プライベートなノードを含むソーシャルネットワークの統計量推定

中嶋 一貴† 首藤 一幸†

†東京工業大学 情報理工学院 数理・計算科学系

あらまし 本研究では、プライベートなノードを含むソーシャルネットワークに対する、ランダムウォークを用いた 統計量推定アルゴリズムを議論する.近年では、ランダムウォークを用いた、ソーシャルネットワークの統計量の効 率的な推定アルゴリズムが多く提案されている.現実のソーシャルネットワークは隣接ノードが非公開であるプライ ベートなノードが一定の割合で存在するが、既存の推定アルゴリズムはこれらの存在を考慮していない.特に、平均 次数とノード数に対する既存の推定アルゴリズムは、グラフ内のプライベートなノードの割合*p*が増加するにつれて、 推定値の誤差が増加する問題点がある.本論文では*p*にほぼ依存しない、グラフの平均次数とノード数の推定アルゴ リズムをそれぞれ提案する.我々はサンプルされた各ノードに既存の推定アルゴリズムと異なる重みをかけることで、 *p*に対して頑健な推定アルゴリズムを実現する.実験では、提案アルゴリズムが*p*にほぼ依存せずに、高い精度で平 均次数とノード数を推定できることを示す.

キーワード ソーシャルネットワーク, グラフサンプリング, ランダムウォーク, プライバシ

1 はじめに

近年,Online Social Networks (OSNs)の解析を目的とする 研究が盛んに行われている.OSNsの解析は平均次数やノード 数などのグラフの統計量や,中心性などのノードの統計量を推 定することである [1-4].しかし,ほとんどのOSN はそのグラ フデータへのアクセスを制限しており,推定することは容易で はない.巨大なOSN の全てのグラフデータを取得して統計量 を厳密計算することはほぼ不可能であり,また事前にノードの 分布を知ることもできないためランダムサンプリングに基づく 推定アルゴリズムを適用することも難しい [5,6].

この問題を打開するために、ランダムウォークを用いた統 計量推定アルゴリズムが多く提案されてきた [7-14]. 多くの OSN のアプリケーション・プログラミング・インターフェー ス (API) は、クエリを実行したノードの隣接ノード情報を提供 する. API から提供された隣接ノードのインデックスをランダ ムに 1 つ選び遷移することを繰り返すことで、OSN 上でラン ダムウォークを実行してインデックスなどのグラフデータのサ ンプリングが可能である.また、ランダムウォークのマルコフ 性から各ノードのサンプルの分布を計算できるため、サンプル ノードに適切な重み付けを行うことで不偏な推定量を求めるこ とができる [6,4,15].既存の統計量推定アルゴリズムは、ごく 限られた API へのクエリ回数で統計量を高精度に推定するこ とを目標として設計されてきた.

OSN に既存の統計量推定アルゴリズムを適用する上で,ユー ザのプライバシ保護によるランダムウォークの制限は無視でき ない問題である.ほとんどの OSN では,各ユーザは友人リス トの公開・非公開を設定することができる.Dey らは,ニュー ヨークの Facebook ユーザの 52.6%が友人リストを非公開にし ていたと報告した [16]. 本論文では隣接ノード集合を公開・非 公開にするノードをそれぞれパブリック・プライベートなノー ドと呼ぶ. ランダムウォークでプライベートなノードに遷移し た場合は, API がその隣接ノードを提供しないためアルゴリズ ムを続行できない.

Gjoka らは,実際に Facebook 上でランダムウォークを実行 したとき,パブリックな隣接ノードに遷移する対処をとった [4]. このランダムウォークは,グラフ全体ではなく,グラフ内のパ ブリックなノードのみから成る連結な部分グラフ上で実行され ている.このため,Gjoka らのランダムウォークを既存の統計 量推定アルゴリズムに適用すると,その部分グラフの統計量を 推定することになる.

平均次数とノード数に対する既存の推定アルゴリズムは,プ ライベートなノードの割合 p が増加するにつれて推定値の誤差 が大きくなる問題点がある.これは,ランダムウォークが実行 される部分グラフの平均次数とノード数が,p が増加するにつ れて真値との誤差が大きくなるためである.

本研究の貢献は, p に対して頑健な, 平均次数とノード数の 推定アルゴリズムを提案することである.まず, Gjoka らのラ ンダムウォーク [4] の定常分布を導出し,各ノードがそのパブ リックな隣接ノード数に比例した分布でサンプルされることを 示す.各ノードのサンプルの分布を計算することで,各サンプ ルノードへの適切な重み付けを議論することが可能になる.そ して我々は,各サンプルノードに既存の推定アルゴリズムと異 なる重みをかけることで, p に対して頑健な推定を可能にする. 本研究では,OSN の API が提供する隣接ノード情報の2つの モデルを考え,それぞれのモデルでクエリ回数を考慮した推定 アルゴリズムを設計する.実際のOSN のグラフに対する実験 では,いずれのモデルに対しても,提案アルゴリズムがpにほ ぼ依存せずに,ごく限られたクエリ回数で平均次数とノード数 を高精度に推定できることを示す.

This is an unrefereed paper.



(a) API が隣接ノードのプライバシ (b) API が隣接ノードのプライバ
情報を提供するモデル.
シ情報を提供しないモデル.

図 1 パブリックなノード v_i にクエリを実行したとき, API が提供 する隣接ノード情報の, 2 つのモデル: 青色のノードはパブリッ ク,赤色のノードはプライベート, 白色のノードはプライバシ情 報が不明なノードを示す.

2 準 備

本研究では無向で連結なグラフG = (V, E)を扱う. ソー シャルネットワークでは、ノードはユーザ、エッジはユーザ間 の関係を表す. グラフGはセルフループと多重辺を持たず辺 に重みはないとする. ノード数n = |V|,エッジ数m = |E|と し、グラフGのノード集合を $V = \{v_1, v_2, ..., v_n\}$,エッジ集合 をE, ノード v_i の次数を d_i と表記する. グラフGの平均次数 は以下のように定義される:

定義 1. $d_{avg} = \frac{\sum_{v_i \in V} d_i}{n} = \frac{2m}{n}$.

研究対象とするネットワークはランダムウォークモデル [9] を満たすとする.事前にグラフの情報は持たず,OSNの API が提供する隣接ノード情報を辿ってサンプリングを行う.ま た,推定アルゴリズムを実行している間,グラフは静的である とする.

グラフ G の各ノード v_i はパブリックまたはプライベートの いずれかのラベルを持つ. ノード v_i に対してクエリを実行し たとき, v_i がパブリックであれば,その隣接ノード情報を取得 でき, v_i がプライベートであれば,その隣接ノード情報は取得 できないとする^{1:}.

本研究では、図1のように、OSNのAPIが提供する隣接ノー ド情報の2つのモデルを考える. Gjoka らが Facebook 上でラ ンダムウォークを実行した事例は、図1(a)のモデルに該当す る[4]. しかし、図1(b)のように、隣接ノードのインデックス のみ提供される場合も考えられる.

ランダムウォークを用いた推定アルゴリズムの目標は、少 ないクエリ回数で OSN のグラフの統計量を推定することであ る [5,17]. 一度クエリを実行したノードの隣接ノードの情報は 保存する.

 $P = \{p_{ij}\}_{i,j\in S}$ を有限な状態空間 S におけるマルコフ連鎖 の遷移確率行列とする. Pの定常分布について以下の定理が成 り立つ.

定理 1. *[18]***P** がエルゴード的であるとき,その定常分布 π が 一意に存在する.

1:出力はエラーメッセージや空集合など、OSN の API に依存する.

3 プライベートなノードを考慮したランダムウ ォーク

プライベートなノードを含むグラフ上のランダムウォークを 議論する.ランダムウォークは、あるノードの隣接ノードから ランダムに1つのノードを選択して遷移することを繰り返し、 遷移したノードのインデックスをはじめとする情報をサンプ リングする.プライベートなノードに遷移した場合はその隣接 ノードの情報が取得できないため、ランダムウォークを続行す ることができない.

このため、Gjoka らが提案したパブリックな隣接ノードをラ ンダムに選び遷移するランダムウォークが適切である [4]: ラ ンダムウォークのiステップ目のパブリックなサンプルノード v_{x_i} の隣接ノード集合からランダムに 1 つ選んだノード $v_{x_{i+1}}$ のプライバシ情報を確認し、 $v_{x_{i+1}}$ がパブリックなノードであ ればそのノードに遷移し、プライベートなノードであれば v_{x_i} の隣接ノードを選択し直す. API が隣接ノードのプライバシ情報を提供しないモデルでは、 $v_{x_{i+1}}$ にクエリを実行してその出 力からプライバシ情報を判別する.

Gjoka らのランダムウォークは、元のグラフ G からプライ ベートなノードを削除してできるいくつか部分グラフのうち、 初期ノードが属する部分グラフ上で実行される。 グラフ G か ら全てのプライベートなノードとそれに接続するエッジを削除 したあとに分割される k 個の連結な部分グラフを C_i とする。 それぞれの部分グラフ C_i をクラスタと呼ぶ [19,20]. この中で ノード数が最大のクラスタを $C^* = (V^*, E^*)$ とする。 C^* に対 して $n^* = |V^*|, m^* = |E^*|$ とおく。最大クラスタ C^* の平均次 数を以下のように定義する:

定義 2. $d_{avg}^* = \frac{\sum_{v_i \in V^*} d_i^*}{n^*} = \frac{2m^*}{n^*}.$

本研究ではランダムウォークの初期ノードは最大クラスタ に属するノードを選ぶとする. ほとんどの OSN はスケールフ リーネットワークに分類され [1-4], スケールフリーネットワー クはランダムなノードの除去に対して頑健であることが知られ ている [19,20]. すなわち, グラフにプライベートなノードがラ ンダムに分布していると仮定すると, ほとんどのパブリックな ノードは最大クラスタに属される. このため, 既存研究でも最 大クラスタでサンプリングが可能であった [1-4].

3.1 最大クラスタ上のランダムウォーク

最大クラスタ上のランダムウォークでは、各ノードはそのパ ブリックな隣接ノード数に比例した分布でサンプルされる.こ れは最大クラスタ上のランダムウォークの定常分布を導出する ことでわかる. C^* 上のランダムウォークは有限な状態空間 V^* におけるマルコフ連鎖でありエルゴード性を満たす.よって、 定理1より定常分布 π^* が一意に存在する. C^* における各ノー ド $v_i \in V^*$ の次数を d_i^* と表す. C^* の各ノードの次数は、その ノードのGにおけるパブリックな隣接ノード数である.ここ で V^* の各ノードの添字集合を $I = \{i_1, i_2, ..., i_n^*\}$ とおく. C^* 上のランダムウォークの定常分布を $\pi^* = \{\pi_i^*\}_{1 \le j \le n^*}$ とする



図 2 プライベートなノードを含むグラフの例. 青色のノードはパブ リック, 赤色のノードはプライベートなノードである.

と、各ノード $v_{i_j} \in V^*$ に対して以下が成り立つ:

$$\pi_j^* = \frac{d_{i_j}^*}{2m^*}$$

例えば図 2 のグラフで $v_i = i$ ($0 \le i \le 9$) とする.以下の 3 つのクラスタ C_1, C_2, C_3 に分割され,最大クラスタは $C^* = C_1$ である:

- $C_1 = (\{3, 4, 5, 6, 8\}, \{(3, 4), (4, 5), (4, 6), (4, 8)\})$
- $C_2 = (\{7,9\},\{(7,9)\})$
- $C_3 = (\{2\}, \{\}).$

添字集合 $I = \{i_1, i_2, i_3, i_4, i_5\}$ を $i_1 = 3, i_2 = 4, i_3 = 5, i_4 = 6, i_5 = 8$ とすると、 $d_3^* = 1, d_4^* = 4, d_5^* = 1, d_6^* = 1, d_8^* = 1$ である。よって C^* 上のランダムウォークの定常分布は $\pi^* = (\frac{1}{8}, \frac{1}{2}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8})$ である。

4 パブリックな隣接ノード数の計算

平均次数とノード数の推定において、各サンプルノード v_{x_i} に重み付けをするために、そのパブリックな隣接ノード数 $d_{x_i}^*$ を計算する必要がある。APIが提供する隣接ノード情報のそれ ぞれのモデル (図 1) に対して、クエリ数を考慮した計算方法を 述べる.

4.1 API が隣接ノードのプライバシ情報が提供するモデル

API が v_{x_i} の隣接ノードのプライバシ情報を提供する場合, $d_{x_i}^*$ を厳密に計算できる.実際に、Gjoka らが Facebook上で ランダムウォークを実行した事例では、このモデルに該当す る [4].

4.2 API が隣接ノードのプライバシ情報が提供するモデル

この場合, $d_{x_i}^*$ の厳密値を計算するには v_{x_i} の全ての隣接ノードに対してクエリを実行する必要があり,現実的ではない.こ のため,ランダムウォークでサンプルした各ノードに対して, パブリックな隣接ノードを選択できた割合から近似計算する方 法を提案する.各サンプルノード v_{x_i} に対して,パブリックな 隣接ノードを選択した合計回数を a_{x_i} ,隣接ノードを選択した 合計回数を b_{x_i} とする.各サンプルノード v_{x_i} のパブリックな 隣接ノード数 $d_{x_i}^*$ の近似値 $\hat{d}_{x_i}^*$ を以下のように定義する:

$$\hat{d}_{x_i}^* = \begin{cases} d_{x_i} \frac{a_{x_i}}{b_{x_i}} & (b_{x_i} \neq 0) \\ \\ d_{x_i} & (otherwise) \end{cases}$$

定理 2. 近似値 $\hat{d}_{x_i}^*$ の期待値は真値 $d_{x_i}^*$ に等しい.

 $E[\hat{d}_{x_i}^*] = d_{x_i}^*$

証明. i < x > y > y目のサンプルノードの隣接ノードを1つラン ダムに選んだとき、確率 $\frac{d_{x_i}^*}{d_{x_i}}$ でパブリックな隣接ノードを選択 する. よって、 $E[\hat{d}_{x_i}^*] = d_{x_i} \frac{d_{x_i}^*}{d_{x_i}} = d_{x_i}^*$ が成り立つ.

5 プライベートなノードを考慮した統計量推定

本論文ではグラフ G の平均次数 davg とノード数 n を推定す るアルゴリズムを議論する. グラフのノード数とエッジ数は最 も基本的なグラフの統計量である. 定義1よりノード数と平均 次数の推定値からエッジ数の推定値を求めることができる.

既存の推定アルゴリズムは C* の平均次数とノード数を推定 しており, pが増加するにつれて推定値の誤差が増加する問題 点がある.提案アルゴリズムは,既存アルゴリズムのサンプル ノードへの重みを変更することで,pの増加に対する推定値の 誤差の増加を大幅に小さくする.

5.1 平均次数の推定

Dasgupta らはグラフの平均次数を推定するアルゴリズム Smooth を提案した [9]. Smooth は、入力として平均次数の 大まかな推定値 cを与え、それより正確な推定値 d_{avg}^{Smooth} を出 力する。ランダムウォークで遷移した各ノードにセルフループ を c本追加して次の遷移先のノードを選択する。

本研究ではグラフの事前知識が無いことを仮定するためc = 0とする.入力として与える定数cは平均次数に対する大まかな 近似値とされるが、グラフに対する事前知識が無い場合はその近 似値を求めることにクエリの追加実行を必要とする.Dasgupta らは定数c=0, 1, 5, 50に対して推定精度の評価実験をしてお り、0や1の小さい値が望ましいとしている [9].特にc = 0の 場合は Gjoka らが Facebook の平均次数を推定したアルゴリズ ムと一致する [4].

C^{*} 上で *Smooth* を適用すれば推定値は以下のように計算される:

$$d_{avg}^{Smooth} = \frac{r}{\sum_{i=1}^{r} \frac{1}{d_{x_i}^*}}$$

定理 3. 推定値 d_{avg}^{Smooth} は d_{avg}^* に収束する.

証明. *R*は有限な状態空間 *V** におけるエルゴード的なマルコ フ連鎖であるから,定常分布 π* を持つ. $\Phi_{avg}^{Smooth} = \frac{1}{r} \sum_{i=1}^{r} \frac{1}{d_{x_i}^*}$ とおく. $E[\Phi_{avg}^{Smooth}] = \frac{1}{d_{avg}^*}$ を示す.

$$E[\Phi_{avg}^{Smooth}] = E[\frac{1}{d_{x_k}^*}] = \sum_{v_i \in V^*} \pi_i^* E[\frac{1}{d_{x_k}^*} | x_k = i]$$
$$= \sum_{v_i \in V^*} \frac{d_i^*}{2m^*} \frac{1}{d_i^*} = \frac{1}{d_{avg}^*}$$

1 つ目の等号は期待値の線型性より、2 つ目の等号は繰り 返し期待値の法則より成り立つ. $d_{avg}^* = \frac{1}{E[\Phi_{ang}^{smooth}]}$ である. Φ_{avg}^{Smooth} は期待値に収束するので d_{avg}^{Smooth} は d_{avg}^{*} に収束する.

本論文では各サンプルノード v_i にかける重みを $\frac{1}{d_i}$ に変更した推定アルゴリズムを提案する.提案アルゴリズムによる d_{avg} の推定値 \hat{d}_{avg} を以下のように定義する:

$$\hat{d}_{avg} = \frac{r}{\sum_{i=1}^{r} \frac{1}{d_{x_i}}}$$

定理 4. \hat{d}_{avg} は $\tilde{d}_{avg} = rac{2m^*}{\sum_{v_i \in V^*} rac{d_i^*}{d_i}}$ に収束する.

証明. $\hat{\Phi}_{avg} = \frac{1}{r} \sum_{i=1}^{r} \frac{1}{d_{x_i}}$ とおく. 定理 3 の証明と同様にして,

$$E[\hat{\Phi}_{avg}] = \sum_{v_i \in V^*} \frac{d_i^*}{2m^*} \frac{1}{d_i} = \frac{1}{\tilde{d}_{avg}}$$

 $\tilde{d}_{avg} = \frac{1}{E[\Phi_{avg}]}$ である. $\hat{\Phi}_{avg}$ は期待値に収束するので \hat{d}_{avg} は \tilde{d}_{avg} に収束する.

5.2 ノード数の推定

Katzir と Hardiman はサンプルノードの衝突に基づく, グ ラフのノード数を推定するアルゴリズムを提案した [7]. この 推定アルゴリズムは, ランダムウォークでサンプルしたノード 集合 $R = \{v_1, v_2, ..., v_r\}$ に対してある定数 m 以上離れたイン デックスのペアの集合 I を調べる:

$$I = \{ (k, l) \mid m \le |k - l| \land 1 \le k, l \le r \}$$

インデックスを *m* だけ離すことで,それぞれのノードが定常 分布 π^{*} から独立にサンプルしたとみなすことができる [7].

各ノードペア (v_{x_k}, v_{x_l}) に対して $v_{x_k} = v_{x_l}$ のとき 1 を返し そうでなければ 0 を返す変数 $\phi_{k,l}$ を定義する:

$$\phi_{k,l} = \mathbf{1}_{\{x_k = x_l\}}$$

ここで各ノードペア $(v_{x_k}, v_{x_l}) \in I$ に対して重み付けした $d_{x_k}^* \phi_{k,l} \geq \frac{(d_{x_k}^*)^2}{d_{x_l}^*}$ のそれぞれの平均値 Φ_n^{KH}, Ψ_n^{KH} を以下のよう に定義する:

$$\Phi_n^{KH} = \frac{1}{|I|} \sum_{(k,l)\in I} d_{x_k}^* \phi_{k,l}$$
$$\Psi_n^{KH} = \frac{1}{|I|} \sum_{(k,l)\in I} \frac{(d_{x_k}^*)^2}{d_{x_l}^*}$$

 C^* 上で Katzir と Hardiman の推定アルゴリズムを適用すれば、その推定値 n^{KH} は以下のように定義される:

$$n^{KH} = \frac{\Psi_n^{KH}}{\Phi_n^{KH}}$$

定理 5. n^{KH} は n^{*} に収束する.

証明. まず
$$E[\Phi_n^{KH}] = \sum_{v_j \in V^*} d_j^* \left(\frac{d_j^*}{2|E^*|}\right)^2$$
 を示す.
 $\Phi_n^{KH} = E[d_{x_k}^* \phi_{k,l}] = \sum_{v_j \in V^*} d_j^* \left(\frac{d_j^*}{2|E^*|}\right)^2$

一つ目の等号は期待値の線型性より成り立つ.2つ目の等号は $x_k = j, x_l = j$ である確率がともに $\frac{d_j^*}{2|E^*|}$ であり、インデック スが*m*だけ離れていることでそれぞれの事象は独立に起こる ことより成り立つ.

次に
$$E[\Psi_n^{KH}] = n^* \sum_{v_j \in V^*} d_j^* \left(\frac{d_j^*}{2|E^*|}\right)^2$$
を示す.

$$\begin{split} E[\Psi_n^{KH}] &= E[\frac{(d_{x_k}^*)^2}{d_{x_l}^*}] = \sum_{v_i \in V^*} \sum_{v_j \in V^*} \frac{(d_j^*)^2}{d_i^*} \frac{d_i^*}{2|E^*|} \frac{d_j^*}{2|E^*|} \\ &= n^* \sum_{v_j \in V^*} d_j^* \left(\frac{d_j^*}{2|E^*|}\right)^2 \end{split}$$

 $n^{*} = \frac{E[\Psi_{n}^{KH}]}{E[\Phi_{n}^{KH}]}$ である. $\Phi_{n}^{KH}, \Psi_{n}^{KH}$ は期待値に収束するので n^{KH} は n^{*} に収束する.

我々は Φ_n^{KH} と Ψ_n^{KH} の重み付けを変更した $\hat{\Phi}_n$ と $\hat{\Psi}_n$ から 推定値を求める.

$$\hat{\Phi}_n = \frac{1}{|I|} \sum_{(k,l) \in I} d_{x_k} \phi_{k,l}$$
$$\hat{\Psi}_n = \frac{1}{|I|} \sum_{(k,l) \in I} \frac{(d_{x_k})^2}{d_{x_l}^*}$$

ノード数の推定値 *î*を以下のように定義する.

$$\hat{n} = \frac{\hat{\Psi}_n}{\hat{\Phi}_n}$$

定理 6. \hat{n} は $\tilde{n} = \frac{n^* \sum_{v_j \in V^*} d_j^* (d_j)^2}{\sum_{v_j \in V^*} d_j (d_j^*)^2}$ に収束する.

証明. 定理5の証明と同様にして以下が示される.

$$E[\hat{\Phi}_n] = \sum_{v_j \in V^*} d_j \left(\frac{d_j^*}{2|E^*|}\right)^2$$
$$E[\hat{\Psi}_n] = n^* \sum_{v_j \in V^*} d_j^* \left(\frac{d_j}{2|E^*|}\right)^2$$

 $\tilde{n} = \frac{E[\hat{\Psi}_n]}{E[\hat{\Phi}_n]}$ である. $\hat{\Phi}_n, \hat{\Psi}_n$ は期待値に収束するので \hat{n} は \hat{n} に 収束する.

API が隣接ノードのプライバシ情報を提供しないモデルで は、それぞれの推定アルゴリズムにおける各サンプルノード v_{x_i} の重みに用いる厳密値 $d_{x_i}^*$ の代わりに近似値 $\hat{d}_{x_i}^*$ を用いる. $\hat{d}_{x_i}^*$ は4章に基づいて近似計算する.このモデルでも、定理2 より、定理4と定理6が成り立つ.

6 実 験

実際の OSN のグラフを用いて提案アルゴリズムを評価する. 対象とするグラフが有向グラフであれば反対向きのエッジを追 加することで無向グラフとし,非連結なグラフはそのグラフの 最大連結成分を扱った.また,セルフループと多重辺は削除し た.本論文で扱う各グラフのノード数,平均次数を表6に示す. API が提供する隣接ノード情報のそれぞれのモデル(図1)

Network	n	d_{avg}
LiveMocha [21]	104103	42.13
douban $[21]$	154908	4.22
gowalla [21]	196591	9.67
academia [21]	200167	10.22
Pokec [22]	1632803	27.32

表1 データセット

において,提案アルゴリズムの頑健性の評価と,実際にプライ ベートなノードを含む Pokec グラフ [24] における推定精度の 評価を行う.推定精度の評価に正規化平方二乗誤差 (NRMSE) を用いる.NRMSE は推定値の誤差と分散を評価できる指標で あり,関連する多くの研究で用いられている [15,7,11,23].

6.1 提案アルゴリズムの頑健性

提案アルゴリズムの頑健性を評価する。グラフ内のプライ ベートなノードの割合が増加したとしても、NRMSE の増加が 小さいことが望ましい。グラフ内のプライベートなノードの割 合を増加させた時のそれぞれの推定アルゴリズムの NRMSE を 比較する。一方で、プライバシ保護などの観点からプライベート なノードを含む公開データセットは我々の知る限り Pokec のみ であった [24]. このため本論文では 4 つのグラフ LiveMocha, douban, gowalla, academia に対してパラメータ p を与えて, 割 合pのノードをランダムに選びそれらをプライベートなノード とする前処理を行った.そして、その最大クラスタからランダ ムに選んだ初期ノードとクエリ数を与えてランダムウォークを 実行し、それぞれの推定値を計算する. 上記のシミュレーショ ンを独立に 1000 回行い, それぞれのアルゴリズムの NRMSE を比較する。 パラメータ p は 0.0 から 0.30 まで 0.01 ごとに 変化させた. p = 0.0 のときそれぞれのアルゴリズムは等しい NRMSE を持つ.また、ノード数推定アルゴリズムにおける定 数 m は、サンプル列の長さの 2.5%と設定した.

まず, API が隣接ノードのプライバシ情報を提供するモデ ルにおける提案アルゴリズムの頑健性を議論する. 平均次数の 推定アルゴリズムの頑健性を評価する. 図 3 と図 4 は各グラ フにおけるそれぞれクエリ数をノード数の 1%, 5%とした時の Dasgupta らのアルゴリズム (Smooth) [9] と提案アルゴリズ ム (Ours) のパラメータ p に対する NRMSE を示す. いずれ のクエリ数に対しても,全てのグラフにおいて Smooth はパ ラメータ p が増加するにつれて NRMSE が増加している. 一 方で, Ours の NRMSE は p が増加してもほとんど増加してい ない.

ノード数の推定アルゴリズムの頑健性を評価する.図5と図 6 はそれぞれクエリ数をノード数の1%,5%とした時の Katzir と Hardiman のアルゴリズム (KH)[7]と提案アルゴリズム (Ours)のパラメータpに対する NRMSE を示す.ノード数衝 突に基づくアルゴリズムはクエリ数が極端に少ない場合,サン プル内のノードの衝突がほとんど生じないため推定精度が低く なる.このため,クエリ数が1%のとき両方のアルゴリズムの

	収束値	相対誤差
\mathbf{Smooth}	19.49	0.287
Ours	28.31	0.0362
KH	1080278	0.338
Ours	1627218	0.00342

表 2 Pokec におけるそれぞれの推定アルゴリズムの収束値

NRMSE が高くなっている.また、pが 0.0 から 0.2 まででは **Ours** の NRMSE は **KH** より高いが、pが 0.2 から 0.3 まで では **Ours** の NRMSE は **KH** より低い.一方で、クエリ数を 5%としたとき NRMSE に顕著な差が見られる.平均次数推定 と同様に、提案アルゴリズムの NRMSE はpの増加に対して ほとんど増加していない.

図 7 は LiveMocha における, API が隣接ノードのプライバ シ情報を提供しないモデルにおける結果を示す. この場合も, 提案アルゴリズムは既存アルゴリズムと比較して p が増加し ても NRMSE はほとんど増加していない. Douban, gowalla, academia でも同様の結果が得られた.

6.2 Pokec グラフにおける提案アルゴリズムの推定精度

実際にプライベートなノードを含むソーシャルネットワーク である Pokec グラフにおいて,提案アルゴリズムの推定精度を 評価する. Pokec グラフは Takac らがプライベートなノードを 含む当時の全グラフデータを取得し,各ノードのプライバシ情 報も付与されている [22,24]. プライベートなノードは 552525 ノード (約 33.8%) 存在する.

Pokec グラフに対して,その最大クラスタからランダムに選 んだ初期ノードとクエリ数を与えてランダムウォークを実行 し,それぞれの推定値を計算する.このシミュレーションを独 立に 1000 回行い,それぞれのアルゴリズムの NRMSE を比較 する.クエリ数はノード数の 0.5%から 5%まで 0.25%ごとに変 化させた.

図8は、APIが隣接ノードのプライバシ情報を提供するモデ ルにおける, Pokec の平均次数とノード数推定における既存ア ルゴリズム (Smooth,KH) と提案アルゴリズム (Ours) のク エリ数に対する NRMSE を示す. 既存アルゴリズムの NRMSE は、クエリ数が増加してもほとんど減少していない、一方で提 案アルゴリズムはクエリ数が増加するにつれて NRMSE が減少 している. 表2は定理 3,4,5,6を用いてそれぞれの推定アルゴ リズムの収束値と真値との相対誤差を計算したものである。提 案アルゴリズムの収束値の相対誤差は既存アルゴリズムと比較 して非常に小さいことがわかる。いずれの推定においても、提 案アルゴリズムによる推定値は既存アルゴリズムより総じて低 い NRMSE を示している。1%未満のクエリ数でのノード数推 定では提案アルゴリズムの推定値の誤差が大きくなっており, この改善は今後の課題である。図9はAPIが隣接ノードのプ ライバシ情報を提供しないモデルにおける結果を示す。この場 合も、提案アルゴリズムは総じて低い NRMSE を示している。







図 4 *p* に対する平均次数の推定値の NRMSE (API が隣接ノードのプライバシ情報を提供する モデル, クエリ数: 0.05*n*)



図 5 *p* に対するノード数の推定値の NRMSE (API が隣接ノードのプライバシ情報を提供する モデル, クエリ数: 0.01*n*)



図 6 p に対するノード数の推定値の NRMSE (API が隣接ノードのプライバシ情報を提供する モデル, クエリ数: 0.05n)



(a) 平均次数推定, クエリ数: 0.01n (b) 平均次数推定, クエリ数: 0.05n (c) ノード数推定, クエリ数: 0.01n (d) ノード数推定, クエリ数: 0.05n





図8 Pokec におけるクエリ数に対する推定値の NRMSE (API が隣 接ノードのプライバシ情報を提供するモデル)



図 9 Pokec におけるクエリ数に対する推定値の NRMSE (API が隣 接ノードのプライバシ情報を提供しないモデル)

7 関連研究

オンラインソーシャルネットワークのプライバシを考慮した研究はいくつかある [25–27]. Bonneau らは Facebook の公開されている隣接関係に基づくグラフ統計量の近似を議論し

た [25]. 各ユーザの友人のうちわずか 8 ユーザが公開されてい るだけで、その部分グラフ上のノードの次数や媒介中心性、コ ミュニティ検出などの統計量を十分に近似できることを実験的 に示した.しかし、ランダムウォークを用いて統計量を推定す る本研究とは異なる. Ye らは OSN のプライベートなノードが クローリングに与える影響を言及している [26]. 実際の OSNs のデータセットに対して 10 万ノードをランダムにプライベー トなノードとしてクローリングで取得できるノード数とエッジ 数の減少率を調べており、大きな減少は見られないとしている. しかし、この研究は統計量推定におけるプライベートなノー ドの影響は論じていない. Chierichetti らはソーシャルネット ワークにおける public-private model を提案した [27]. このモ デルでは全体のグラフGをパブリックなグラフと呼び、パブ リックなグラフに属する各ノード u はプライベートなグラフ G_{u} を持つ. パブリックなグラフはどのノードからでも見える グラフであり、プライベートなグラフ *G*_u はノード *u* のみが見 えるグラフである. 各ノードuに対して $G \cup G_u$ の統計量を効 率よく計算するアルゴリズムを提案している。しかし、この研 究はグラフの全データにアクセスできることを前提にしている ため本研究とは異なる.

8 ま と め

本研究ではプライベートなノードを含むソーシャルネット ワークに対する、ランダムウォークを用いた統計量推定アルゴ リズムを議論した.まず、Gjoka らのパブリックなノードのみ を遷移するランダムウォークの定常分布を導出し、各サンプル ノードはそのパブリックな隣接ノード数に比例した分布でサン プルされることを示した.平均次数とノード数に対する既存の 推定アルゴリズムは、プライベートなノードの割合が増加する につれて推定値の誤差が大きくなる問題点があった.我々は、 既存アルゴリズムのサンプルノードへの重みを変更すること で、プライベートなノードの割合に対して頑健な推定を可能に した.本研究では、OSN の API が提供する隣接ノード情報の 2つのモデルを考え、それぞれのモデルにおける実験で、提案 アルゴリズムがプライベートなノードの割合にほぼ依存せず、 高精度に平均次数とノード数を推定できることを示した.今後 は,統計量の推定精度の向上だけでなく,プライベートなノー ドの割合に対して頑健な推定を達成するアルゴリズムを設計す るべきである.

謝 辞

本研究の一部は、国立研究開発法人新エネルギー・産業技術 総合開発機構(NEDO)の委託業務として行われました.本研 究は JSPS 科研費 16K12406 の助成を受けたものです.

文 献

- Yong-Yeol Ahn, Seungyeop Han, Haewoon Kwak, Sue Moon, and Hawoong Jeong. Analysis of topological characteristics of huge online social networking services. In *Proceedings of the 16th international conference on World Wide Web*, pp. 835–844. ACM, 2007.
- [2] Alan Mislove, Massimiliano Marcon, Krishna P Gummadi, Peter Druschel, and Bobby Bhattacharjee. Measurement and analysis of online social networks. In Proceedings of the 7th ACM SIGCOMM conference on Internet measurement, pp. 29–42. ACM, 2007.
- [3] Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon. What is twitter, a social network or a news media? In Proceedings of the 19th international conference on World wide web, pp. 591–600. AcM, 2010.
- [4] Minas Gjoka, Maciej Kurant, Carter T Butts, and Athina Markopoulou. Practical recommendations on crawling online social networks. *IEEE Journal on Selected Areas in Communications*, Vol. 29, No. 9, pp. 1872–1892, 2011.
- [5] Flavio Chiericetti, Anirban Dasgupta, Ravi Kumar, Silvio Lattanzi, and Tamás Sarlós. On sampling nodes in a network. In *Proceedings of the 25th International Conference on World Wide Web*, pp. 471–481. International World Wide Web Conferences Steering Committee, 2016.
- [6] Rong-Hua Li, Jeffrey Xu Yu, Lu Qin, Rui Mao, and Tan Jin. On random walk based graph sampling. In 2015 IEEE 31st International Conference on Data Engineering (ICDE), pp. 927–938. IEEE, 2015.
- [7] Liran Katzir and Stephen J Hardiman. Estimating clustering coefficients and size of social networks via random walk. *ACM Transactions on the Web (TWEB)*, Vol. 9, No. 4, p. 19, 2015.
- [8] Liran Katzir, Edo Liberty, Oren Somekh, and Ioana A Cosma. Estimating sizes of social networks via biased sampling. *Internet Mathematics*, Vol. 10, No. 3-4, pp. 335–359, 2014.
- [9] Anirban Dasgupta, Ravi Kumar, and Tamas Sarlos. On estimating the average degree. In Proceedings of the 23rd international conference on World wide web, pp. 795–806. ACM, 2014.
- [10] Bruno Ribeiro and Don Towsley. On the estimation accuracy of degree distributions from graph sampling. In *Decision and Control (CDC), 2012 IEEE 51st Annual Conference on*, pp. 5240–5247. IEEE, 2012.
- [11] Xiaowei Chen, Yongkun Li, Pinghui Wang, and John Lui. A general framework for estimating graphlet statistics via random walk. *Proceedings of the VLDB Endowment*, Vol. 10, No. 3, pp. 253–264, 2016.
- [12] Kenta Iwasaki and Kazuyuki Shudo. Estimating the clustering coefficient of a social network by a non-backtracking random walk. In *Big Data and Smart Computing (Big-Comp), 2018 IEEE International Conference on*, pp. 114– 118. IEEE, 2018.

- [13] Toshiki Matsumura, Kenta Iwasaki, and Kazuyuki Shudo. Average path length estimation of social networks by random walk. In *Big Data and Smart Computing (Big-Comp), 2018 IEEE International Conference on*, pp. 611– 614. IEEE, 2018.
- [14] Kazuki Nakajima, Kenta Iwasaki, Toshiki Matsumura, and Kazuyuki Shudo. Estimating top-k betweenness centrality nodes in online social networks. In *Proceedings of the IEEE ISPA-IUCC-BDCloud-SocialCom-SustainCom 2018*, pp. 1128–1135. IEEE, 2018.
- [15] Chul-Ho Lee, Xin Xu, and Do Young Eun. Beyond random walk and metropolis-hastings samplers: why you should not backtrack for unbiased graph sampling. In ACM SIGMET-RICS Performance evaluation review, Vol. 40, pp. 319–330. ACM, 2012.
- [16] Ratan Dey, Zubin Jelveh, and Keith Ross. Facebook users have become much more private: A large-scale study. In Pervasive Computing and Communications Workshops (PERCOM Workshops), 2012 IEEE International Conference on, pp. 346–352. IEEE, 2012.
- [17] Kenta Iwasaki and Kazuyuki Shudo. Comparing graph sampling methods based on the number of queries. In Proceedings of the IEEE ISPA-IUCC-BDCloud-SocialCom-SustainCom 2018, pp. 1136–1143. IEEE, 2018.
- [18] David A Levin and Yuval Peres. Markov chains and mixing times, Vol. 107. American Mathematical Soc., 2017.
- [19] Réka Albert, Hawoong Jeong, and Albert-László Barabási. Error and attack tolerance of complex networks. *nature*, Vol. 406, No. 6794, p. 378, 2000.
- [20] Reuven Cohen, Keren Erez, Daniel Ben-Avraham, and Shlomo Havlin. Resilience of the internet to random breakdowns. *Physical review letters*, Vol. 85, No. 21, p. 4626, 2000.
- [21] Ryan A. Rossi and Nesreen K. Ahmed. The network data repository with interactive graph analytics and visualization. In *Proceedings of the Twenty-Ninth AAAI Conference* on Artificial Intelligence, 2015.
- [22] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. http://snap.stanford. edu/data, June 2014.
- [23] Pinghui Wang, John Lui, Bruno Ribeiro, Don Towsley, Junzhou Zhao, and Xiaohong Guan. Efficiently estimating motif statistics of large networks. ACM Transactions on Knowledge Discovery from Data (TKDD), Vol. 9, No. 2, p. 8, 2014.
- [24] Lubos Takac and Michal Zabovsky. Data analysis in public social networks. In International Scientific Conference and International Workshop Present Day Trends of Innovations, 2012.
- [25] Joseph Bonneau, Jonathan Anderson, Ross Anderson, and Frank Stajano. Eight friends are enough: social graph approximation via public listings. In *Proceedings of the Second ACM EuroSys Workshop on Social Network Systems*, pp. 13–18. ACM, 2009.
- [26] Shaozhi Ye, Juan Lang, and Felix Wu. Crawling online social graphs. In Web Conference (APWEB), 2010 12th International Asia-Pacific, pp. 236–242. IEEE, 2010.
- [27] Flavio Chierichetti, Alessandro Epasto, Ravi Kumar, Silvio Lattanzi, and Vahab Mirrokni. Efficient algorithms for public-private social networks. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 139–148. ACM, 2015.