

# 日本語 Web ページを対象としたタイトルタグの傾向と分析

秋岡明香\*  
電気通信大学/早稲田大学

下山剛司†  
早稲田大学

村岡洋一‡  
早稲田大学

## 1 はじめに

ブログの普及やネットワークサービスの大衆化に伴い、HTML 文書が爆発的に増えている。例えば NetCraft 社の Web Server Survey によると、2007 年 12 月時点で世界の Web サーバ数は 1.5 億台を超えるとされている [1]。Lawrence らの調査によると、1999 年 2 月現在に発信されていた Web ページは約 8 億ページであり [2]、NetCraft 社の調査では当時の Web サーバ数が約 430 万台であることから、単純な比例計算によると、現在では約 280 億ページが公開されている計算になる。このような背景を踏まえ、HTML 文書を知識データベースとして利用する研究も盛んに行なわれているが、膨大なデータをいかに効率良く処理し必要な情報のみを取り出すかは、重要な研究課題である。

一方で HTML にはタイトルタグが定義されており、そのページで最も重要なキーワードや、そのページの要約を含むことが期待されている。したがって、タイトルタグに含まれるキーワードの分析を行ない傾向を正しく把握することは、昨今の Web ページの傾向を把握するために役立つだけでなく、Web コンテンツから必要な情報のみを取り出す知識フィルタを作成する上でも多いに有用である。そこで本稿では、特に日本語で記述された 1 億 URL を超える Web ページを対象に、タイトルタグの傾向と分析を行なった結果を報告する。

## 2 タイトルタグの全体的な傾向調査

本稿では、2007 年 11 月に収集した、日本語で記述された HTML 文書 110,925,216 件のタイトルタグに対して、MeCab [3] と Lingua::JA::Summarize [4] による形態素解析を適用し、抽出したキーワード 15,649,393 (重複を含めて 349,827,579) 語の全体的な傾向を示す。なお、Lingua::JA::Summarize での抽出単語は最大 5 とした。

\*Sayaka Akioka, The University of Electro-Communications/Waseda University

†Takeshi Shimoyama, Waseda University

‡Yoichi Muraoka, Waseda University

図 1 は、タイトルタグに含まれるキーワード数の集計結果である。タイトルタグに含まれるキーワード数で最も一般的なのは 5 語であり、このようなページは全体の約 27% を占め、次いで多かったのは 2 語の約 20% であった。また、タイトルタグから有効なキーワードを抽出できなかった URL が約 1% 存在した。

図 2 に、キーワードの出現回数の分布を示す。また、表 1 は、出現頻度が上位 10 位までのキーワードの一覧である。表 1 の上位 8 語は出現頻度が 100 万回を超える。しかし、このように頻出するキーワードは極めて少数であり、抽出したキーワードのうち 93% は出現頻度が 10 回未満である。またこれらの頻出キーワードには一般的なキーワードが多く、コンテキストマイニング等の目的において重要な単語であるとは言い難い。さらに、約 6% のキーワードは出現頻度が 10 回以上 100 回未満であり、全体として 99% の語は出現頻度が 100 回未満である。

なお本稿で抽出したキーワードは厳密な意味での単語に区切られたものではない。例えば、「連続放火」、「連続殺人犯」、「連続掲載」といったフレーズは、各々独立のキーワードとして扱われており、このような単語の切り出し方が出現頻度 100 回未満のキーワードを増加させた一因であると考えられる。しかし HTML を分類・分析する上で、これらの語を関連づけることはあっても、全く同種の文書として扱われるケースは稀であると推測され、ここでは大きな問題はないと考える。

## 3 出現キーワードに関する考察

本章では、第 2 章で切り出したキーワードのうちで出現頻度が 100 回以上のキーワード 213,290 件について、いくつかの事件等と関連を考察する。

本稿で用いたデータは 2007 年 11 月現在での収集データであるが、この時期には防衛省事務次官と山田洋行の癒着疑惑が発覚した。そこでこの事件に関連するキーワードに注目すると、守屋 (1022 回)、防衛省 (723 回)、山田洋行 (439 回)、次官問題 (432 回)、守屋氏 (248

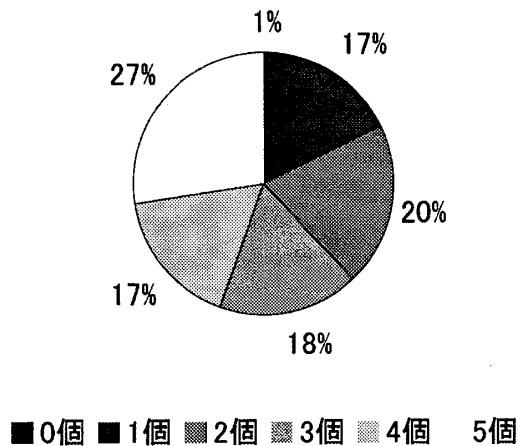


図 1: タイトルタグに含まれるキーワード数の傾向

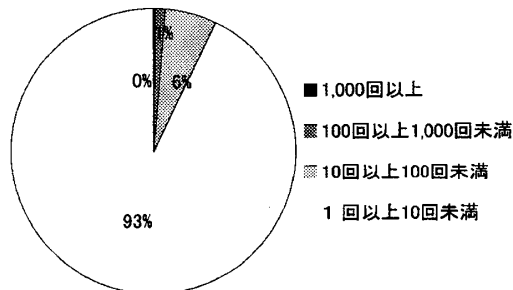


図 2: キーワードの出現回数の傾向

回), 東京地検 (108 回) などが含まれていた。

一方で約半年前の 2007 年 6 月には, ミートーホープ社の食肉偽装事件が起こっているが, 2007 年 11 月時点の収集データにもこの事件に関連すると思われる偽装 (879 回), 食肉 (767 回), ミンチ (250 回), ミートーホープ (156 回) などが含まれていた。

これらの他に, 2007 年流行語大賞である「ハニカミ王子」は 203 回, 2006 年流行語大賞である「イナバウアー」は 337 回含まれていた。

まとめると, 社会現象となる流行語や大きな事件に関連する単語は, 100 回以上の出現頻度で概ね捉えることができると予想されるが, このような語は本稿で抽出したキーワードの 1% 程度であり, 極めて少数部分であると言える。また, 今回得た知見を知識フィルタリング等に応用するためには, この 1% 程度の文書を如何に効率良く切り分けるかが大きな課題になると考えられる。さらに, 今回は 2007 年 11 月のみの収集データを用いたが, 時間経過とキーワードの出現頻度の変化の動向にも着目するべきであり, さらに大規模なデータの解析が必要であると思われる。

表 1: 出現頻度が上位のキーワード 10 件

順位	キーワード	出現回数
1	修正	7,169,906
2	登録データ	7,100,549
3	削除	6,006,136
4	パスワード認証画面	2,411,132
5	amp	2,253,723
6	管理者	1,831,204
7	blog	1,328,197
8	通知画面	1,247,247
9	search	985,481
10	日記	979,606

## 4 まとめと今後の課題

本稿では, 日本語で記述された 1 億 URL を超える HTML 文書を対象に, タイトルタグに出現するキーワードについての傾向調査と分析を行なった。今後は, 出現キーワードの時系列分析等のさらに詳細な調査を行ない, 知識フィルタ開発等へ応用する予定である。

## 謝辞

本研究は文部科学省リーディングプロジェクト e-Society 「基盤ソフトウェアの総合開発」プロジェクトのサブプロジェクト「インターネット上の知識集約を可能にするプラットフォーム構築技術」の支援により実施された。

## 参考文献

- [1] Netcraft: December 2007 Web Server Survey, [http://news.netcraft.com/archives/2007/12/29/december\\_2007\\_web\\_server\\_survey.html](http://news.netcraft.com/archives/2007/12/29/december_2007_web_server_survey.html).
- [2] S. Lawrence and C. K. Giles, "Accessibility of Information on the Web", Science, Vol. 400, pp. 107 - 109, 1999.
- [3] MeCab: Yet Another Part-of-Speech and Morphological Analyzer, <http://mecab.sourceforge.net/>.
- [4] Lingua::JA::Summarize, <http://search.cpan.org/~kazuho/Lingua-JA-Summarize-0.07/>.