

Paper Reading

Quick Detection of High-degree Entities in Large Directed Networks

K. Avrachenkov
Inria
k.avrachenkov@inria.fr

N. Litvak
University of Twente
n.litvak@utwente.nl

L. Ostroumova Prokhorenkova
Yandex
ostroumova-la@yandex.ru

E. Suyargulova
Yandex
siyargul@yandex.ua

(from ICDM '14)

Jan 2015

Speaker: Kazuhiro Inaba

内容

- Social graph から
- Top-k high in-degree nodes を
- **Sublinear time で**
- 見つける

より具体的な問題設定 (1)

- Twitter の人気ユーザを知りたい
- しかし Twitter API には回数制限が！
- **頂点数 \approx 1.5G より遙かに少ない API 呼出で、人気ユーザを割り出したい**

より具体的な問題設定 (2)

- できること

- ユーザを uniformly random に選ぶ
- API呼出1回で、ユーザ1人の入次数を得る
- API呼出1回で、ユーザ1人の出辺を5000本まで得る
(※仮定: 5000 ≧ ほとんどのユーザの出次数)

提案手法

(n_1+n_2 回のAPI呼出で n_2 個の<頂点,次数>対を返す)

```
for  $w$  in  $W$  do
```

```
   $S[w] \leftarrow 0;$ 
```

```
for  $i \leftarrow 1$  to  $n_1$  do
```

```
   $v \leftarrow \text{random}(N);$ 
```

```
  foreach  $w$  in  $\text{OutNeighbors}(v) \subset W$  do
```

```
     $S[w] \leftarrow S[w] + 1;$ 
```

```
 $w_1, \dots, w_{n_2} \leftarrow \text{Top}_{n_2}(S)$  //  $S[w_1], \dots, S[w_{n_2}]$  are  
the top  $n_2$  maximum values in  $S$ ;
```

```
for  $i \leftarrow 1$  to  $n_2$  do
```

```
   $d_i \leftarrow \text{InDegree}(w_i);$ 
```

n_1 個の頂点を
ランダムに選ぶ

そこから多く指されている
 n_2 個の頂点について

全体からの入次数を取得。
(Top-K が欲しければここからK個取る)

実験

- 実際に Twitter に対して
 $n_1 + n_2 = 1000$ 回APIを呼んで実験
- “Ground Truth” との近さで評価
 - **Fraction:** 真の Top-K のうち何割を見つけられたか
 - **First-Error Index:** 見つけた n_2 個から漏れた最上位

実験: “Ground Truth”

- Twitter社は公開していない
- twittercounter.com という第三者サービスを参照しようとした



The screenshot shows a web browser window displaying the Twitter Counter website. The browser's address bar shows the URL "twittercounter.com/pages/100". The website header includes the "TWITTER COUNTER" logo and a "Sign in with Twitter" button. The main content area is a table listing the top 100 most followed Twitter users. The table has two columns: "Twitter users" and "Followers". The top three entries are:

	Twitter users	Followers
1	 KATY PERRY @katyperry	62,765,047
2	 Justin Bieber @justinbieber	58,821,179
3	 Barack Obama @BarackObama	52,509,744

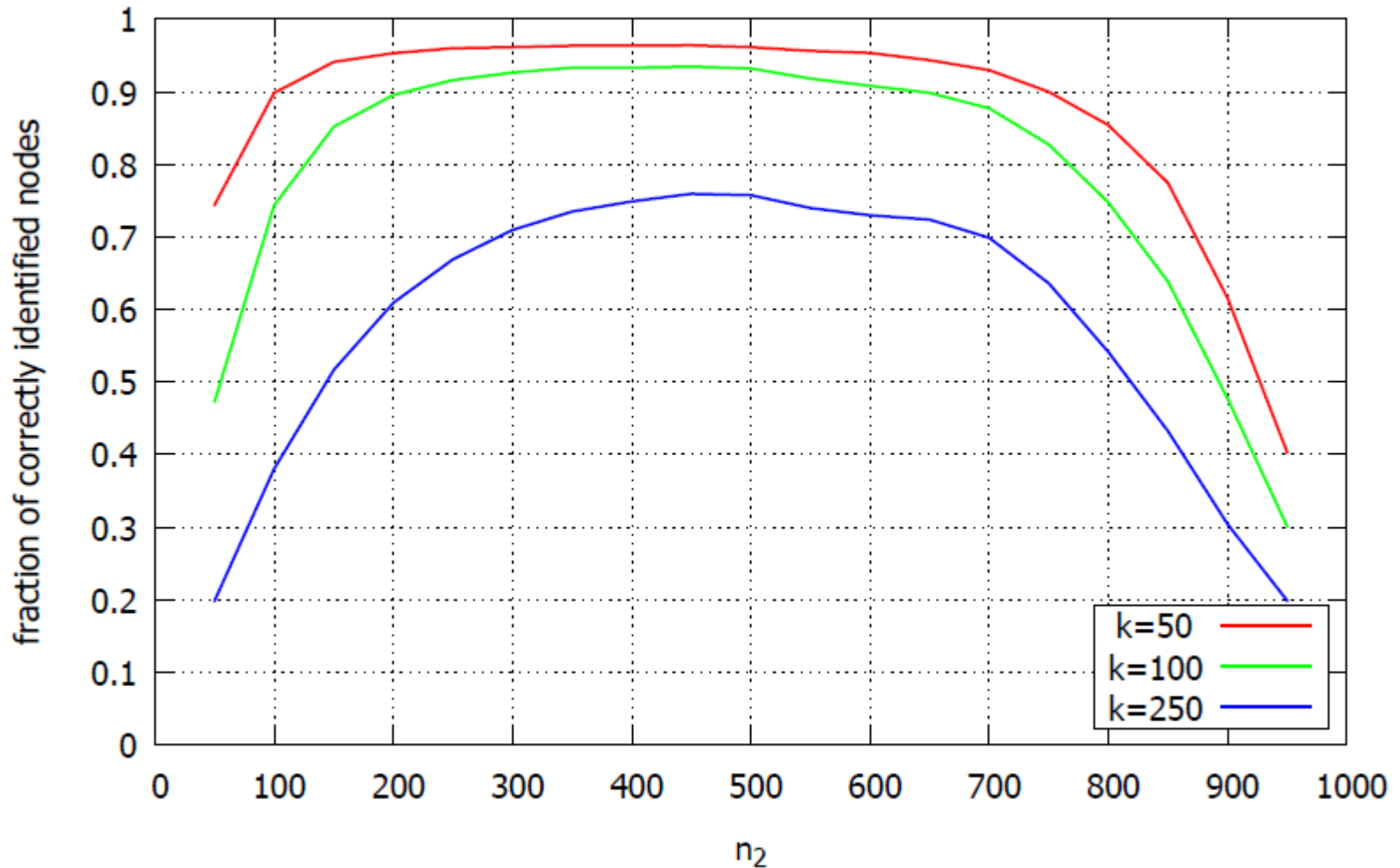
“Ground Truth”

- Twitter社は公開していない
- twittercounter.com という第三者サービスを参照しようとした

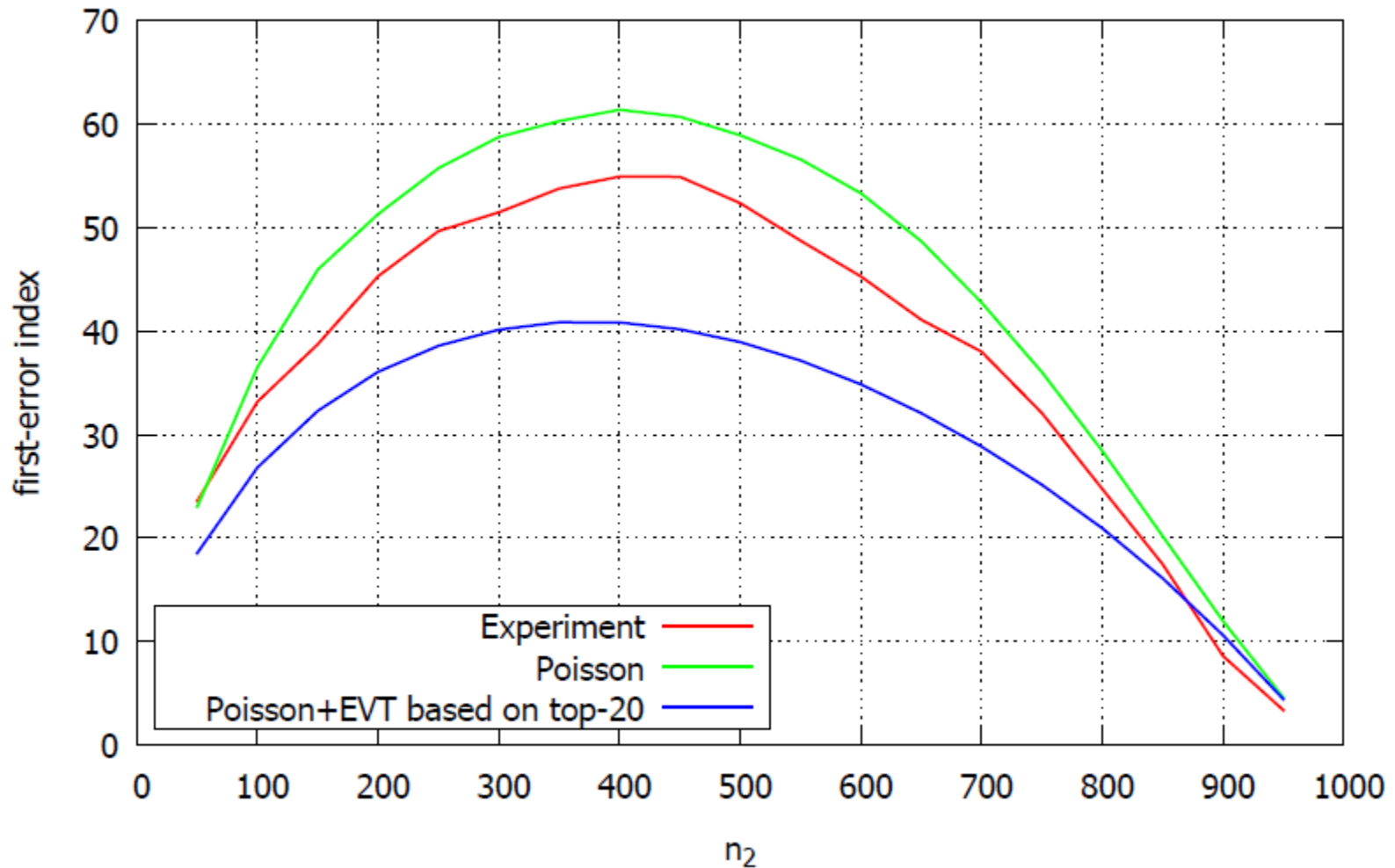
!!!しかし!!!

- 提案手法 [$n_1=n_2=20,000$] の方が精度がよかった
→ [$n_1=n_2=500,000$] をGround Truthとして使用

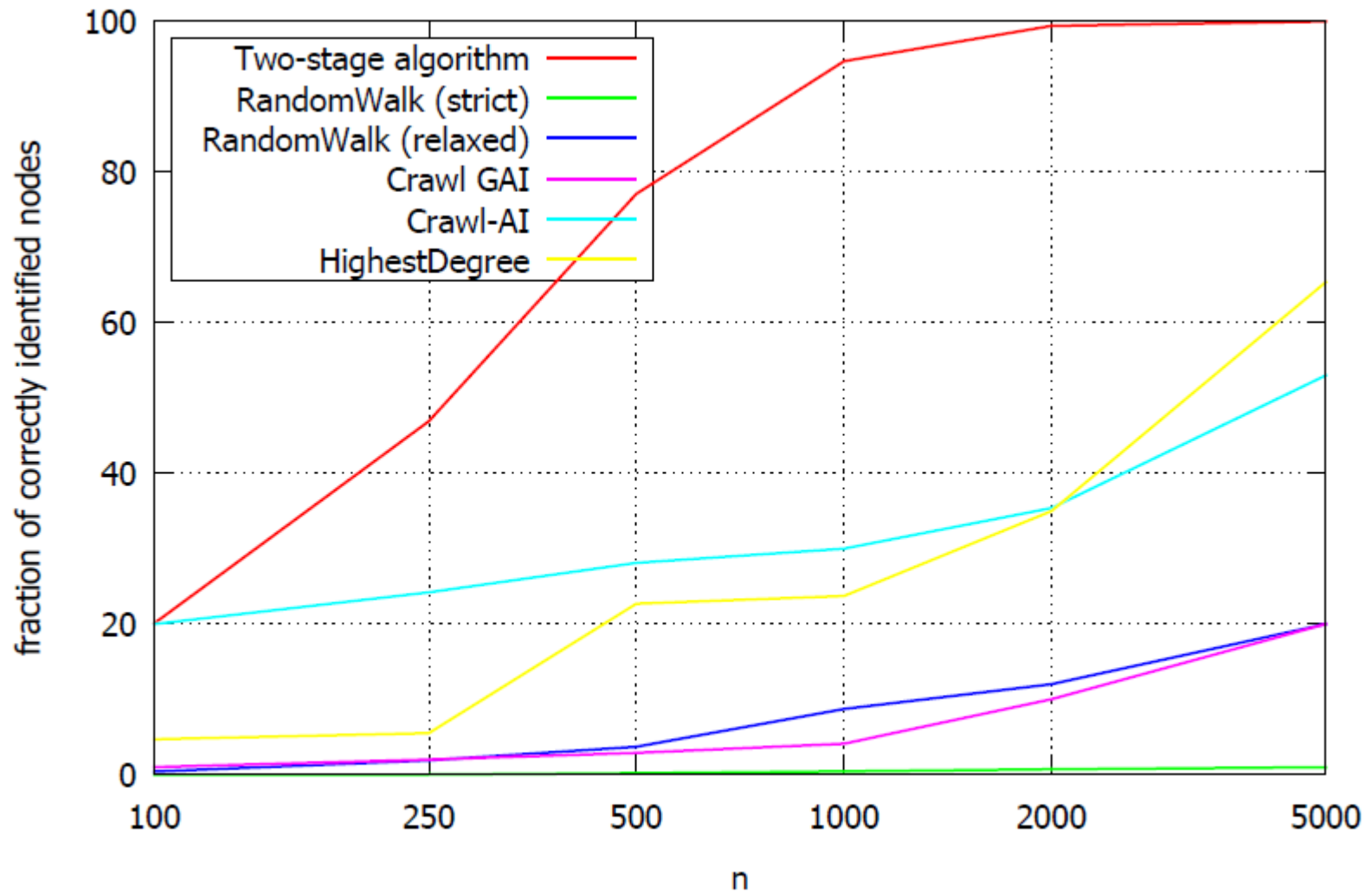
Twitterでの実験結果 (Fraction)



Twitterでの実験結果 (First-Error)



他の手法との比較 (Top-100; Fraction)



その他の応用

- 提案手法は「フォローする側」と「される側」が分かれている場合でもそのまま適用可能
 - 例：SNSで人気のある“グループ”は？
 - ロシアのSNS VK.com (～200M users) で実験。
($n_1=700$, $n_2=300$) で top-100 の73.2%を発見

解析

- n_1 と n_2 の最適なバランスは？
- $n = n_1 + n_2$ をどのくらい増やせば十分な精度が得られる？

n_1 と n_2 の最適なバランスは？

Proposition 1. *It is optimal to choose n_1 such that $n = O(n_1)$.*

= “ n_1 に比べてあまり n_2 を大きくしても意味がない”

証明: そもそもfirst stageで

平均次数 $\times n_1$ 頂点くらいしか選ばれないので。

どのくらい増やせば十分な精度？

Proposition 2. For large enough n_1 , the inequality

$$Z_k(n_1) := \sqrt{\frac{n_1}{N}} \frac{F_k - F_{n_2}}{\sqrt{F_k + F_{n_2}}} \geq z_{1-\varepsilon}, \quad (14)$$

where $z_{1-\varepsilon}$ is the $(1 - \varepsilon)$ -quantile of a standard normal distribution, guarantees that the mean fraction of top- k nodes in W identified by Algorithm 1 is at least $1 - \varepsilon$.

“この不等式が成り立つように n_1, n_2 を取れば 確率 $1-\varepsilon$ で 第 k 位のノードが出力に含まれる”

where

$P_k(n_1) :=$ 真の第 k 位ノードが選ばれる確率

$z_{1-\varepsilon} :=$ 平均0分散1の正規分布の $(1-\varepsilon)$ -quantile

$\gamma :=$ Scale free 性を仮定。べき分布の指数(の逆数)

$F_k :=$ InDeg(k) $\propto k^{-\gamma}$

例: $n_2=300, k=100, \varepsilon=10\%$, Twitter に対して計算すると $n_1 \geq 1300$

$$P_k(n_1, n_2)$$

$$= P[S_k \geq S_{n_2}^+]$$

$$\cong P[S_k \geq S_{n_2}]$$

$$\cong P[\text{二項分布}(n_1, F_k/N) \geq \text{二項分布}(n_1, F_{n_2}/N)]$$

正規分布で近似

$$\cong P[\text{正規分布}(n_1 F_k/N, n_1 F_k/N(1-F_k/N)) \geq \text{正規分布}(\dots)]$$

$$\cong P[\text{正規分布}(n_1 F_k/N, n_1 F_k/N) \geq \text{正規分布}(\dots)]$$

$$= P[\text{正規分布}(n_1(F_k - F_{n_2})/N, n_1(F_k + F_{n_2})/N) \geq 0]$$

正規分布の差

$$= P[\text{正規分布}(\sqrt{n_1/N} \cdot (F_k - F_{n_2})/\sqrt{F_k + F_{n_2}}, 1) \geq 0]$$

分散を1に

P_k := 第k位のノードが出力に入る確率

S_k := 第k位のノードの得票数。二項分布になる

F_k := 第k位のノード入次数

Proposition 2. For large enough n_1 , the inequality

$$Z_k(n_1) := \sqrt{\frac{n_1}{N}} \frac{F_k - F_{n_2}}{\sqrt{F_k + F_{n_2}}} \geq z_{1-\varepsilon}, \quad (14)$$

where $z_{1-\varepsilon}$ is the $(1 - \varepsilon)$ -quantile of a standard normal distribution, guarantees that the mean fraction of top- k nodes in W identified by Algorithm 1 is at least $1 - \varepsilon$.

n_1 と n_2 の最適なバランスは？

Theorem 1. Assume that $k = o(n)$ as $n \rightarrow \infty$, then the maximizer of the probability $P_k(n - n_2)$ is

$$n_2 = (3\gamma k^\gamma n)^{\frac{1}{\gamma+1}} (1 + o(1)),$$

with γ as in (1).

証明: さきほど導出した P_k の近似式を元に頑張っって計算。

どのくらい増やせば十分な精度？

Theorem 2. *Let the in-degrees of the entities in W be independent realizations of a regularly varying distribution G with exponent $1/\gamma$ as defined in (1), and $F_1 \geq F_2 \geq \dots \geq F_M$ be their order statistics. Then for any fixed $\varepsilon, \delta > 0$, Algorithm 1 finds the fraction $1 - \varepsilon$ of top- k nodes with probability $1 - \delta$ in*

$$n = O(N/a(M))$$

API requests, as $M, N \rightarrow \infty$, where $a(M) = l(M)M^\gamma$ and $l(\cdot)$ is some slowly varying function.

In the case $M = N$, as in our experiments on Twitter, Theorem 2 states that the complexity of the algorithm is roughly of the order $N^{1-\gamma}$, which is much smaller than linear

- 証明: 計算。

まとめ

- 次数の高いノードを少ないクエリ回数で発見
- 非常にシンプルなアルゴリズム
- Extreme Value Theory を用いた解析

```
for  $w$  in  $W$  do
   $S[w] \leftarrow 0$ ;
for  $i \leftarrow 1$  to  $n_1$  do
   $v \leftarrow \text{random}(N)$ ;
  foreach  $w$  in  $\text{OutNeighbors}(v) \subset W$  do
     $S[w] \leftarrow S[w] + 1$ ;
 $w_1, \dots, w_{n_2} \leftarrow \text{Top\_}n_2(S)$  //  $S[w_1], \dots, S[w_{n_2}]$  are
the top  $n_2$  maximum values in  $S$ ;
for  $i \leftarrow 1$  to  $n_2$  do
   $d_i \leftarrow \text{InDegree}(w_i)$ ;
```

