

# ブログデータ集合からの頻出なコミュニティ抽出手法

Extraction Method of Frequent Communities from Blog Data Sets

高木 允<sup>♦</sup> 森 康真<sup>♦</sup> 田村 慶一<sup>♦</sup>  
黒木 進<sup>♦</sup> 北上 始<sup>♦</sup>

Makoto TAKAKI Yasuma MORI  
Keiichi TAMURA Susumu KUROKI  
Hajime KITAKAMI

近年、ブログ空間からの知識発見に関する研究が様々行われている。本研究では、プログラーに焦点を当て、プログラーをノード、トラックバックによる繋がりを辺としたグラフから、数ヶ月に渡って頻出するコミュニティを発見する手法を提案する。提案手法では、各グラフの辺へラベル（識別 ID）付けを行い、ラベルをアイテムとみなし、頻出アイテム集合を抽出し、頻出部分グラフを抽出する。さらに、得られた頻出部分グラフを Newman らによって提案されているアルゴリズムを用いてクラスタリングし、コミュニティとなり得るクラスタを抽出する。実際に複数ヶ月に渡りブログデータを収集し、作成したグラフに提案手法を適用した。抽出されたコミュニティ内部を解析した結果、共通の興味・関心を持つ頻出なコミュニティを発見できた。

Recently, many studies on blog spaces have been published. In this paper, we propose a technique to discover frequent blogger communities across the multiple months in bloggers graphs. In our method, a node of bloggers graph represents a blogger and an edge represents a trackback connection. The discovery of frequent communities leads to the identification of groups of bloggers with common interests. The technique, which uses the clustering algorithm proposed by Newman, extracts and clusters frequent subgraphs. After experimentally collecting blog articles and creating bloggers graphs, we applied the proposed technique and discovered the frequent communities.

## 1. はじめに

ウェブログ（ブログ）の登場によりウェブに関する深い知識を持たない人々も容易に情報を発信できるようになっている。ブログは個人の意見を反映したものが多く、世の中の動きを知る上でブログ空間から有益な知識を発見することが重要な課題となっている。

著者らは、ブログ記事ではなく、ブログの書き手であるプログラーに着目し、プログラーをノードとみなし、記事のトラッ

クバックに基づくプログラー同士の繋がりを辺とみなしたグラフ構造に着目している。その中で、ある一定の期間ごとに発生するグラフの集合を時系列グラフと呼び、その時系列グラフから頻出なコミュニティを発見することを目標としている。時系列グラフから抽出される頻出な部分グラフの中に存在するクラスタが特定の話題に偏ったブログ記事を持つとき、そのクラスタを頻出なコミュニティであると定義する。

本稿では、前述の時系列グラフから頻出なコミュニティを発見するために、以下の2つの処理を用いた方法を提案する。

- (1) 時系列グラフから頻出部分グラフ抽出する。
- (2) 頻出部分グラフに存在する全てのコミュニティを見つけ出す。

提案手法では、各グラフの辺にラベル付けを行い、ラベルをアイテムとしたトランザクションデータベースを作成する。トランザクションデータベースから極大頻出アイテム集合（辺の集合）を抽出し、得られた辺の集合からグラフを復元する。復元されたグラフは頻出部分グラフであり、Newmanらの手法[1]（以下、*Newman*）を用いてクラスタリングし、頻出なコミュニティを発見する。

頻出部分グラフをクラスタリングすることで、単一グラフをクラスタリングする場合と比べ、よりコアなコミュニティを正確に発見できる。提案手法では、長期的に形成されているコミュニティを抽出できるため、頻出なコミュニティに向けた効率的なマーケティング、特定の話題に特化したブログ検索、プログラーへの情報推薦などへの応用も期待できる。

抽出されるコミュニティの特徴としては、特定の話題について長期間議論しているコミュニティであり、辺の意味を考慮せずクラスタリングするため、特定の話題で繋がっているだけでなく、コミュニティ内でなんらかの交流関係や社会的な繋がりを持つコミュニティであると考えられる。

論文の構成は以下の通りである。2章で関連研究について述べ、3章で提案手法について説明し、例を用いて頻出部分グラフ抽出とクラスタリングについて説明する。4章でデータ収集の説明と実際のデータに提案手法を適用した結果を示し、5章でまとめる。

## 2. 関連研究

### 2.1 コミュニティに関する研究

従来のコミュニティに関する研究[2,3,4]では、Web上の記事やブログ記事をノードとするコミュニティに着目しているのに対して、本研究では、人（記事の投稿者）をノードとし、時間経過とともに頻りに現れる人のコミュニティを見つけ出すことに着目している点が大きく異なる。即ち、本研究では、消滅しやすいコミュニティよりも長期的に安定した繋がりを持つ、人のコミュニティの抽出に着目している。

### 2.2 グラフマイニングに関する研究

与えられたグラフデータベース  $D=\{G_1, \dots, G_n\}$  から頻出な部分グラフを抽出する研究が様々行われている[5,6,7,8]。文献[5,6,7]では、同一ラベルを持つノードや辺が複数存在する一般グラフから同型な頻出部分グラフを抽出する方法が提案されている。一般グラフから、同型な部分グラフを生成するには多大な計算時間を要するため、これらの文献では、各  $G_i$  のノード数が 15~40 程度の疎なグラフを扱っている。本研究で扱うグラフは、ノードと辺のどちらにも同一ラベルを許さない時系列グラフであるので、一般グラフとは異なる。また、各  $G_i$  のノード数が数千ノード規模であり、12ヶ月間に渡り収集した時系列グラフを想定している。以上により、

<sup>♦</sup> 学生会員 広島市立大学大学院情報科学研究科博士後期課程 / 日本学術振興会 DC

[makoto@db.its.hiroshima-cu.ac.jp](mailto:makoto@db.its.hiroshima-cu.ac.jp)

<sup>♦</sup> 正会員 広島市立大学大学院情報科学研究科

[{mori,ktamura,kuroki,kitakami}@its.hiroshima-cu.ac.jp](mailto:{mori,ktamura,kuroki,kitakami}@its.hiroshima-cu.ac.jp)

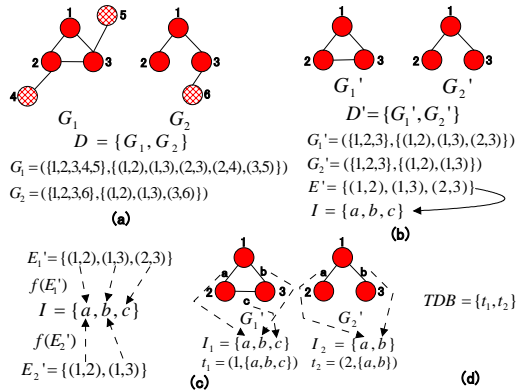


図1 辺へのラベル付けとTDB作成例

Fig.1 An example of labeling and creation of TDB

文献[5,6,7]の手法を我々が想定しているコミュニティ発見の問題に応用するには、計算時間の面で不向きである。

文献[8]で提案されているCODENSEは、分子生物学の分野で開発されたアルゴリズムである。これは、カットを用いて密な頻出部分グラフを高速に抽出する方法として知られているが、著者らがブログ空間に応用した実験では、本来1つになるべきコミュニティがサイズの小さな頻出部分グラフに分割されてしまうという問題が生じている。また、スター構造のようなコミュニティは、密な頻出部分グラフとして抽出できないという問題も生じている。

以上の問題点に対して、本研究では、2つの処理を用いて、コミュニティを高速に発見する方法について提案している。

### 3. 提案手法

#### 3.1 アルゴリズム

ブロッカーをノード、トラックバックによる繋がりを辺とした重みなし・無向単純グラフの集合であるグラフデータベースを  $D = \{G_1, \dots, G_n\}$  とする。以後、 $u, v$  はひとつのノード、 $(u, v)$  は辺を表す記号とする。図1(a)に  $D$  の例を示す。 $D$  中のグラフ  $G_i$  は  $G_i = G(V_i, E_i)$ ,  $E_i \subset V_i \times V_i$ ,  $(u, v) \in E_i$ ,  $1 \leq i \leq n$  と定義される。 $V_i$  はノードの集合、 $E_i$  は辺の集合である。提案手法のアルゴリズムは以下の3つのステップから成る。

- (1) 全てのグラフに共通して存在するノードを抽出
- (2) グラフの辺にラベルを付与し、頻出部分グラフを抽出
- (3) Newmanを用いて頻出コミュニティを抽出

以下に、各ステップにおけるアルゴリズムの詳細を示す。

##### (1) 全てのグラフに共通して存在するノードを抽出

頻出部分グラフを抽出するために  $D = \{G_1, \dots, G_n\}$  から全ての  $G_i$  について共通しているノードを取り出したグラフデータベース  $D' = \{G_1', \dots, G_n'\}$  を作成する。図1(b)に  $D'$  の例を示す。 $G_i' = G(V', E_i')$  と定義すると、ノード集合  $V'$  は全ての  $G_i'$  に対して同一である。つまり、 $V' = V_1 \cap V_2 \cap \dots \cap V_n$  である。また  $E_i' = \{(u', v') \mid u', v' \in V', (u, v) \in E_i\}$  と定義できる。 $E_i'$  は  $V'$  に含まれるノードのペアのみで構成されている。

##### (2) グラフの辺にラベルを付与し、頻出部分グラフを抽出

$E'$  を全ての  $E_i'$  の和集合とする。 $|E'|$  個のラベルを要素としたラベル集合  $I$  を作成する (図1(b))。  $E'$  から  $I$  への全単射  $f$  と定義すると、 $f$  は以下のように表現できる。

$$f: E' \rightarrow I \text{ または } I = f(E')$$

$f$  を用いて  $E_i'$  からラベル集合  $I_i = \{\text{label}_{i1}, \dots, \text{label}_{i|E_i'|}\}$  を作成する (図1(c))。このラベル集合をアイテム集合とみなし、トランザクションデータベース  $TDB = \{t_1, \dots, t_n\}$  を作成する。ここで、 $t_i = (i, I_i)$  である (図1(d))。

TDB から極大頻出アイテム集合を高速に抽出するために、文献[9]で提案されている手法を用いた。最小支持数を  $\text{min\_sup}$  とし、極大頻出アイテム集合を抽出する関数を  $\text{EX\_MAX}(\text{min\_sup}, TDB)$  とする。 $\text{EX\_MAX}$  で得られた極大頻出アイテム集合を  $\text{MAX\_PAT} = \{\text{PAT}_1, \dots, \text{PAT}_l\}$  とする。ここで、 $\text{PAT}_i = \{\text{label}_{i1}, \dots, \text{label}_{im}\}$ ,  $1 \leq i \leq l$ ,  $1 \leq m$  である。 $\text{PAT}_i$  から、辺の集合  $FE_i$  へ、 $f$  の逆写像  $f^{-1}$  を用いて変換する。

$f^{-1}$  を用いてアイテム集合から辺集合へ変換し、辺集合からノード集合を得る。得られた辺集合とノード集合から頻出部分グラフ  $\text{FSG}$  を復元する。頻出部分グラフ集合を、 $\text{FR\_SUBG} = \{\text{FSG}_1, \dots, \text{FSG}_l\}$  と定義する。 $\text{FSG}_i = G(\text{FV}_i, \text{FE}_i)$ ,  $\text{FE}_i = \{(u_i, v_i) \mid u_i, v_i \in V'\}$ ,  $1 \leq i \leq l$  であり、 $\text{FV}_i$  は  $\text{FE}_i$  を構成する全てのノードの集合である。

##### (3) Newmanを用いて頻出コミュニティを識別

Newmanを  $\text{FSG}_i$  に適用する。Newmanを適用して得られる出力  $\text{RESULT}$  が最終的なクラスタリング結果である。

本手法の特長は、ブロッカーのグラフに直接 Newman を適用するのではなく、頻出部分グラフに Newman を適用することで、一過性のノイズを除去したコミュニティを発見することである。一過性のノイズの中にも有益な情報が含まれているが、本研究では、頻出なコミュニティを発見するという立場をとっているため、今回は一過性のノイズを除去した頻出なコミュニティの発見を行っている。

Newman を重みなしグラフに適用すると、グラフサイズに対して多数の非常に小さなクラスタと少数の非常に大きなクラスタに分割されるという問題点が文献[10]で示唆されている。文献[10]では、ブログの重みなしグラフに Newman を適用した場合、クラスタ内のブログの内容が多様であり、共通する話題を見つけることが困難であると報告している。しかしながら、辺に重み付けを行った場合、非常に小さなクラスタと非常に大きなクラスタの数が減少し、クラスタのサイズが平均化され、辺への重み付けを行うことにより Newman を用いても的確なクラスタリングが可能となることが示唆されている。

本研究においては、重みなし・無向単純グラフに Newman を適用しているが、頻出な辺のみを取り出したグラフへ Newman を適用する。辺に、頻出という概念で重み付けを行っているため、正確なクラスタリングができる。

#### 3.2 頻出部分グラフ抽出とクラスタリングの例

本節では、例を用いて図2に示すグラフデータベース  $D' = \{G_1', G_2', G_3', G_4'\}$  からの頻出部分グラフ抽出とクラスタリングの例を示す。

まず、図2(a)に示す  $D'$  から頻出部分グラフ  $\text{FR\_SUBG}$  を抽出する。全ての  $E_i'$  の和集合  $E'$  を求める。 $|E'|$  個のラベルを持ったラベル集合  $I = \{a, b, c, d, e, f, g, h, i, j, k, l\}$  を得る。 $f$  を用いて、 $E'$  をラベル付けする。さらに、 $f$  を用いて  $E_i'$  に対応するラベルをアイテムとしたトランザクション  $I_i$  を生成する。そして、 $TDB$  を作成する (図2(b))。このとき、 $TDB = \{t_1, t_2, t_3, t_4\}$  であり、 $t_1 = (1, I_1)$ ,  $t_2 = (2, I_2)$ ,  $t_3 = (3, I_3)$ ,  $t_4 = (4, I_4)$  である。ただし、 $I_1 = \{a, b, c, d, e, f, g\}$ ,  $I_2 = \{a, b, f, h, i\}$ ,  $I_3 = \{a, b, c, d, e, f, g\}$ ,  $I_4 = \{b, e, j, k, l\}$  である。

そして、図2(c)に示すように、 $TDB$  から  $\text{EX\_MAX}$  を用いて極大頻出アイテム集合を抽出する。最小支持数は2としている。 $\text{EX\_MAX}(2, TDB)$  によって得られた極大頻出アイテム集合は  $\text{MAX\_PAT} = \{\text{PAT}_1\}$ ,  $\text{PAT}_1 = \{a, b, c, d, e, f, g\}$  である。抽出された  $\text{MAX\_PAT}$  を辺集合へ変換し (図2(d))、辺から元のグラフ  $\text{FR\_SUBG} = \{\text{FSG}_1\}$  を復元する。ここで、 $\text{FSG}_1 = G(\text{FV}_1,$

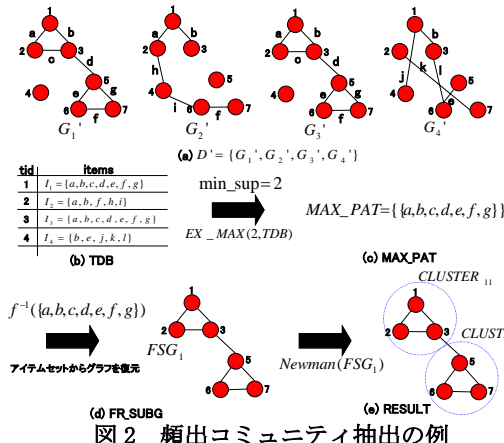


図2 頻出コミュニティ抽出の例

Fig. 2 An example of frequent community extraction  $FE_1$ ,  $FV_1 = \{1, 2, 3, 5, 6, 7\}$ ,  $FE_1 = \{(1, 2), (1, 3), (2, 3), (3, 5), (5, 6), (5, 7), (6, 7)\}$ である。最後に  $FSG_1 \sim Newman$  を適用する。  
 Newmanによって得られる結果は,  $RESULT = \{\{CLUSTER_{11}, CLUSTER_{12}\}\}$  となり,  $CLUSTER_{11} = G(CV_{11}, CE_{11})$ ,  $CV_{11} = \{1, 2, 3\}$ ,  $CE_{11} = \{(1, 2), (1, 3), (2, 3)\}$ ,  $CLUSTER_{12} = G(CV_{12}, CE_{12})$ ,  $CV_{12} = \{5, 6, 7\}$ ,  $CE_{12} = \{(5, 6), (5, 7), (6, 7)\}$  である (図 2 (e)).

4. 評価

本章では, 収集データと評価の説明を行う。2. で示した本研究で想定するグラフデータの規模は  $|D|=12$  であったが, 今回の実験では収集できた4つのグラフを基に実験を行った。

4.1 データ収集とグラフ作成

収集を行ったブログ記事は2006年6月1日から2006年6月30日までの記事のように1ヶ月単位で, 2006年6月から2006年9月までの4か月分を収集した。つまり,  $D = \{G_1, G_2, G_3, G_4\}$  となる。各  $G_i$  の詳細を表1に示す。

データ収集開始時において, 始点となる記事(ブロガー)をランダムに選択する。記事からトラックバックを抽出し, トラックバックを辿ることにより記事を収集する。この動作をトラックバックがなくなるまで行う。データ収集が終了した時点で出来上がるグラフのノード数は収集したブロガー数と等しい。各月において他の月に収集したグラフのデータは全く参照しない。つまり, 6月においてあるノード  $u, v$  間に形成されていた辺があったとしても, 7月のデータを収集する際には全てのノードと辺の情報は削除されており, 新規にグラフを作成していく。

4.2 評価実験

実際に記事を収集して作成した4か月分のグラフデータベース  $D$  から, 4ヶ月に渡って共通して出現しているブロガーを抽出し, グラフデータベース  $D' = \{G_1', G_2', G_3', G_4'\}$  を作成する。全ての月に存在していたブロガーは319人であった。つまり,  $|V|=319$  であり,  $|E_1|=1,650$ ,  $|E_2|=1,695$ ,  $|E_3|=1,697$ ,  $|E_4|=1,230$  であった。まず, グラフデータベース  $D$  中の  $G_i$  に直接  $Newman$  を適用した結果について示す。続いて, 提案手法を適用した結果を示す。

4.2.1  $G_i$  のクラスタリング結果

$D$  中の6月のブロガーのグラフ  $G_1$  に  $Newman$  を適用した結果, 43個のクラスタが識別された。クラスタサイズは最大で1195, 最小で2であり, サイズが10未満のクラスタが29個, サイズが400以上のクラスタが6個, 残りのクラスタはサイズが16から156であった。極端に小さなクラスタが多数存在し, 極端に大きなクラスタと中間サイズのクラ

表1 収集したデータの詳細  
 Table 1 Details of collected data

グラフ $G_i$	$ V_i $ (ノード数)	$ E_i $ (辺の数)	記事数
$G_1$ (2006年6月)	4,432	28,733	15,569
$G_2$ (2006年7月)	3,147	18,986	9,878
$G_3$ (2006年8月)	3,951	23,966	9,715
$G_4$ (2006年9月)	2,288	10,760	3,406

タは少数であった。

各クラスタについて  $tf-idf$  を用いた解析を行った。上位にランキングされたキーワード同士は関連がなく, 抽出されたクラスタがコミュニティであると判断できない。実際に複数のクラスタについてブロガーの記事を調査したところ, ある野球チームの話題を主としているブロガー, 政治の話題を主としているブロガーのように, 様々なブロガーが混在していた。これは, ある月にだけトラックバックを張ったブロガーが多く存在しており, 抽出されるべきコミュニティにうまく分割できなかったために起こった現象である。

例えば, 2006年6月にはサッカーワールドカップが開催されていたため, 野球や政治の話題を主に扱っていたとしてもサッカーワールドカップにも興味があれば, サッカーワールドカップの記事を書き, サッカーを主に話題にしているブロガーにトラックバックを張るといった現象が起こる。しかしながら, このトラックバックは一過性のものであり, 毎月トラックバックを張り続ける関係ではない。

このように, あるイベントの発生時のみの一過性のトラックバックが多く混在すると, 共通の興味・趣味を持ったコミュニティの発見が困難となる。

単一のグラフではなく, 全ての  $G_i$  の辺の和を取り, 指定した頻度でフィルタリングしたグラフを作成し,  $Newman$  を用いてクラスタリングするという手法も考えられるが, クラスタ間の境界面が不鮮明になり, 複数のコミュニティが1つのコミュニティとして抽出される結果となってしまう, 精度の高いコミュニティ抽出が不可能となる結果となった。

4.2.2 提案手法の結果

頻出部分グラフを抽出するための最小支持数は  $|D|$  の半分である2と設定した。頻度を  $|D|$  の半分以上とすることで, 一時的な興味などによりトラックバックを張っているような繋がりを除去することができる。抽出された頻出部分グラフは全部で6個あった。ここでは, 抽出された頻出部分グラフである  $FSG_1$  について説明する。 $FSG_1$  をクラスタリングした結果, 13個のクラスタが識別された。

図3に,  $FSG_1$  のクラスタリングを行い,  $CLUSTER_{14}$  を取り除いて可視化したものを示す。 $CLUSTER_{14}$  中のブロガーの記事を調査したところ, 主にスポーツ全般の話題を中心としており, 各クラスタからのトラックバックが多数張られていた。 $CLUSTER_{14}$  においては, クラスタ間の境界面に, 特定の話題に興味を持つコミュニティに属すべきノードが多数含まれていた。これは, クラスタリングに使用した  $Newman$  がクラスタ間の境界面をうまく識別できなかったことが理由として挙げられる。しかしながら, 本実験で収集したデータからは,  $CLUSTER_{14}$  のようなハブクラスタが存在し, その周辺にコアなコミュニティが存在しているというクラスタ間の構造を確認できた。頻出部分グラフを抽出することで一過性の繋がりを除去し, 長期的に存在しているコミュニティの発見ができた。

次に,  $FSG_1$  内のクラスタリングされたノード集合がどのような話題を中心としているのかを調べるために  $tf-idf$  を用い



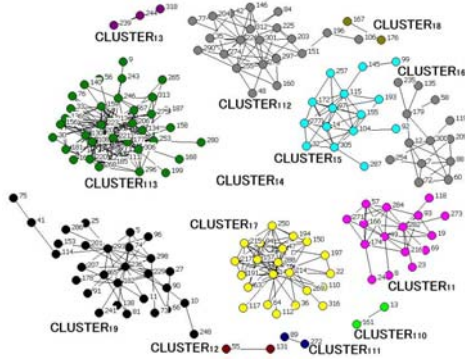


図3 識別されたコミュニティ  
Fig.3 Extracted communities

て重要語句を抜き出した。表2に  $FSG_1$  をクラスタリングした結果に  $tf-idf$  を適用した結果を示す。  $tf-idf$  の値が高い上位3件を示している。図3中の  $CLUSTER_j$  と表2中の  $C_{ij}$  がそれぞれ対応している。

表2から、 $tf-idf$  によって抽出されたキーワード上位3件はそれぞれ容易に連想できるキーワードとなっている(例えば、表2の  $CLUSTER_{11}$  はプロ野球球団、広島東洋カーブについての記事を扱っている集団である)。さらに、手作業でクラスターに属しているブロガーのブログを確認したところ、 $CLUSTER_{113}$  では、阪神タイガースファンのブロガーが90%を占めていた。ファンの判断基準としては、ブログの題名やプロフィールなどに、自ら阪神ファンであることを記述しているブロガーをファンとした。このように、クラスタリング結果とクラスターを解析した結果が強い相関を持っているのは、頻出部分グラフを抽出することで一過性のトラックバックによるブロガー間の繋がりを除去することができ、より繋がりの強いブロガー同士の繋がりのみを抽出し、クラスタリングできたためである。

このように、グラフデータベース  $D$  から  $D'$  を作成し、頻出部分グラフを抽出してクラスタリングを行うことで、ノイズを除去したよりコアなコミュニティを抽出できた。また、頻度を  $|D|$  の半分である2としても精度の高いコミュニティを抽出することができた。

## 5. おわりに

本論文では、時系列グラフから頻出コミュニティを発見するために、辺へのラベル付けを行って頻出部分グラフを抽出し、個々の頻出部分グラフをクラスタリングする手法を提案した。実際に収集した時系列グラフから抽出された頻出なコミュニティを解析した結果、数ヶ月に渡って同じ興味・関心を持つコミュニティ、特に、野球・政治に関しての話題を深く議論しているコミュニティを発見できた。単一グラフをクラスタリングした結果と比較すると、ノイズを含まない正確なコミュニティを発見できた。

今後は、トラックバックだけではなくコメントやブログ本文中のブログ記事へのリンクも考慮したグラフからのコミュニティ抽出を行う予定である。また、4.2.2で述べたクラスター間の境界面における問題を解決するために、現在、重複を許したクラスタリング手法の検討を行っている。

## 【謝辞】

本研究の一部は、日本学術振興会・特別研究員奨励費(課題番号:18・0205)、日本学術振興会・科学研究費補助金(基

表2  $tf-idf$  値上位3件のキーワード  
Table 2 Top 3 key words of  $tf-idf$

クラスター	$tf-idf$ 値上位3件		
	1	2	3
$C_{11}$	カーブ	広島	日本
$C_{12}$	投手	楽天	野球
$C_{13}$	楽天	イーグルス	野球
$C_{14}$	日本	ブラジル	ドイツ
...	...	...	...
$C_{111}$	日本	国家	国旗
$C_{112}$	ホークス	鷹	投手
$C_{113}$	阪神	虎	阪神タイガース

盤研究(C)(一般)、課題番号:17500097)の支援により行われた。

## 【文献】

- [1] A. Clauset, M. E. J. Newman, and C. Moore. Finding Community Structure in Very Large Networks. *Physical Review E*, 70:066111, 2004.
- [2] G. W. Flake, S. Lawrence, C. L. Giles, and F. Coetzee. Self-Organization of the Web and Identification of Communities. *IEEE Computer*, 35(3):66-71, 2002.
- [3] 豊田 正史, 喜連川 優. 日本におけるウェブコミュニティの発展過程. *日本データベース学会 Letters Vol.2, No.1*, pp.35-38, 2003.
- [4] D. Gruhl, R. V. Guha, D. Liben-Nowell, and A. Tomkins. Information Diffusion Through Blogspace. In *WWW*, pages 491-501, 2004.
- [5] A. Inokuchi, T. Washio, and H. Motoda. An Apriori-Based Algorithm for Mining Frequent Substructures from Graph Data. In *PKDD*, pages 13-23, 2000.
- [6] M. Kuramochi and G. Karypis. Frequent Subgraph Discovery. In *ICDM*, pages 313-320, 2001.
- [7] X. Yan and J. Han. gSpan: Graph-Based Substructure Pattern Mining. In *ICDM*, pages 721-724, 2002.
- [8] H. Hu, X. Yan, Y. Huang, J. Han, and X. J. Zhou. Mining Coherent Dense Subgraphs Across Massive Biological Networks for Functional Discovery. In *ISMB*, pages 213-221, 2005.
- [9] T. Uno, M. Kiyomi, and H. Arimura. LCM ver. 2: Efficient Mining Algorithms for Frequent / Closed / Maximal Itemsets. In *FIMI*, 2004.
- [10] 安藤潤, 吉井伸一郎. WWW ナビゲーション向けコミュニティ分割手法に関する一考察. *情報処理学会研究報告 知能と複雑系*, pp.115-122, January 2006.

## 高木 允 Makoto TAKAKI

広島市立大学院情報科学研究科博士後期課程在学中。2006年より日本学術振興会 DC。

## 森 康真 Yasuma MORI

広島市立大学大学院情報科学研究科助教。1994年北陸先端科学技術大学院大学情報科学研究科博士前期課程修了。

## 田村 慶一 Keiichi TAMURA

広島市立大学大学院情報科学研究科助教。2000年九州大学大学院システム情報科学研究科修士課程修了。博士(情報科学)。

## 黒木 進 Susumu KUROKI

広島市立大学大学院情報科学研究科准教授。1990年東京大学工学系研究科計数工学博士前期課程修了。博士(工学)。

## 北上 始 Hajime KITAKAMI

広島市立大学大学院情報科学研究科教授。1976年東北大学大学院工学研究科博士前期課程修了。博士(工学)。