

## リンクの共起関係を用いた URL 関連度計算手法の一検討

3X-04

杉崎 正之 牧野 俊朗 稲垣 博人

日本電信電話株式会社 NTT サイバーソリューション研究所

## 1 はじめに

インターネットなどに代表されるコンピュータネットワークの普及により、大量のテキスト情報が不特定多数に対して公開されている。その代表格が HTML ファイルであり、分散して存在する多量の HTML ファイルの中から欲しい情報を容易に取り出せるようにするため、検索や分類、ナビゲーション等、多くの方式が研究されている。これらを効果的に行なうには、単語と文書、あるいは文書と文書の関連付けが重要である。本稿では、ハイパーリンクを用いることにより、Web ページを関連付ける手法について検討する。

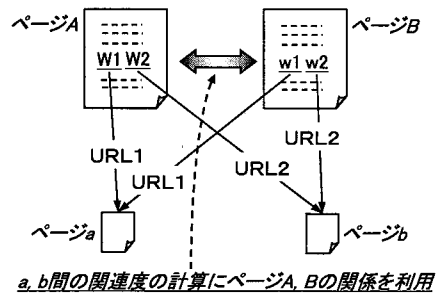


図 1: ハイパーリンクの共起の考え方

## 2 リンクの共起に基づく関連付け

## 2.1 概要

基本的な考え方は、単語間の関連度を求めるのに文内での共起関係を利用するのと同様な方法である。すなわち、同一の Web ページに複数のハイパーリンクが存在するとき、それらのリンク先の Web ページどうしは互いに関連しているとみなす (図 1 の「ページ a, b」が該当)。

WWW の特性として、自分の興味に合致し、かつ何度もアクセスするページをリンク集として作成公開しているユーザが多く存在する。その場合、同じ分野へのリンクはまとまっていることが多く、さらにそれらが恒常的な情報の場合には多くの「リンク集」Web ページで内容が書き換えられることがない。また、恒常的な情報でなくても企業間の提携情報などが Web ページで公開される場合に、関連 URL として各企業へのリンクが同一ページ内に記述されたりする。

このような特性を生かし、同一ページ内での共起情報を利用することで、関連のある URL が抽出できることを確認した [1]。さらに、Web ページ内において、出現位置情報を利用して共起を調べる範囲を設定することにより、検索サイトのように複数の Web ページに大量に存在する可能性のあるリンク先との関連度が大き

くなってしまう問題 (トピックドリフト問題) もある程度解消できることも分かった。

しかし、上記の手法の問題点として、もともと関連のない 2 つの URL を同一ページ内に距離的に近く記述したページを大量に作成しておくことで、意図的に関連度を高くすることが可能となってしまう (図 1 でいえば「ページ A」「ページ B」に当たる)。また、意図的でないにしても「問い合わせ先」や「ヘルプページ」、「トップページ」へのハイパーリンクが記述された大量の Web ページを所有するサイトから単純に共起関係を抽出すると、偏った関連度が計算されてしまう。

そこで、それらを避けるために 2 つの URL 間の関連度を定義している Web ページどうしの情報の利用を検討する。

## 2.2 改良手法

まず、2 つの URL ( $u, v$  とする) に対して次のような関数を定義する。

$$f^s(u, v) = \begin{cases} 1 & \text{if } u, v \text{ が同一ページ上の} \\ & s \text{ バイト以内} \\ 0 & \text{それ以外} \end{cases}$$

従来は、関数  $f^s(u, v)$  の単純な総和とした。今回は、 $a, b$  が共起している Web ページ、すなわち  $f^\infty(a, b) = 1$  となる Web ページの集合を  $S(a, b)$  として、集合  $S(a, b)$  の

処理したWebページ数 : 1,035,753  
 抽出されたリンク数 : 1,016,684  
 関連度が確認されたURL : 235,148  
 関連度が確認されたURLの組: 4,393,662

図 2: 抽出結果

要素数を  $|S|$  とし、ユニークなドメイン数を  $u\_dom(S)$  として、2つの URL( $a, b$  とする)の関連度  $rel(a, b)$  を、

$$rel(a, b) = \sum f^s(a, b) * \frac{u\_dom(S)}{|S|} \quad (1)$$

とする。上記の式は、互いに異なる多くのページから Web ページから参照されている場合に値を大きくすることで、上記の問題点を解消しようとしている。

### 3 実験と考察

実際に Web ページを収集し実験を行なった。2001 年 11 月に収集した日本語の Web ページのうちの約 100 万ページを利用した。同一ドメイン内の相対リンクや CGI へのリンクは削除することにし、抽出時の  $s$  の値は 800 とした。抽出結果を図 2 に示す。

関連度の高い URL の抽出結果の違いについて検討する。図 3 をみると、従来手法では「Yahoo! Japan (<http://www.yahoo.co.jp>)」と企業間で関連があるサイト (2, 3) や自社内のサービスの URL(7) との関連度が高かった。一方、改良手法ではいわゆる「検索サイト」と呼ばれる URL が上位を占めた。また、図 4 をみると、従来手法では「goo (<http://www.goo.ne.jp>)」を運営している会社のサイト (1) や自社内のサービスの URL(3, 4, 6) との関連度が高かった。一方、改良手法では図 3 の時と同様な URL が上位を占めた。

今回抽出元となった Web ページ内には図 3 や図 4 の従来手法に出現したサイト内から取得したページも含まれており、運営会社へのリンクや同社が行なっている別サイトでのサービスとの共起が多くなるという問題が実際に生じている。が、改良手法を用いることで、多くの人が関連があると考えている Web ページを抽出できていることが確認できた。

また、図 5 は「日本書籍出版協会 (<http://www.books.or.jp>)」に対する関連 URL を関連度の高い順に並べたものであるが、従来手法と改良手法での出力に大きな違いは生じなかった。これは、関数  $f^s(a, b) = 1$  となる Web ページの集合においてそれぞれのページが別のドメイン上に存在しており、この場合は今回の改良による悪い影響が出なかった。

従来手法	改良手法
1. <a href="http://www.goo.ne.jp/">http://www.goo.ne.jp/</a>	1. <a href="http://www.goo.ne.jp/">http://www.goo.ne.jp/</a>
2. <a href="http://www.zdnet.co.jp/">http://www.zdnet.co.jp/</a>	2. <a href="http://www.infoseek.co.jp/">http://www.infoseek.co.jp/</a>
3. <a href="http://www.softbank.co.jp/">http://www.softbank.co.jp/</a>	3. <a href="http://www.lycos.co.jp/">http://www.lycos.co.jp/</a>
4. <a href="http://www.infoseek.co.jp/">http://www.infoseek.co.jp/</a>	4. <a href="http://www.excite.co.jp/">http://www.excite.co.jp/</a>
5. <a href="http://www.lycos.co.jp/">http://www.lycos.co.jp/</a>	5. <a href="http://www.yahoo.com/">http://www.yahoo.com/</a>
6. <a href="http://www.excite.co.jp/">http://www.excite.co.jp/</a>	6. <a href="http://www.google.com/">http://www.google.com/</a>
7. <a href="http://chat.yahoo.co.jp/">http://chat.yahoo.co.jp/</a>	7. <a href="http://navi.ocn.ne.jp/">http://navi.ocn.ne.jp/</a>

図 3: 「<http://www.yahoo.co.jp/>」と関連がある URL

従来手法	改良手法
1. <a href="http://www.ntx.co.jp/">http://www.ntx.co.jp/</a>	1. <a href="http://www.yahoo.co.jp/">http://www.yahoo.co.jp/</a>
2. <a href="http://www.yahoo.co.jp/">http://www.yahoo.co.jp/</a>	2. <a href="http://www.infoseek.co.jp/">http://www.infoseek.co.jp/</a>
3. <a href="http://community.goo.ne.jp/">http://community.goo.ne.jp/</a>	3. <a href="http://www.lycos.co.jp/">http://www.lycos.co.jp/</a>
4. <a href="http://shop.goo.ne.jp/">http://shop.goo.ne.jp/</a>	4. <a href="http://www.excite.co.jp/">http://www.excite.co.jp/</a>
5. <a href="http://www.infoseek.co.jp/">http://www.infoseek.co.jp/</a>	5. <a href="http://www.nikkei.co.jp/">http://www.nikkei.co.jp/</a>
6. <a href="http://channel.goo.ne.jp/">http://channel.goo.ne.jp/</a>	6. <a href="http://japan.infoseek.com/">http://japan.infoseek.com/</a>
7. <a href="http://www.nikkei.co.jp/">http://www.nikkei.co.jp/</a>	7. <a href="http://www.google.com/">http://www.google.com/</a>

図 4: 「<http://www.goo.ne.jp/>」と関連がある URL

### 4 今後の課題

今回は関連度の計算に「ユニークなドメイン数」を用いたが、同一ドメインで別の個人の Web ページから張られているリンク情報がうまく利用されておらず、今後「ドメイン単位」ではなく「サイト単位」での関連度の計算を検討したい。また、 $s$  の値による URL のカバー率への影響や、HITS[2] などの再帰的な手法との比較検討を行なっていきたい。

### 参考文献

- [1] 大久保, 杉崎, 田中: リンクの共起関係を用いた Web ページ分類方式の検討情処第 59 回全大 (3), pp.81-82, 1999.9
- [2] J.Kleinberg: Authoritative sources in a hyperlinked environment, Proc. of 9th ACM-SIAM Symposium in Discrete Algorithms, pp.668-677,1998

従来手法	改良手法
1. <a href="http://www.maruzen.co.jp/">http://www.maruzen.co.jp/</a>	1. <a href="http://www.maruzen.co.jp/">http://www.maruzen.co.jp/</a>
2. <a href="http://bookweb.kinokuniya.co.jp/">http://bookweb.kinokuniya.co.jp/</a>	2. <a href="http://bookweb.kinokuniya.co.jp/">http://bookweb.kinokuniya.co.jp/</a>
3. <a href="http://webcat.nacsis.ac.jp/">http://webcat.nacsis.ac.jp/</a>	3. <a href="http://www.trc.co.jp/trc-japa">http://www.trc.co.jp/trc-japa</a>
4. <a href="http://www.trc.co.jp/trc-japa">http://www.trc.co.jp/trc-japa</a>	4. <a href="http://webcat.nacsis.ac.jp/">http://webcat.nacsis.ac.jp/</a>
5. <a href="http://www.ndl.go.jp/">http://www.ndl.go.jp/</a>	5. <a href="http://www.ndl.go.jp/">http://www.ndl.go.jp/</a>
6. <a href="http://www.amazon.com/">http://www.amazon.com/</a>	6. <a href="http://www.amazon.com/">http://www.amazon.com/</a>
7. <a href="http://www.asahi.com/">http://www.asahi.com/</a>	7. <a href="http://www.asahi.com/">http://www.asahi.com/</a>

図 5: 「<http://www.books.or.jp/>」と関連がある URL