

キーワードに着目したブログ空間のロコミパターン抽出

平井 佑一*1
Yuichi Hirai

熊野 雅仁*1
Masahito Kumano

木村 昌弘*1
Masahiro Kimura

1. はじめに

近年、ブログは Web 上で誰もが容易に情報を発信できるメディアとして急速に普及し、ブログ空間は新たなコミュニケーションの場として活性化しており、ブログユーザー同士の情報伝搬現象を分析することが注目されている。本研究では、ブログ空間におけるロコミ的な情報伝搬現象分析システムの構築を目指し、キーワードに着目した分析システムを試作した。本システムでは、ブログ文書ストリームデータが与えられたとき、まず、Gruhl らの手法[1]を用いてそのデータからキーワード群を抽出する。次に、各キーワードのロコミ的伝搬パターンを、ブログネットワーク情報とそのキーワードが記載されていたブログ文書の作成時間情報を用いて抽出する。さらに、得られた各キーワードのロコミ的伝搬パターンに対して、その伝搬経路長と伝搬速度を表示する。

2. ロコミパターン抽出機能

2.1 概要

本システムでは、まずブログ文書ストリームに対して形態素解析を行い単語に分解する。そして、活発に文書化されるトピックを特徴づける単語を抽出し、キーワードとする。次に、キーワードの記述があったブログ文書からキーワードの伝搬に関与した可能性のあるブログネットワークを同定し、その伝搬経路からロコミパターンを抽出する。

2.2 キーワード抽出

キーワードの抽出法は、Gruhl らの手法を用いた。まず、 i 日目に単語 t が出現した回数 $tf(i)$ を算出し、 $tf(i) > T_n$ を満たす t を抽出する。次に、抽出された t の 0 日目から $i-1$ 日目までの出現回数 $\sum_{j=0}^{i-1} tf(j)$ と $tf(i)$ から $tfidf(i) > Th$ を満たす単語をキーワードとし抽出する。

$$tfidf(i) = (i-1)tf(i) / \sum_{j=0}^{i-1} tf(j)$$

2.3 伝搬経路抽出

キーワードが記述されたブログ文書群から、そのブログ文書を記述したユーザーと時間を同定する。次に、ブログネットワークから同じキーワードの伝搬に関与した可

能性のあるユーザーのペア群を抽出する (図 1 参照)。このペア群の中から、まず、どのユーザーからも伝搬を受けてないペアから順に、伝搬先がなくなるまで伝搬するペアを探索する。これにより、すべてのキーワードの伝搬経路を抽出する。

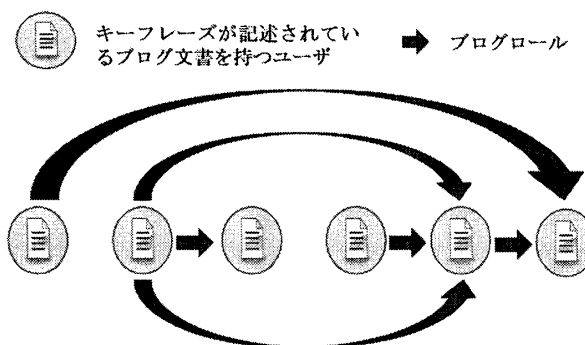


図 1: キーワード伝搬への関与可能性があるユーザーペア

3. 伝搬分析機能

3.1 伝搬経路長

各キーワードの一つの伝搬経路ごとに、一番初めに情報伝搬させたユーザーから、最後に伝搬を受けたユーザーまでに経由したユーザーの総数を算出する。

3.2 伝搬速度

各キーワードの一つの伝搬経路ごとに、一番初めに情報伝搬させたユーザーから、最後に伝搬を受けたユーザーまでの伝搬にかかった時間と、伝搬経路長に基づいて伝搬速度を算出する。

4. 実験

4.1 実験データ

実験では、Doblog^{※1} のデータ^{※2} を、分析対象の文書ストリームデータとして使用した。ブログ総記事数 52, 525、ブログロール総数(ブックマーク総数) 115, 552 であった。

4.2 実験設定

キーワード抽出における形態素解析エンジンには

※1 株式会社 NTT データ

※2 (株) ホットリンクと (株) NTT データの共同事業契約に基づき、(株) ホットリンクより提供。2003 年 10 月から 2005 年 6 月のデータを利用

*1 龍谷大学 理工学部 電子情報学科

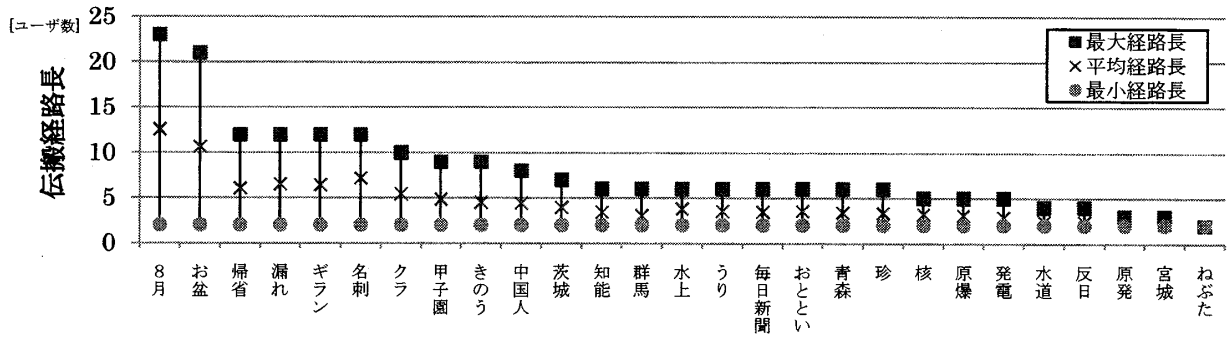


図 2: 伝搬経路長の最大経路長に基づくランキング (8月 10 日のキーワード)

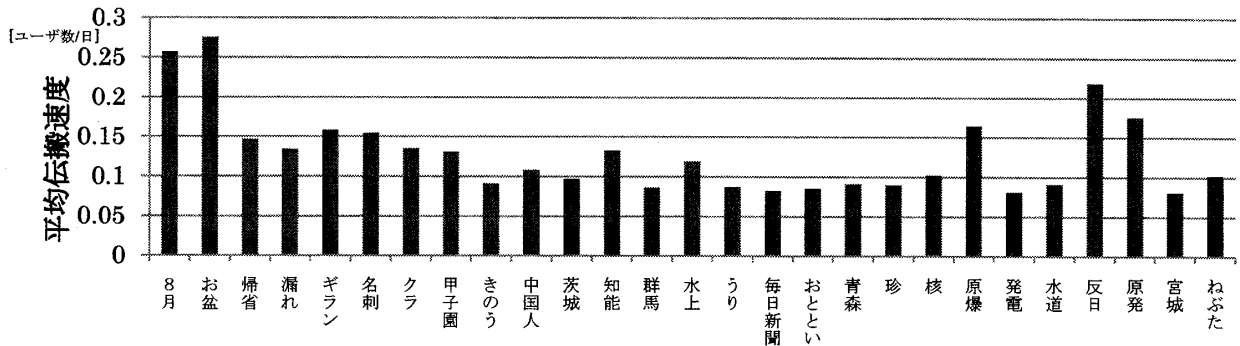


図 3: 平均伝搬速度に基づく分析結果 (8月 10 日のキーワード)

MeCab を使用した。ブログ文書ストリームには、Doblog の 2004 年 6 月 1 日から 2004 年 9 月 10 日までのデータを使用し、ブログネットワークとしては、ブログロールを用いた。キーワード抽出のパラメータは、Gruhl らと同じ $T_n=10$ とし、 $Th=3$ とした。キーワードとなる単語には、形態素解析した結果のうち名詞だけを用いた。

4.3 実験結果

キーワード抽出には、2004 年 6 月 1 日から 2004 年 8 月 9 日までのデータに基づいて、2004 年 8 月 10 日のキーワードを抽出した。27 個のキーワードが抽出された (図 2, 図 3 参照)。本実験では、2004 年 7 月 10 日から 2004 年 9 月 10 日までの期間におけるキーワード伝搬を調べた。図 2 は、伝搬経路の最大経路長に基づいて、抽出されたキーワードをランキングした結果である。図 3 は、各キーワードの平均伝搬速度の出力結果である。キーワード抽出では、「8月」や「お盆」のようなロコミ情報との関連が薄いと考えられるものも抽出されたが、「ギラン」や「クラ」のようなロコミ性の強い情報も抽出できた。ロコミパターン抽出機能を高度化して、ロコミ的情報のキーワード群を抽出することは、今後の重要な課題である。

「ギラン」の情報伝搬パターンでは、総経路数 2624、最大経路長 12、平均経路長 7.3、最小経路長 2、平均伝搬速度 0.15 であった。平均伝搬速度からみると、あるユーザ

から別のユーザへ伝搬するのに、平均約 1 週間かかったことになる。よって、「ギラン」の情報伝搬では、多くの伝搬経路が存在し、伝搬経路長もある程度長く、ゆっくりとした伝搬であったことが示唆される。一方、「反日」の情報伝搬は、伝搬経路長は短い、速い伝搬であったことが見て取れる。

5. まとめ

ブログ空間におけるロコミ的な情報伝搬現象分析システムの構築を目指し、キーワードに着目した分析システムを試作した。大規模ブログデータを用いた実験により、本システムの性能を調べた。キーワードの抽出法にはいくつかの問題点が生じた。

謝辞

Doblog データは (株) NTT データおよび (株) ホットリンクより提供を受けた。記して感謝致します。

参考文献

[1] Gruhl, D., Guha, R., Liben-Nowell, D., and Tomkins, A.: Information diffusion through blogspace, in WWW'04, pp. 107-117 (2004)