

Boilerplate removal and Content Extraction from Dynamic web pages

By

Pan Ei San

Affiliation

University of computer studies (Yangon)

Correspondences

Address: Ziwaka hostel, Hlaing Twonship, Yangon, Myanmar

Email: paneisan1985@gmail.com

ABSTRACT

*Web pages not only contain main content, but also other elements such as navigation panels, advertisements and links to related documents. To ensure the high quality of web page, a good boilerplate removal algorithm is needed to extract only the relevant contents from web page. Main textual contents are just included in HTML source code which makes up the files. The goal of content extraction or boilerplate detection is to separate the main content from navigation chrome, advertising blocks, and copyright notices in web pages. The system removes boilerplate and extracts main content. In this system, there are two phases: Feature Extraction phase and Clustering phase. The system classifies the noise or content from HTML web page. Content Extraction algorithm describes to get high performance without parsing DOM trees. After observation the HTML tags, one line may not contain a piece of complete information and long texts are distributed in close lines, this system uses **Line-Block concept** to determine the distance of any two neighbor lines with text and **Feature Extraction** such as text-to-tag ratio (TTR), anchor text-to-text ratio (ATTR) and new content feature as Title Keywords Density (TKD) classifies noise or content. After extracting the features, the system uses these features as parameters in threshold method to classify the block are content or non-content.*

Keywords: content extraction, line-block, TKD, TTR, ATTR

1. INTRODUCTION

Today, the internet matures, thus the amount of data available continues to increase. The artifacts of this ever-growth media provide interesting new research opportunities that explore social interactions, language, art, and politics and so on. In order to effectively manage this ever-growing and ever-changing media, content extraction methods have been developed to remove extraneous information from web pages. Extracting useful or relevant information from Web pages thus becomes an important task. Also irrelevant information is contained in these Web pages. A lot of researches on WWW need the main contents of web pages to be gathered and processed efficiently. Web page content extraction technology is a critical step in many technologies. Content Extraction (CE) is just the technique to clean the documents from extraneous information and to extract the main contents.

Nowadays, web pages become much more complex than before, so CE becomes more difficult and nontrivial. Template based algorithms and template detection algorithms also perform poorly because of web page's structure being changed more frequently and web page's being generated dynamically. Traditionally, Document Object Model (DOM) based algorithms and vision based algorithms may get

better results but they always consume a lot of computing resource. Parsing DOM tree is a time consuming task. Vision based algorithms need to imitate browsers to render HTML documents, which will consume much more time. This system is implemented to remove noises or boilerplate based on Line-Block concept, content features and uses d threshold method to classify whether the block is content or not.

2. MOTIVATION

For human, the behavior can be done relatively fast and accurate because they can use their knowledge, visual representation and layout of the web pages to distinguish the main content from other parts. WWW rapidly grow as it is accessible for public use through the web browser. Typically, a modern web document comprises of different kinds of content. Elimination of noisy and irrelevant contents from web pages has many applications,

- web page classification, clustering, web featuring,
- proper indexing of search engines,
- efficient focused crawlers,
- Cell phones and PDA browsing.

Usually, apart from the main content blocks, web pages usually have such blocks as navigation bars, copyright and privacy notices, relevant hyperlinks, and advertisements, which are called noisy blocks. Modern web pages have largely abandoned the use of structural tags within a web page and adopted an architecture which makes use of style sheets and <div> or tags for structural information. Most current content extraction techniques make use of particular HTML cues such as tables, fonts, size and line, etc., and since modern web pages no longer include these cues, many content extraction algorithms have begun to perform poorly. One difference between our approach and other related work is that no assumption about the particular structure of a given webpage, nor does look for particular HTML cues. In our approach, the system uses the Line-Block concept to improve preprocessing step. And then, the system calculates the content features as Text-to-Tag Ratio (TTR), Anchor-Text to Text Ratio and the new feature Title Keyword Density (TKD). This state is called featured extraction phase. After feature extraction, the system use this features' values to classify the block is content or not by using threshold method. The system's objectives are followed:

- To develop a web content extraction method that given an arbitrary HTML document
- To extract the main content and discard all the noisy content
- To get high performance of noises detection without parsing DOM trees
- To decrease the consuming time of preprocessing step such as noise detection and classification of blocks.
- To enhance accuracy of information retrieval of Web data

Contributions: Four main contributions can be claimed in our paper:

1. To propose Extended Content Extraction algorithm that contains line block concepts, boilerplate detection and extraction of main content block
2. To reduce the web page's preprocessing time that used Line-Block concept.
3. To reduce the loss of important data by adding the new feature Title Keyword Density (TKD)
4. To retrieve the more important blocks that use threshold method.

The paper is structured as follows. After shortly reviewing related work in Section 3, we discuss background theory in section 4. Next, in section 5 we describe our proposed system in detail. In Section 6 we give our evaluation and experiments other CE algorithms. Finally we offer the conclusion with a discussion of further work in Section 7.

3. RELATED WORK

The term Content Extraction was introduced by Rahman et.al. [1] in which the authors describe a basic content extraction algorithm. Shortly thereafter Finn et al. [2] introduced the Body Text Extraction (BTE) algorithm wherein the authors extract content-text by identifying the single, continuous region which contains the most words and the least amount of HTML tags. Gottron [20] applied the Document Slope Curves (DSC) [3] extension to the BTE algorithm to create Advanced DSC (ADSC) in which a windowing technique is used to locate document regions in which word tokens are more frequent than tag tokens.

Debnath et al. developed the FeatureExtractor (FE) [10] and K-FeatureExtractor (KFE) [2] approaches based on block segmentation of the HTML document. Each block is analyzed for particular features like the amount of text, images, script code etc. Content text is extracted by selecting blocks which meet some criteria, e.g. most text. Gottron present an approach most similar to CETR by way of Content Code Blurring (CCB) [19], wherein content regions are identified by homogeneously formatted source code character sequences.

An attempt to combine different content extraction methods into one system was made by the Crunch framework [18, 19, and 17]. Crunch showed that a combination of different methods can provide better results than a single approach on its own. A more recent ensemble method called the CombineE framework [15] was recently developed to more easily configure ensembles of content extraction algorithms. Yet another approach is to induce a wrapper from labeled examples. One such approach was studied by Kushmerick [22], however this approach could not handle complex or in expected structures.

The Visual Page Segmentation (VIPS) [6] heuristically segments documents into a tree where the nodes are visually grouped blocks. The major problem with this approach is that the result of the VIPS algorithm does not label the nodes as content or non-content. The results presented in later sections show that if the best possible parameters are selected and a perfect mechanism is provided to label the nodes then VIPS can extract article text with a high degree of accuracy. However, there exists no such labeling mechanism; furthermore, VIPS must partially render a page in order to analyze it including retrieving all external style sheets, etc. Therefore, compared to other techniques, VIPS is very resource intensive.

Template detection algorithms [8, 11] are a different approach to content extraction in which collections of training documents based on the same template are used to learn a common structure. Specifically, Bar-Yossef et al. present an approach which automatically detects templates from the largest pagelet (LP) [3]. In general template detection algorithms find the main content by removing identical parts across all web documents. This is an accurate approach but has been found to be too time consuming and burdensome because a model must be built for each individual website and therefore for each sites multiple pages known to have the same templates are required. In the CleanEval content extraction competition only a few pages are available from the same site thus mandating a more general approach.

Wninger et al. introduced the Content Extraction via Tag Ratio (CETR) algorithm, a method to extract content text from diverse web pages using the HTML document's tag ratio [2]. The approach computes tag ratios on a line-by-line basic and then clusters the resulting histogram into content and noise area. This is a laconic and efficient algorithm, however vulnerable to the page's source code style changes.

4. BACKGROUND THEORY

There are non-informative parts outside of the main content of a web page. Navigation menus or advertisement are easily recognized as boilerplate, for some other elements it may be difficult to decide whether they are boilerplate or not in the sense of the previous definition. The **CleanEval guidelines** instruct to remove boilerplate types such as [10]

- Navigation

- Lists of links
- Copyright notices
- Template material, such as header and footers
- Advertisements
- Web spam, such as automated postings by spammers
- Forms
- Duplicate material, such as quotes of the previous posts in a discussion forum

4.1. Related Concept of Line

In this section, we describe the some concepts about the line of HTML source documents and content-feature that we use in our system.

A. Line

A HTML tag is a continuous sequence of characters surrounded by angle brackets like <html> and <a>. Hyperlink is one tag of HTML tag set. A complete Hyperlink tag has two markups: <a> as the open tag and as the close tag. A line is a HTML source code sequence from original HTML documents with texts and complete HTML tags (especially Hyperlinks tags). Anchor text of a line is the text between hyperlink tag's opening tag '<a>' and closing tags ''. Text of a line is the plain text of a line. It is all the continuous sequence characters between angle brackets '>' and '<'. If a line has no angle brackets, then all characters in this line are text of this line.

B. Define Line-Block

Line-Block is a line or some continuous lines, in which the distance of any two neighbor lines with text. Block means that it defines between open tag <>and end tag</>. In this paper, we define and use the important block tags as p, div, h1, h2 and so on.

C. Content Features

Definition (Text-to-tag ratio (TTR)): TTR is the ratio of the text length in the block is divided by the total sum of tags in this block.

$$TTR = \frac{text.length}{sum(tag)} \quad (1)$$

Definition (Anchor text-to-text ratio (ATTR)): ATTR is the ratio of the length of the anchor text is divided by the text length in the block.

$$ATTR = \frac{anchortext.length}{text.length} \quad (2)$$

Definition (Title Keyword Density (TKD)): A web page title is the name or heading of a Web Site or a Web Page. If there is more number if title words in a certain block, then it means that the corresponding block is of more importance.

$$TKD_k = 1 - \left[\frac{m_k}{m_k + \sum_{i=1}^{|m_k|} F(m_k)} \right] \quad (3)$$

Where, m_k is Number of Title keywords and $F(m_k)$ means Frequency of title keyword m_k in the block.

4.2. Selecting Content from Threshold method

Finally, we get three feature values for each line block and need the best the parameters as thresholds to remove the noise block. It is increasing of precision and a sharp decrease of recall. τ is threshold. $\tau = \lambda \sigma$, where λ is constant parameter and σ is standard deviation. Following steps are to find the mean value,

find the variance and find the standard deviation for calculates to get σ . Here, different web pages may have different kinds of content, so if we set the thresholds as constants it will lead to skew determinations. We assign the TTR threshold as 30 and ATTR's threshold as 0.2 and TKD' threshold as 2.

5. PROPOSED SYSTEM

In our proposed system include the four steps. They are defined as follows:

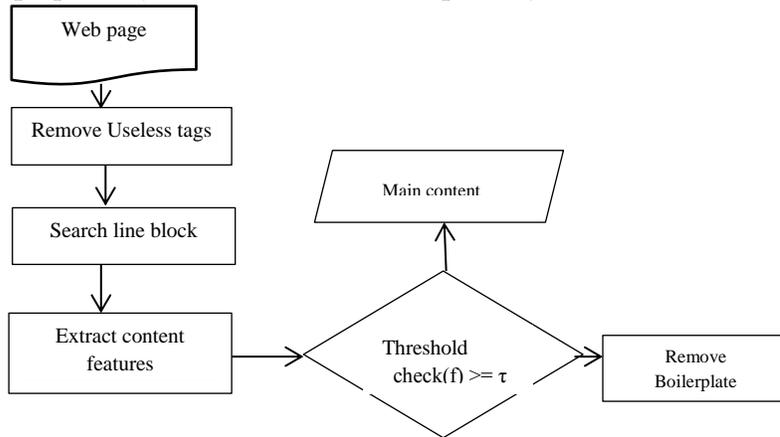


Figure1. Proposed system for content extraction

Step1. Preprocessing the web page tags

The tags filtered in this step, contains `<head>`, `<script>`, `<style>`, Remark and so on.

Step2. Define Line-Block

Line-block is a line or some continuous lines, which the distance of any two neighbor lines with text. The system reads line and makes the block using line-block concept. By sampling merging the lines, the system gets the line-blocks.

Step3. Feature Extraction

Next, the system calculates features for each block to determine whether they are content or not. TTR and ATTR are calculated as their formulas. For TKD, the system uses the title keywords in a block. Title Keyword Density (TKD) calculates to solve the loss of important information. There may be possibility that tag with less density also has the some important information. To remedy this the system a list of keywords from the title of the page and check if keyword density is greater than the threshold then the system add it the output block.

Step 4: Clustering Main contents or not

After calculating the content features, the system determines whether the block is content or not based on these features values. In this step, the system uses threshold methods to classify the main content or non - content and analyze the results. The threshold method uses standard derivation method. Threshold methods use three thresholds for TTR, ATTR and TKD. If $TTR > TTR$'s threshold and $ATTR < ATTR$'s threshold and $TKD \geq TKD$'s threshold then the block is main content. Otherwise, the block is noise block. Finally, the system extracts more accurately main contents.

5.1. Proposed Algorithm

- Input: D
- Output: mC
 - DF \leftarrow filter_useless_tags (D)
 - DB \leftarrow break_original_lines (DF)
 - DL \leftarrow get_lines (DB)
- LB \leftarrow get_line_blocks (DL)
- **For all block in LB do**
- f \leftarrow get_feature (block)
- **If threshold_check(f) \geq τ then**
- mC.append (block.text)
- **End for**

6. EXPERIMENTAL RESULTS

In this paper proposes a new content extraction algorithm. It differentiates noisy blocks and main content blocks. We present here the experimental results to testify the effect of algorithm. In many web page, so many links the main content that they can produce enough noise. At the same time, so many links in the text reduces the weight of the text, but Title Keyword density (TKD) effectively supplements the weight of main text. In this example the original web page has 48.4 KB to reduce when removing the boilerplate blocks; the testing page has only 8.82 KB. So, our proposed system can be reduced the storage space than original file size.



Figure2. Original Web page

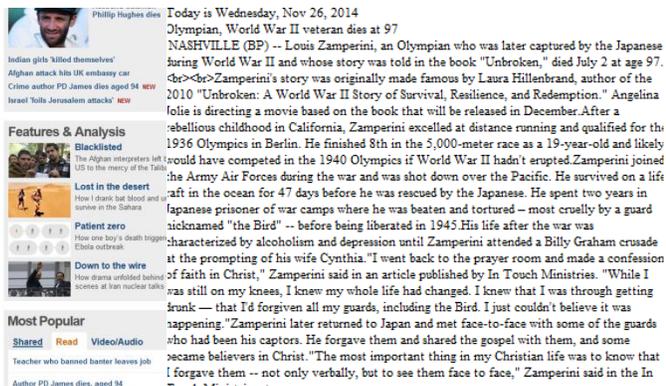


Figure3. Content Result

6.1. Data sets

The test data sets we use are from development and evaluation data sets from the CleanEval competition. They both hand-labeled gold standard set of main contents files; the amount of documents in each source are total of 606 web pages. In this dataset contains the following web site as BBC, nytimes and so on. CleanEval [10] is a shared task and competitive evaluation on the topic of cleaning arbitrary web pages. Besides extracting content, the original CleanEval competition also asked participants to annotate the structure of the web pages: identify lists, paragraphs and headers. In this paper, we just focus on extracting content from arbitrary web pages and use the 'Final dataset'. It is a diverse data set, only a few pages are used from each site, and the sites use various styles and structure.

7. CONCLUSION

The structures of webpages become more complex and the amount of data to be processed is very large, so Content Extraction (CE) remains a hot topic. We propose a simple, fast and accurate CE method.

We do not parse the DOM trees to get a high performance. We can get the main contents from HTML documents and research can be done on the original files, which widens the direction of CE research. However, our approach uses some parameters and depends on the logic lines of HTML source code. In the future plan, we continue to classify the web page and search engine for information retrieval.

REFERENCES

- A.F.R.Rahman, H.Alam and R.Hartono.(2001) Content Extraction from html documents. InWDA, pages 7-10
- A.Finn,N.Kushmerick, and B.Smyth.(2001) Fact or fiction: Content classification for digit libraries. In DELOS Workshop: Personalization and Recommender Systems in Digital libraries
- D.Pinto, M.Branstein. R.Coleman, W.B.Corft, M.King, W.Li, and X.Wei.(2002) Quasm: a system for question answering using semi-structured data. In JCDL, pages 46-55. ACM.
- D.Cai, S.Yu, J-R,Wen, and W-Y.Ma (2003). Extracting content structure for web pages based on Visual representation. In APWeb, volume 2642 of lecture Notes in Computer Science, pages 406-417. Springer
- D.de.Castro Reis, P.B.Golgher, A.S.da Silva amd A.H.F.Laender.(2004) Automatic web news extraction using tree edit distance. In WWW, pages 502-511. ACM
- H.Y.Kao, S.H.Lin, J.M.Ho and M.S.Chen.(2004) Mining web informative structures and contents based on entropy analysis. IEEE Trans. Knowl. Data Eng., 16(1):41-55.
- L.Chen, S.Ye and X.Li.(2006) Template detection for large scale search engines. In SAC, pages 1094-1098. ACM.
- L.Yi, B.Liu and X.Li. (2003) Eliminating noisy information in web pages for data mining. In KDD, pages 296-305, ACM.
- M.Marek, P.Pecina, and M.Spousta. (2007) Template detection through conditional random fields. In WAC3, 007.
- M.Baroni, S.Sharoff (2007)https://cleaneval.sigwac.org.uk/annotation_guidelines.html, Jan
- N.Kushmerick, (2000)Wrapper induction: efficiency and expressiveness. Artificial Intelligence, 118 (1-2): 15-68,.
- R.Cathet, L.Ma, N.Goharian and D.A.Grossman. (2003) misuse detection for information retrieval systems. In CIKM, apges 183-190. ACM.
- S.H.Lin and J.M.Ho. (2002) Discovering informative content blocks from web documents. In KDD, pages 588-593. ACM.
- S.Debnath, P.Mitra, and C.L.Giles.(2005) Automatic extraction of informative blocks from webpages. In SAC, pages 1722-1726. ACM.
- S.Debnath, P.Mitra, and C.L.Giles.(2005) Identifying content blocks form web documents. In ISMIS, volume 3488 of Lecture Notes in Computer Science, pages 285-293. Springer,
- S.Gupta, G.E. Kaiser, D.Neistadt, and P.Grimm.(2003) Dom-based content extraction of html documents. In WWW, pages 207-214
- S.Gupta, G.E. Kaiser and S.J.Stolfo. Extracting context to improve accuracy for html content extraction. In WWW (special interest tracks and posters), pages-1114-1115, ACM,2005.
- S.Gupta, G.E.Kaiser, P.Grimm, M.F.Chiang, and J.Starren.(2005) Automating content extraction of html documents. World Wide Web, 8(2):179-224.

T.Gottron. Content code blurring: A new approach to content extraction. In DEXA workshops [1], pages 29-33.

T.Gottron.(2007) Evaluating content extraction on html documents. In ITA, pages 123-132,

T.Gottron.(2008) Combining content extraction heuristics: the Combine System. In iiWAS, pages 591-595. ACM.

T.Weninger,(2010) W.H.Hsu and J.Han. CETR-Content Extraction via tag ratios. In proceedings of WWW'10, pages 971-980, New York, NY, USA.