

# HTML5のセクショニング要素自動付与システムの提案

今泉 智博<sup>†</sup> 齋藤 裕佑<sup>††</sup> 西山 裕之<sup>†</sup>

東京理科大学理工学部経営工学科<sup>†</sup>

東京理科大学大学院理工学研究科<sup>††</sup>

## 1 はじめに

1990年代後半以降、情報の発信及び収集のための手段としてインターネットが普及している。総務省によれば、2011年末の国内インターネット利用者人口は9,610万人で、普及率は79.1%である[1]。また別の調査では一日に30分以上インターネットを利用する人が全体の85%に達しているとの結果がある[2]。またインターネットを構成する主要な要素にWorld Wide Web(以下Web)がある。1993年に世界全体で130件しか存在しなかったウェブサイトの数は、2012年には約7億件にまで増加した。

そこで近年Webページから主要コンテンツを抽出するニーズが高まっている。ここでの主要コンテンツは、ニュースサイトやブログの投稿ページにおける本文部分のことを指す。その理由として、既に述べたように日々大量の情報がWeb上に蓄積されていく一方で、1人のユーザーが到達できる情報量には限界があるという背景が挙げられる。例えばニュースやブログの本文を抽出することで、内容を要約し短時間で把握が可能となる。また検索エンジンを利用する際、広告やナビゲーションなどの不要なコンテンツを検索結果の中に反映することがなくなり、検索結果の質の向上が期待される。さらに通信時の負担を削減するために、主要なコンテンツのみを優先して取得することで、Webページの表示待ち時間を削減することが可能である。

主要コンテンツ抽出手法についての既存の研究では、吉田ら[3]が同一Webサイトにおいて他のWebページには共通して現れない部分を本文箇所として推定する手法を提案している。またPappasら[4]は多様なWebページで主要コンテンツ部分を推定するため、独自アルゴリズムによる手法を提案している。

## 2 HTML5とセクショニング要素について

Webページを構成するのはHTMLタグと呼ばれる識別子によって装飾された文章である。本研究では主要コンテンツの推定に有用な情報として、HTML5のセクショニング要素と呼ばれるHTMLタグに注目する。HTML5とは、Web利用者にとってのユーザビリティ向上を目的として、2014年10月28日に行われたHTMLの5回目の大幅な仕様変更である。またセクショニング要素はHTML5にて新しく定義されたHTMLタグで、Web

ページ内の文章に対する意味付けを行うことを目的としている。具体的には次の4つのHTMLタグを指す。

- article タグ: Webページの中で完全もしくは自己完結した構造を表す。これは、フォーラムの投稿、雑誌や新聞の記事、ブログのエントリ、ユーザーの投稿コメントなどが考えられる。つまりarticleタグは、主要コンテンツそのものを表すHTMLタグと言える。
- section タグ: 文書またはアプリケーションの一般的なセクション(主題を表すコンテンツのグループ)を表す。通常1つのarticleタグが0個以上のsectionタグを内包している場合が多い。
- nav タグ: 他のページへのリンク、またはナビゲーションリンクをもつ部分を表す。
- aside タグ: 周囲のコンテンツとわずかに関連し、かつそのコンテンツから分離すると見なすことができる部分を表す。具体的にはリード文またはサイドバーのような部分や広告などを指す。

上記のように、セクショニング要素は主要コンテンツの推定に大いに役立つと考えられる。しかし正式勧告されたばかりの仕様のため、セクショニング要素を持たないWebページも未だ数多く存在する。それらを手動で更新、修正するのは作業者にとって負担である。

そこで本研究では、HTML5への更新作業時の負荷軽減とWebにおけるユーザビリティの向上を目的とし、セクショニング要素を持たないWebページに対してそれらを自動付与するシステムを提案する。

なお以降でセクショニング要素を含むWebページをHTML5ページ、それらを持たないページを非HTML5ページと表現する。また本研究で利用するセクショニング要素は、4つのうちarticleタグのみに限定する。

## 3 システム設計

本システムの構成は図1の通りである。具体的な処理の流れは以下の通りである。

1. 予め指定したURLのHTML5ページと非HTML5ページのページデータを取得する。
2. 取得したWebページは用意したデータベースへテキストデータとして保存する。
3. 各HTML5ページについてDOMツリーを生成する。
4. セマンティック要素のノードと、その兄弟ノードや祖先の兄弟ノードについて特徴量を算出する。例えば図2のようなDOMツリーの場合、緑のノードを正事例、赤のノードを負事例とする。

Proposal for a system to automatically append HTML5 sectioning elements to Web pages

Tomohiro Imaizumi<sup>†</sup>, Yusuke Saito<sup>††</sup>, Hiroyuki Nishiyama<sup>†</sup>

<sup>†</sup>Faculty of Science and Technology, Tokyo University of Science

<sup>††</sup>Graduate School of Science and Technology, Tokyo University of Science

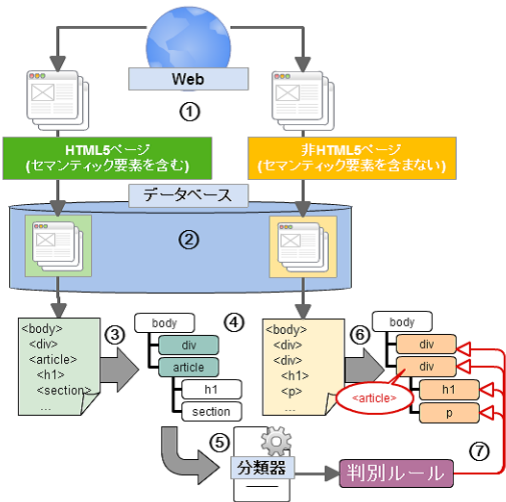


図1 セクショニング要素自動付与システムの構成図

5. 前述の特徴量から、ノードがセマンティック要素のノードか否かを判別するルールを分類器で生成する。
6. 非HTML5ページでもDOMツリーを生成する。またその後、それぞれの特徴量も求める。
7. 非HTML5ページの全てのDOMノードに対して5で求めたルールを適用し、特徴量にもとづいてセマンティック要素のノードとして分類されたノードに対応するタグをarticleに置き換える。

分類器についてはSVMやナイーブベイズを用い、分類器の違いによる性能の比較も行。また上記の4で求める特徴量を求める手順は以下のとおりである。まず次のようなノード  $N$  に対するスコア  $X_{score}$  を定義する。

- $N_{text}$ : ノード  $N$  に含まれる文字数
- $N_{timeinfo}$ : ノード  $N$  が時間表現を含んでいるか (含んでいる場合1, 含まない場合0)
- $N_{puncs}$ : ノード  $N$  に含まれる句読点数
- $N_{tag(E)}$ : ノード  $N$  の子要素に現れるノード  $E$  の数

上記のスコアを注目するノード、そのノードが含まれるDOMツリーの根ノード  $R$  についてそれぞれ求める。そして、 $N_{time}$  を除く各スコアについて  $N_{text}/R_{text} * 100$  のように、 $R$  のスコアに対する割合を100分率にし、これをノード  $N$  の特徴量とする。なお  $N_{time}$  についてはそのままの値を特徴量として用いる。

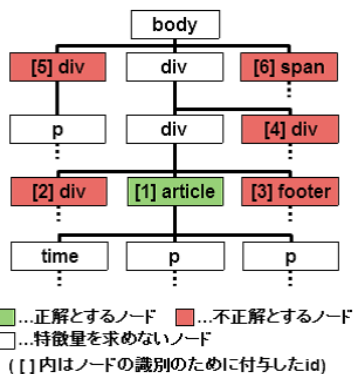


図2 DOMツリーで特徴量をカウントするノード

表1 各ノードのスコア

id	正解	$N_{text}$	$N_{timeinfo}$	$N_{puncs}$	$N_{tag(p)}$	$N_{tag(div)}$
1	T	4752	1	148	35	12
2	F	742	0	38	3	0
3	F	15	0	0	0	0
4	F	22	0	0	0	0
5	F	44	0	0	1	0
6	F	0	0	0	0	1

表2 得られる特徴量の例

id	正解	$N_{text}$	$N_{timeinfo}$	$N_{puncs}$	$N_{tag(p)}$	$N_{tag(div)}$
1	T	70.8	1	65.5	89.7	70.6
2	F	11.1	0	16.8	7.7	0.0
3	F	0.2	0	0.0	0.0	0.0
4	F	0.3	0	0.0	0.0	0.0
5	F	0.7	0	0.0	2.6	0.0
6	F	0.0	0	0.0	0.0	5.9

模式的なDOMツリーを図2に示す。今このDOMツリーのうちスコアを計算した結果が表1であったとする。またこのDOMの根ノードのスコアが  $R_{text} = 6710$ ,  $R_{timeinfo} = 1$ ,  $R_{puncs} = 226$ ,  $R_{tag(p)} = 39$ ,  $R_{tag(div)} = 17$  であった場合、これらから表2の特徴量を得る。

このようにして得られた特徴量からセクショニング要素を示すノードと、それ以外のノードの分類ルールを生成する。分類ルールの生成にはサポートベクターマシンやナイーブベイズ分類器などを用いる。

その後生成したルールは非HTML5ページから生成したDOMツリーの全てのノードに対して適用し、特徴量により各ノードをそれぞれセクショニング要素とそれ以外に分類することで、最終的にセクショニング要素に相当するノードの推定を行う。

#### 4 おわりに

本研究では、HTML5への更新作業負荷軽減とWeb利用者にとってのユーザビリティの向上を目的として、HTML5のセクショニング要素を自動付与するシステムを設計した。提案したシステムでは、セマンティック要素を持たないページに対し、適切な位置にセマンティック要素を付与することに成功した。本研究では付与できるセクショニングタグがarticleタグに限定されているが、今後は他の3つのタグについても同様に自動付与できるシステムの作成を目指す。

#### 参考文献

- [1] 総務省 平成24年版 情報通信白書, 2012
- [2] 一般財団法人経済広報センター 情報源に関する意識・実態調査報告書, 2013
- [3] 吉田光男, 山本幹雄. (2009). 教師情報を必要としないニュースページ群からのコンテンツ自動抽出. 日本データベース学会論文誌, 8(1), 29-34.
- [4] Pappas, N., Katsimpras, G., Stamatatos, E. (2012, September). Extracting informative textual parts from web pages containing user-generated content. In Proceedings of the 12th International Conference on Knowledge Management and Knowledge Technologies (p. 4). ACM.