

Refereed Conference paper

## ネットワークの局所構造における成長パターンの 発見とその応用

甲谷 優<sup>†1</sup> 西田 京介<sup>†1</sup> 藤村 考<sup>†1</sup>

ネットワークの局所構造における成長パターンに着目したネットワークを分析する手法を提案する。本手法によりこれまでの分析では理解することが困難であったネットワークの生成過程に関する性質が明らかになり、さらにどのノード間に将来エッジが発生するかを予測することができる。エッジの発生の予測は「リンク予測」問題と呼ばれ、近年のソーシャルネットワーキングサービス等で友人推薦のコア技術となっており注目されている。本稿ではソーシャルメディアを用いて実験を行い、提案法により明らかになるソーシャルネットワークの性質と、リンク予測アルゴリズムとしての性能について報告する。

## Discovering Evolutionary Patterns of Substructures in Complex Networks and its Applications

YUTAKA KABUTOYA,<sup>†1</sup> KYOSUKE NISHIDA <sup>†1</sup>  
and KO FUJIMURA<sup>†1</sup>

We propose an entirely new approach to understanding complex networks; called dynamic network motifs (DNMs), it is significantly different from existing network analysis approaches in terms of focusing on the evolution of substructures in a graph. The concept has the potential to uncover previously unknown properties of graphs and allows us to infer whether new edges among nodes the present in a snapshot of the graph are likely to occur. This inference, called “link prediction”, is significant because it yields various useful applications, such as interaction recommenders for social networking services. Experiments on snapshots of social networks based on real log data sets of social media show the following two significant results: First, DNM analysis has the potentials described above. Next, a new link predictor based on DNMs outperforms existing link predictors in terms of prediction accuracy.

### 1. はじめに

様々な分野でネットワーク分析に関する研究は数多くなされてきたが、近年は計算機の処理能力の向上から、部分ネットワーク (ネットワークのノードの部分集合と、それらの間のエッジからなるネットワーク) に関するマイニングが可能となった<sup>10),13)</sup>。

部分ネットワーク分析の1つに、ネットワークモチーフ分析と呼ばれる手法が存在する。ネットワークモチーフとは、ネットワーク中に特徴的に現れる部分ネットワーク構造と定義されており、ネットワーク理解のための新たな分析手法として注目を集めている。Milo<sup>ら7),8)</sup> はネットワークが時間とともに成長してもそのネットワークモチーフは変わらないことを発見し、ネットワークモチーフはネットワークの成長に関する制約を表すと言及している。しかしながらネットワークモチーフはある時点での静的なネットワークのスナップショットに対する分析であり、ネットワークの時間にもなう成長には着目できていない。

一方で、近年は部分ネットワークの成長パターンに関する研究も始まっている。たとえばInokuchi<sup>ら4)</sup> は効率のよい部分ネットワークの成長パターンの発見に関して研究を行っているが、未だに部分ネットワークの成長パターンに関する統計的分析に関する研究はなされていない。そのため、どの部分ネットワークの成長パターンが重要なのかについて議論することは困難である。

我々は文献<sup>14)</sup>にて部分ネットワーク中のエッジの発生、再発生に関する統計的分析手法を提案した。本稿では、前述の近年の研究の流れを踏まえエッジの再発生は省き成長 (新しいエッジの発生) に関するもののみに着目し、ネットワークの成長過程における統計的に重要な部分ネットワークの成長パターン (動的ネットワークモチーフと定義する) の発見手法を提案する。3ノードの部分ネットワークの場合、エッジの発生と再発生のパターンは54通りであるが、そのうち本稿で着目する成長パターンは24通りである。

また我々<sup>14)</sup>は動的ネットワークモチーフのリンク予測への応用も提案している。リンク予測とはある時点でのネットワークのスナップショットが与えられた際に、そのネットワーク中のどの2ノード間に将来エッジが発生するかを推定する問題である。リンク予測技術はSNSにおける友人推薦機能<sup>\*1</sup>をはじめとしてさまざまなアプリケーションに利用されてい

<sup>†1</sup> 日本電信電話株式会社 NTT サイバソリューション研究所  
NTT Cyber Solutions Laboratories, NTT Corporation

\*1 たとえば、mixi (<http://mixi.jp>) が友人推薦機能を提供している

る。従来のリンク予測アルゴリズムの多くはネットワークの大域的特性に基づいている\*1が、我々の手法はその局所構造に基づいている。本稿では文献<sup>14)</sup>記載のヒューリスティックな手法とは異なり、局所構造の変化パターンのうち成長に関するもののみを素性として2ノード間のエッジの発生を確率モデル化する手法を提案する。理論的な基礎を与えることで、より高精度なリンク予測が可能となると考えたためである。

本稿では様々なソーシャルメディアからネットワークを抽出し、動的ネットワークモチーフ分析と従来のネットワークモチーフ分析の比較を行う。また、リンク予測実験を行い、提案法が従来法よりも高精度にリンク予測できることを示す。

## 2. 従来法：ネットワークモチーフ

提案手法の説明に入る前に、その基礎であるネットワークモチーフ分析について言及する。今、あるネットワークのスナップショットを  $\mathbf{G} = \langle \mathbf{V}, \mathbf{E} \rangle$  とする。ここでネットワーク  $\mathbf{G}$  は有向であってもなくてもかまわない。

まず、サイズ  $n$  パターンを、ノード数が  $n$  で、連結であり、互いに同型でない全ネットワークと定義する。このとき、有向なネットワークの場合サイズ3パターンは13、サイズ4パターンは199得られることが Milo ら<sup>3)</sup>によって発見されている。

ここで与えられたネットワーク  $\mathbf{G}$  においてあるパターンの出現頻度が等価なランダムネットワーク  $\{\hat{\mathbf{G}}_k\}_{k=1}^K$  と比較して統計的に有意に高い場合、そのパターンをネットワークモチーフと定義する。ここで  $\mathbf{G}$  と等価なランダムネットワークは、各ノードの接続するエッジ数、双方向エッジ数が同じになるように発生させる。

Milo ら<sup>7)</sup>の提案している各パターンの統計量は以下のように算出することができる。まず、 $\mathbf{G}$  中のパターン  $i$  の出現頻度を  $m_i$  とする。また、ランダムネットワーク  $\hat{\mathbf{G}}_k$  中のパターン  $i$  の出現頻度を  $\hat{m}_{ik}$  とする。このとき  $\mathbf{G}$  におけるパターン  $i$  の統計量は以下に示すように  $Z$  スコアで与えられる。

$$z_i = \frac{m_i - \mu_i}{\sigma_i}, \quad (1)$$

ただし  $\mu_i = \frac{1}{K} \sum_{k=1}^K \hat{m}_{ik}$ ,  $\sigma_i^2 = \frac{1}{K} \sum_{k=1}^K (\hat{m}_{ik} - \mu_i)^2$ .

このときベクトル  $\mathbf{z} = \{z_i\}_{i=1}^I$  の長さを1とするように正規化することで、ネットワークモチーフプロファイルを得る。

\*1 たとえばクラスタ性に基づいた手法<sup>1),9)</sup>、スケールフリー性に基づいた手法<sup>2)</sup>が存在する

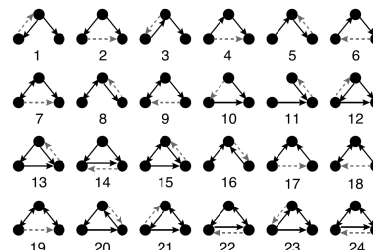


図1 全24通りのサイズ3成長パターン

Fig. 1 All 24 size-three evolutionary patterns.

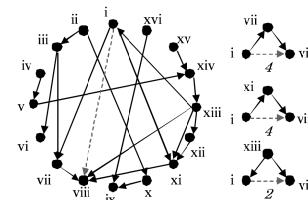


図2 ノードペア  $\langle i, viii \rangle$  は  $\langle i, vii, viii \rangle$  と  $\langle i, xi, viii \rangle$ ,  $\langle i, xiii, viii \rangle$  の3つの成長に対応する

Fig. 2 Pair  $\langle i, viii \rangle$  yields three triad-evolution,  $\langle i, vii, viii \rangle$ ,  $\langle i, xi, viii \rangle$ , and  $\langle i, xiii, viii \rangle$ .

## 3. 提案法：動的ネットワークモチーフ

### 3.1 動的ネットワークモチーフ分析

本研究はネットワークの成長パターンに着目した手法であり、ネットワークが動的であれば期間  $t$  から  $t+1$  内にノード  $\mathbf{V}$  間に発生したエッジが対象となる。本稿では静的なネットワークにも対応するために、それぞれのエッジ  $e \in \mathbf{E}$  がネットワーク  $\mathbf{G}_{\setminus e} = \langle \mathbf{V}, \mathbf{E} - \{e\} \rangle$  に対する新たなエッジの発生とみなす。

まずは  $n$  ノードの連結な部分ネットワーク中における新たな1本のエッジの発生をサイズ  $n$  成長パターンと定義する。既存のエッジを実線で、新たなエッジの発生を点線で表わすと、有向ネットワークにおける全24通りのサイズ3成長パターンは図1のようになる。今扱う成長パターンが  $J$  通りであるとする。

一般に、1本のエッジの発生は複数の部分ネットワークの成長とみなすことができる。例を図2に示す。ここでエッジの発生  $\langle u, v \rangle$  に対応する成長パターン  $j$  の数を  $\langle u, v \rangle$  の  $j$  に関する  $EF$  (Evolution Frequency) と定義し、 $x_{uvj}$  で表わす。図2の場合、 $x_{i, viii, 4} = 2$  となる。サイズ  $n$  成長パターンに着目した場合、 $x_{uvj}$  は  $u, v$  と、 $\mathbf{V}$  中の  $u, v$  以外の  $(n-2)$  ノードの全組合せに対して、それらからなる部分ネットワークを検証することにより得ることができる。このとき  $\mathbf{G}$  内の成長パターン  $j$  の頻度は以下のようにして得られる。

$$n_j = \sum_{\langle u, v \rangle \in \mathbf{E}} x_{uvj}. \quad (2)$$

ここで与えられたネットワーク  $\mathbf{G}$  において等価なランダムネットワーク  $\{\hat{\mathbf{G}}_k\}_{k=1}^K$  と比

較してある成長パターンの頻度が統計的に有意に高い場合、その成長パターンを動的ネットワークモチーフと定義する。等価なランダムネットワークはネットワークモチーフ分析と同様の手法で発生させる。\$\mathbf{G}\_k\$ における成長パターン \$j\$ の頻度を \$\hat{n}\_{jk}\$ とするとき、\$\mathbf{G}\$ における成長パターン \$j\$ の統計的な重みは以下に示すように \$Z\$ スコアで与えられる。

$$\zeta_j = \frac{n_j - \nu_j}{\tau_j}, \quad (3)$$

ただし \$\nu\_j = \frac{1}{K} \sum\_{k=1}^K \hat{n}\_{jk}\$, \$\tau\_j^2 = \frac{1}{K} \sum\_{k=1}^K (\hat{n}\_{jk} - \nu\_j)^2\$.

ベクトル \$\boldsymbol{\zeta} = \{\zeta\_j\}\_{j=1}^J\$ の長さを 1 とするように正規化することで、動的ネットワークモチーフプロフィールを得る。

### 3.2 成長パターンに基づくリンク予測

現在のネットワークが \$\mathbf{G} = \langle \mathbf{V}, \mathbf{E} \rangle\$ であるとき、将来のネットワークの \$\mathbf{V}\$ に関する部分ネットワークを \$\mathbf{G}' = \langle \mathbf{V}, \mathbf{E}' \rangle\$ とする。このとき \$\mathbf{V}\$ の 2 ノードの全順序から \$\mathbf{E}' - \mathbf{E}\$ を予測する問題がリンク予測である。

ここでエッジの発生は対応する EF 値に依存すると仮定する。

$$P(y = 1|u, v) \propto P(y = 1|\mathbf{x}) = P(y = 1|x_1, \dots, x_J), \quad (4)$$

ただし \$\langle u, v \rangle \in \mathbf{E}'\$ の場合 \$y = 1\$、そうでない場合 \$y = 0\$。ここで確率 \$P(y = 1|\mathbf{x})\$ は以下のように推定することができる。

$$P(y = 1|\mathbf{x}) = \frac{|\{\langle u, v \rangle \in \mathbf{E} | x_{uv1} = x_1 \wedge \dots \wedge x_{uvJ} = x_J\}| + \delta}{|\{\langle u, v \rangle \in \mathbf{V} \times \mathbf{V} | x_{uv1} = x_1 \wedge \dots \wedge x_{uvJ} = x_J\}| + 2\delta}, \quad (5)$$

ただし \$\delta\$ はスムージングパラメータである。しかしこのようにモデル化すると、EF 値の最大値が \$W\$ である場合パラメータ数が \$O(W^J)\$ に達し、頑健なパラメータ推定が困難であるという問題点がある。

そこでまずはベイズの定理から \$u\$ から \$v\$ にエッジが発生する確率を、以下に示す尤度で近似する。

$$P(y = 1|\mathbf{x}) = \frac{P(y = 1)P(\mathbf{x}|y = 1)}{\sum_{y'=0}^1 P(y')P(\mathbf{x}|y')} \propto \frac{P(\mathbf{x}|y = 1)}{P(\mathbf{x}|y = 0)}, \quad (6)$$

このとき以下に示す 3 つの方法でパラメータ数を減らし、頑健な推定を行う。

#### (1) 独立性の仮定。

エッジが発生したという条件のもと成長パターン \$j\$ の EF 値と \$j'\$ の EF 値が独立であると仮定すると、下式

$$P(\mathbf{x}|y = 1) = P(x_1, \dots, x_J|y = 1) = \prod_{j=1}^J P(x_j|y = 1), \quad (7)$$

のように分解でき、パラメータ数を \$O(JW)\$ まで減らすことができる。エッジが発生したという条件に対する成長パターン \$j\$ の EF 値の事後確率は以下のように推定できる。

$$P(x_j|y = 1) = \frac{|\{\langle u, v \rangle \in \mathbf{E} | x_{uvj} = x_j\}| + \delta}{|\mathbf{E}| + \delta W}. \quad (8)$$

\$P(\mathbf{x}|y = 0)\$ も同様に頑健に推定することが可能である。

#### (2) ベルヌーイ試行の仮定。

ベルヌーイ試行の仮定は文書分類においてもよく用いられている<sup>3)</sup>。各ノードペアは各成長パターンが出現するか否かのベルヌーイ試行であると仮定すると、

$$P(\mathbf{x}|y = 1) = \prod_{j=1}^J P(b_j = 1|y = 1)^{b_j} (1 - P(b_j = 1|y = 1))^{1-b_j}, \quad (9)$$

ここで \$b\_j\$ は \$x\_j > 0\$ の場合 1、それ以外の場合は 0 をとる。このようにベルヌーイ試行を仮定する場合、各成長パターン \$j\$ の EF 値は考慮せず出現したか否かだけを考慮するため、パラメータ数は \$O(J)\$ まで減らすことができる。エッジの発生に対する成長パターン \$j\$ の出現の事後確率は以下のように推定できる。

$$P(b_j = 1|y = 1) = \frac{|\{\langle u, v \rangle \in \mathbf{E} | x_{uvj} > 0\}| + \delta}{|\mathbf{E}| + 2\delta}. \quad (10)$$

#### (3) 多項分布の仮定。

多項分布の仮定もベルヌーイ試行と同様文書分類においてよく用いられている<sup>11)</sup>。各ノードペアは \$\sum\_j x\_j\$ 回の独立した成長の結果だと仮定する。各成長はその成長パターンの出現確率 \$P(b\_j = 1|y)\$ で生じるものとしている。このとき

$$P(\mathbf{x}|y = 1) \propto \prod_{j=1}^J \frac{P(b_j = 1|y = 1)^{x_j}}{x_j!}. \quad (11)$$

この場合のパラメータ数も \$O(J)\$ で、頑健な推定が可能である。ここでエッジの発生

表 1 ソーシャルネットワークの基本的な統計量

Table 1 General statistics of social networks created from datasets.

	Nodes	Edges	Density	Mutual edges	Clustering coefficient
Blog	3,079	46,695	0.507%	17,469	30.38%
Flickr	3,721	175,402	1.331%	35,938	30.50%
Email	5,122	21,346	0.086%	654	4.511%

に対する成長パターン  $j$  の出現の事後確率は以下のように推定する.

$$P(b_j = 1 | y = 1) = \frac{\sum_{\langle u,v \rangle \in E} x_{uvj} + \delta}{\sum_{j'=1}^J \sum_{\langle u,v \rangle \in E} x_{uvj'} + \delta J}. \quad (12)$$

## 4. 実 験

### 4.1 データセット

提案法を評価するため、以下に示す 3 種類のソーシャルメディアから得たネットワークデータを対象として実験を行った。ただし、いずれのデータからも接続しているエッジ数が 3 未満のノードは省いた。

- (1) *Blog*. 本データセットは 2008 年 9 月 21 日から 2009 年 5 月 1 日の期間クローリングした日本のブログ記事に基づいており、各ノードはブログサイト、各エッジはトラックバックを表す。
- (2) *Flickr*. 本データセットは Flickr API<sup>\*1</sup> を用いてクローリングした Flickr のメタデータのデータベースであり、各ノードはユーザで、各エッジはコメントのやりとりを表す。
- (3) *Email*. 本データセットは Enron 社のメールログであり、研究目的でのみ利用が許可されている<sup>\*2</sup>。各ノードはユーザで、各エッジはメールのやりとりを表す。

表 1 に各ソーシャルネットワークの統計量を示す。Flickr が他のメディアと比較して密なネットワークで、また Email が比較的疎なネットワークであることがわかる。

### 4.2 ネットワークモチーフ分析と動的ネットワークモチーフ分析

今回の実験ではサイズ 3 の部分ネットワークに注目する。

まず、図 3 は 3 つのソーシャルネットワークのネットワークモチーフプロフィールであ

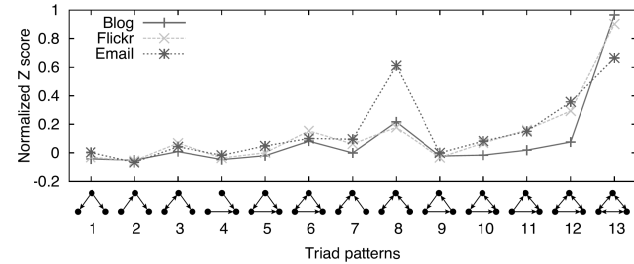


図 3 ソーシャルネットワークのネットワークモチーフプロフィール

Fig.3 Network motif profiles for social networks.

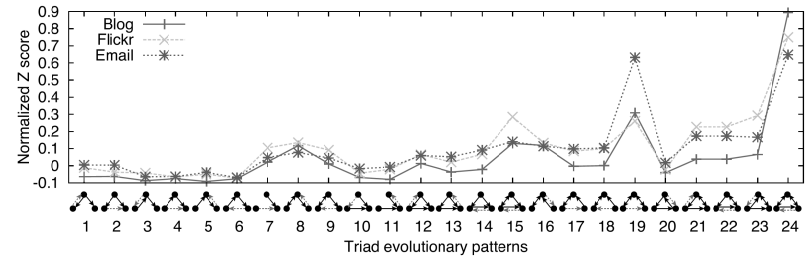


図 4 ソーシャルネットワークの動的ネットワークモチーフプロフィール

Fig.4 Dynamic network motif profiles for social networks.

る。Clique (図 3 中のトライアド 13) に着目すると、Blog, Flickr, Email いずれにおいても 3 人のユーザがお互いにコミュニケーションを取り合っていることが多い。また two mutual dyads (図 3 中のトライアド 8) に着目すると、Email は他の 2 つのメディアに較べて 2 人と相互にインタラクションしているユーザは存在するがその 2 人に直接インタラクションがないことが多い。これについては、例えば社内のメールでは上司と部下は連絡を取り合うが部下の間では連絡を取り合わないというような直観的な説明が可能である。

次に、図 4 は 3 つのソーシャルネットワークの動的ネットワークモチーフプロフィールである。今回は静的なネットワークに対して行った分析であったため、得られた動的ネットワークモチーフプロフィールがネットワークモチーフプロフィールとほとんど変わらないと

\*1 <http://www.flickr.com/services/api>

\*2 <http://www.cs.cmu.edu/~enrcn>

いう結果になった。したがって、時刻  $t$  までのネットワーク  $G$  に対する時刻  $t$  から  $t+1$  までのエッジの発生を対象とした、動的なネットワークに対する分析をする必要がある。

### 4.3 動的ネットワークモチーフに基づくリンク予測

#### 4.3.1 評価方法

動的ネットワークモチーフのリンク予測精度を評価するために、10 分割交差検定を行った。すなわち、データセット中のエッジを 10 分割し、9 つをトレーニングセット  $E$  とし、残りの 1 つをテストセット  $E' - E$  とした。

評価手順としては各リンク予測手法にて  $V \times V - E$  中の各ノードペア  $\langle u, v \rangle$  に対してエッジが発生する確率  $P(y = 1|u, v)$  を算出し、上位  $|E' - E|$  件のノードペア集合  $\hat{Y}$  を得る。

評価尺度は Liben-Nowell ら<sup>9)</sup> と同様で、以下のように各リンク予測器の出力中の正答率をまったくランダムなリンク予測の正答率の理論値と比較することで得られる。

$$\text{Performance} = \frac{|\hat{Y} \cap (E' - E)|}{|\hat{Y}|} / \frac{|V \times V \cap (E' - E)|}{|V \times V|}. \quad (13)$$

#### 4.3.2 比較手法

以下に示す 7 つの手法を比較し、提案手法の有効性を示す。

##### (1) *Our method.* 提案手法.

$$P(y = 1|u, v) = P(y = 1|x_{uv}), \quad (14)$$

右辺は式 (6) によって計算する。今回の実験では扱う部分ネットワークのノード数は 3 とし、スムージングパラメータは  $\delta = 1$  とした。

##### (2) *Common neighbors.* 共通隣接ノード数. Newman ら<sup>9)</sup> は 2 ノードの近接性をそれらに共通して隣接するノード数で与えた。

$$P(y = 1|u, v) \propto |\Gamma_u \cap \Gamma_v|, \quad (15)$$

ただし  $\Gamma_u = \{w \in V | \langle u, w \rangle \in E\} \cup \{w \in V | \langle w, u \rangle \in E\}$ 、すなわち  $u$  に隣接するノード集合を意味する。

##### (3) *Jaccard.* Liben-Nowell ら<sup>6)</sup> は 2 ノード間の近接性を共通隣接ノードの数ではなくジャカード係数により与えることで、次数 (隣接するノード数) が小さいノードであってもリンク予測が可能な手法を提案した。

$$P(y = 1|u, v) \propto \frac{|\Gamma_u \cap \Gamma_v|}{|\Gamma_u \cup \Gamma_v|}. \quad (16)$$

##### (4) *Adamic/Adar.* 次数 (隣接するノード数) が大きいノードはいずれのノードとも隣接している可能性が高い。そこで、Adamic と Adar<sup>1)</sup> は各共通隣接ノードにその次数にペナルティを与えることでより高精度なリンク予測器を実現した。

$$P(y = 1|u, v) \propto \sum_{w \in \Gamma_u \cap \Gamma_v} \frac{1}{\log |\Gamma_w|}. \quad (17)$$

##### (5) *Preferential attachment.* 優先的選択指標. Barabasi ら<sup>2)</sup> と Newman ら<sup>9)</sup> によって提案された手法で、ネットワークのスケールフリー性から導かれる「次数が多いノードほど隣接しやすい」という仮定に基づいている。

$$P(y = 1|u, v) \propto |\Gamma_u| \cdot |\Gamma_v|. \quad (18)$$

##### (6) *Katz.* 2 ノード間にパスが存在すればするほど、それらは近接していると考えられる。また、それらのパス長が短ければ短いほど近接している。そこで Katz<sup>5)</sup> は 2 ノードの近接性を、それらの間のパスについて長さで指数的に減衰する重みの総和で与えた。

$$P(y = 1|u, v) \propto \sum_{l=1}^{\infty} \beta^l |P_{uv}^{(l)}|, \quad (19)$$

ただし  $P_{uv}^{(l)}$  は  $u$  から  $v$  への長さ  $l$  のパスの集合を意味する。重み  $\beta$  は  $\{0.05, 0.005, 0.0005\}$  の 3 つの候補から、トレーニングセットを使ってさらに交差検定を行い、リンク予測の performance を最大にするものを選択した。

##### (7) *Random walk with restart.* Tong ら<sup>12)</sup> が採用している手法で、PageRank の概念に基づいた手法である。 $u$ から出発して、 $1 - \alpha$ の確率で今いるノードに隣接するいずれかのノードに移り $\alpha$ の確率で $u$ に戻るというランダムウォークを繰り返す。このときの $v$ の滞在確率を $P(y = 1|u, v)$ とみなす。 $\alpha$ は $\{0.01, 0.05, 0.15, 0.30, 0.50\}$ の 5 つの候補から、トレーニングセットを使ってさらに交差検定を行い、リンク予測の performance を最大にするものを選んだ。

### 4.3.3 結果

表 2 に比較手法のリンク予測の性能を示す。太字は各ネットワークに対するもっとも高いリンク予測精度を表している。Blog と Email においては、提案法は他の比較手法よりも統計



表 2 ソーシャルネットワークに対する比較手法の性能  
Table 2 Performance of compared methods for social networks.

Predictor	Blog	Flickr	Email
Probability that a random prediction is correct	0.052%	0.136%	0.009%
<i>Our method</i>			
conditional independence	792.2	74.2	391.5
multi-variate Bernoulli	<b>816.2</b>	74.8	<b>409.6</b>
multinomial, EF attributes	635.4	83.9	101.3
<i>Common neighbors</i>	279.6	85.2	234.4
<i>Jaccard's coefficient</i>	69.0	77.8	8.0
<i>Adamic/Adar</i>	289.5	<b>85.5</b>	213.9
<i>Preferential attachment</i>	61.5	76.1	72.9
<i>Katz</i>	$\beta = 0.05$	156.4	51.1
	$\alpha = 0.01$	76.0	7.5
<i>Random walk with restart</i>			318.8

的に有意に高精度であった (符号検定,  $p < .01$ ). 一方で Flickr においては Adamic/Adar が最も高精度であったものの, 提案法の精度も比較的高いという結果になった. このように, 動的ネットワークモチーフに基づくリンク予測は有効であることがわかった.

また提案法に着目すると, ベルヌーイ試行の仮定に基づく手法はどのネットワークに対しても有効だが, 多項分布の仮定に基づく手法は Email のような疎なネットワークでは有効ではないが Blog, Flickr の密なネットワークで有効であることがわかった.

## 5. まとめと今後の課題

動的ネットワークモチーフというネットワークの局所構造の特徴的な成長に着目し, 複雑ネットワークを分析する手法を提案した. ソーシャルネットワークに提案法を適用した結果, 提案法の結果はネットワークモチーフ分析の結果とあまり変わらなかった. これは, 実験の設定が静的なネットワークに対する分析になっていたことが原因であった.

また, 動的ネットワークモチーフのリンク予測への応用をベイズの定理により実現した. 本稿ではソーシャルネットワークを用いた実験を通じて, 提案法が従来法よりも高精度にリンク予測できることを示した.

本稿では静的なネットワークに対してのみ分析を行ったが, 時間とともに成長するような動的なソーシャルネットワークに適用し分析を行う必要がある. 具体的には, SNS におけるユーザ間のつながりが時間とともにどう成長していくのか提案法により分析を行うことを検討している. また, 今回は 3 ノードに関する成長パターンに着目して分析を行ったが, 4 ノード以上の部分ネットワークに関する分析も行っていく予定である.

## 参考文献

- 1) Adamic, L. and Adar, E.: Friends and neighbors on the web, *Social Networks*, Vol.25, No.3, pp.211–230 (2003).
- 2) Barabási, A., Jeong, H., Néda, Z., Ravasz, E., Schubert, A. and Vicsek, T.: Evolution of the social network of scientific collaborations, *Physica A: Statistical Mechanics and its Applications*, Vol.311, No.3–4, pp.590–614 (2002).
- 3) Domingos, P. and Pazzani, M.: On the optimality of the simple Bayesian classifier under zero-one class, *Machine Learning*, Vol.29, No.2, pp.103–130 (1997).
- 4) Inokuchi, A. and Washio, T.: A fast method to mine frequent subsequences from graph sequence data, *Proceedings of the 8th IEEE International Conference on Data Mining*, pp.303–312 (2008).
- 5) Katz, L.: A new status index derived from sociometric analysis, *Psychometrika*, Vol.18, No.1, pp.39–43 (1953).
- 6) Liben-Nowell, D. and Kleinberg, J.: The link-prediction problem for social networks, *Journal-American Society for Information Science and Technology*, Vol.58, No.7, pp.1019–1031 (2007).
- 7) Milo, R., Itzkovitz, S., Kashtan, N., Levitt, R., Shen-Orr, S., Ayzenshtat, I., Sheffer, M. and Alon, U.: Superfamilies of evolved and designed networks, *Science*, Vol.303, No.5663, pp.1538–1542 (2004).
- 8) Milo, R., Shen-Orr, S., Itzkovitz, S., Kashtan, N., Chklovskii, D. and Alon, U.: Network motifs: simple building blocks of complex networks, *Science*, Vol.298, No.5594, pp.824–827 (2002).
- 9) Newman, M.: Clustering and preferential attachment in growing networks, *Physical Review E*, Vol.64, No.2, p.025102 (2001).
- 10) Nijssen, S. and Kok, J.: A quickstart in frequent structure mining can make a difference, *Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp.647–652 (2004).
- 11) Schneider, K.: On word frequency information and negative evidence in Naive Bayes text classification, *Proceedings of the 4th International Conference on Advances in Natural Language Processing*, Berlin, Heidelberg, pp.474–485 (2004).
- 12) Tong, H., Faloutsos, C. and Pan, J.: Random walk with restart: fast solutions and applications, *Knowledge and Information Systems*, Vol.14, No.3, pp.327–346 (2008).
- 13) Yan, X. and Han, J.: gSpan: Graph-based substructure pattern mining, *Proceedings of the 2002 IEEE International Conference on Data Mining*, pp.721–724 (2002).
- 14) 甲谷 優, 川島晴美, 藤村 考: QA コミュニティの成長パターンに基づく回答者への質問推薦, 日本データベース学会論文誌, Vol.8, No.1, pp.89–94 (2009).