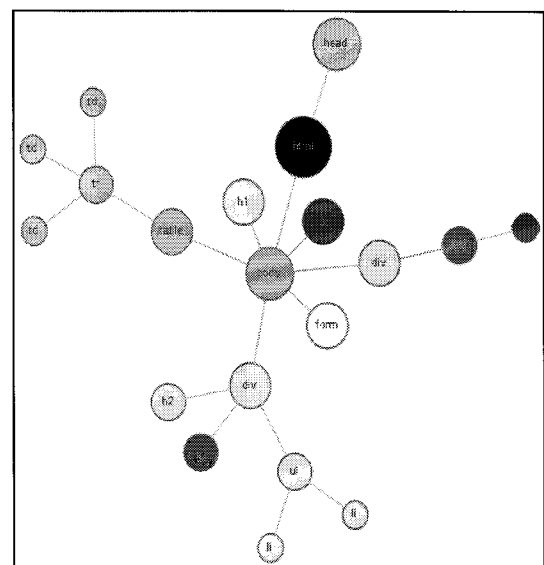


Figure 3. Visual representation of a basic HTML document in (spring layout)





The analysis of Figures 5 and 8, and other web pages in the English Gmarket website revealed that this basic improvement (substituting a div tag for three table tags) can be done in all website. It will increase the speed of website and also reduce considerably the size of the HTML code.

#### 4. Conclusions

The visualization of the HTML structure with this technique provides a good representation about how efficient and well-organized is the HTML programming of a web page.

The basic analysis presented in the study case and other more complex analysis can be done with this tool.

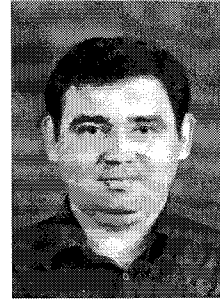
The great advantage of this tool is that developers can identify very fast potential errors in the HTML programming and improvements to have a clear semantic structure so that it will improve the speed and ranking of the E commerce website.

#### 5. References

- [1] Lee Underwood, "The HTML Hierarchy: Thinking inside the box", [http://www.htmlgoodies.com/beyond/article.php/3681551/The HTML Hierarchy Thinking Inside the Box.htm](http://www.htmlgoodies.com/beyond/article.php/3681551/The_HTML_Hierarchy_Thinking_Inside_the_Box.htm)
- [2] Websites as graphs, 2006, [http://www.aharef.info/2006/05/websites\\_as\\_graphs.htm](http://www.aharef.info/2006/05/websites_as_graphs.htm)
- [3] Oleg Burlaca, "HTML as graphs: the HTML2GDL application", 2009, [http://www.burlaca.com/2009/02/html2gdl\\_graphviz/](http://www.burlaca.com/2009/02/html2gdl_graphviz/)
- [4] Graph Visualization Software (GraphVis), <http://www.graphviz.org/>
- [5] Wikipedia, Gmarket, [http://en.wikipedia.org/wiki/G Market](http://en.wikipedia.org/wiki/G_Market)
- [6] <http://english.gmarket.co.kr>
- [7] <http://www.amazon.com>

#### 저 자 소 개

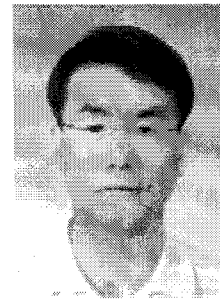
##### 호세 멘토라



ESEN 엘살바도르에서 학사, 명지대학교에서 석사학위를 취득하였고, 관심분야는 정보의 시각화, 시각적 데이터마이닝, MS 등이다.

주소: 경기도 용인시 처인구 남동 명지대학교 공학관 545호

##### 박 정 선



서울대학교에서 학사, 한국과학기술원에서 석사학위를 취득하였고, 미국 텍사스주립대학교 경영학박사를 취득하였으며, 한국전산원에서 선임연구원을 거쳐 현재는 명지대학교 산업경영공학과 교수로 재직 중이다.

주소: 경기도 용인시 처인구 남동 명지대학교 공학관 516호