

# Efficient **PageRank** Tracking in **Evolving** Networks

*KDD'15*

21<sup>st</sup> ACM SIGKDD Conference on  
Knowledge Discovery and Data Mining

大坂 直人 (東京大学 / プロジェクト RA)

前原 貴憲 (静岡大学)

河原林 健一 (NII)

はじめに

# PageRankとPersonalized PageRank

## ■ PageRank [Brin-Page.'98]

Webページの**重要度**の指標

リンク構造だけから決まる

一般化

## ■ Personalized PageRank

[Jeh-Widom.'03]

バイアス付き ⇨ **関連度**

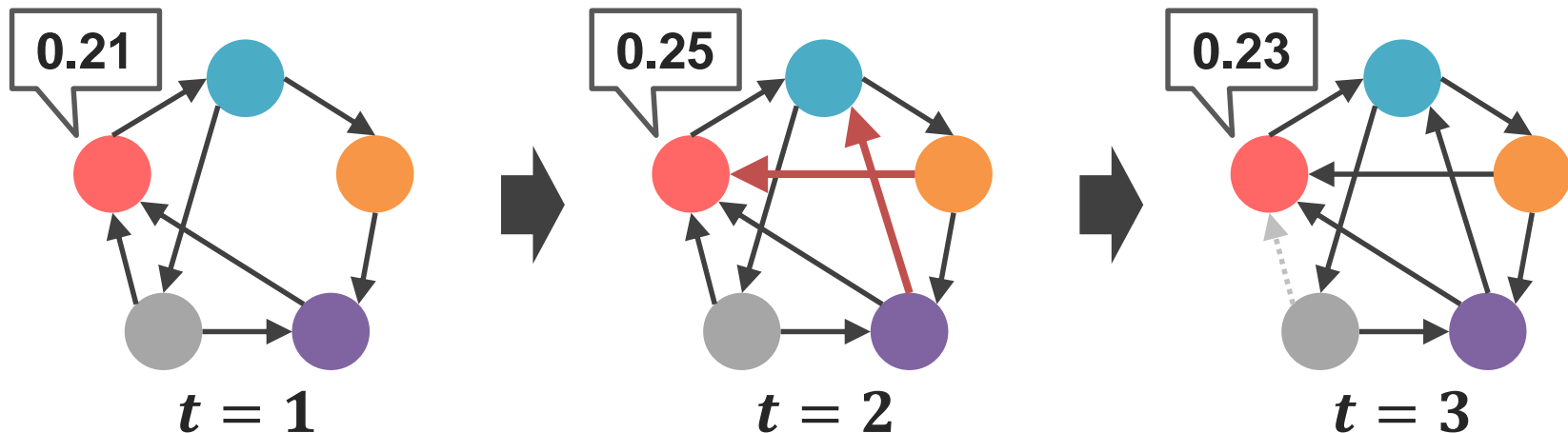


去年の感謝祭…静的グラフ上の高速計算

[Maehara-Akiba-Iwata-Kawarabayashi. PVLDB'14]

はじめに

# Evolving Networks … 動的グラフ



- World Wide Web

新しいページ・リンク      60万ページ/秒

<http://www.internetlivestats.com/>

- ソーシャルネットワーク

新しい友人関係

- マイクロブログ

ユーザ同士のやりとり      5000ツイート/秒

<http://www.technologyreview.com/graphiti/522376/the-many-tongues-of-twitter/>

グラフ全体を見ずに更新したい

# はじめに 関連度としての応用

## スパム検知

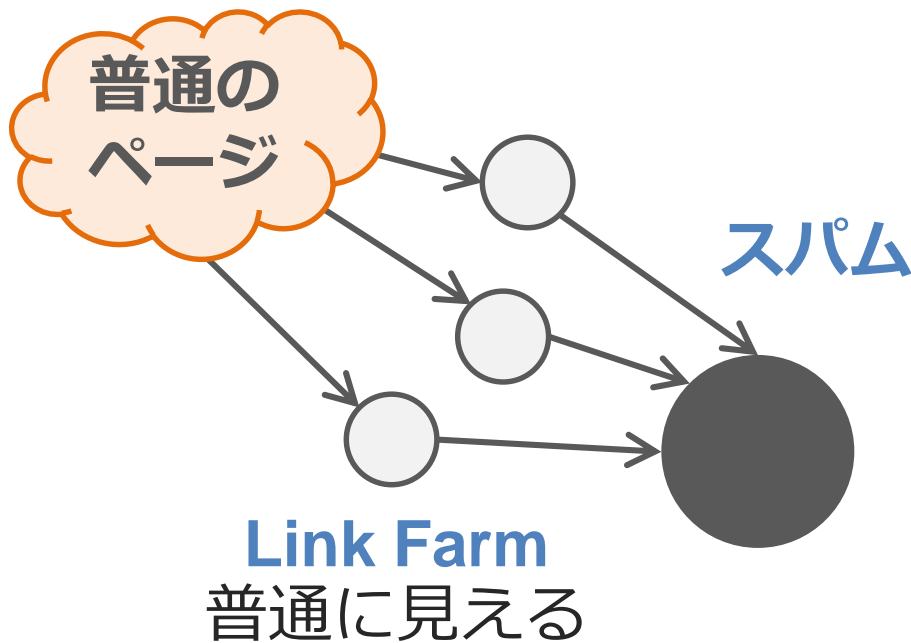
スコアの時間変化を利用

[Chung-Toyoda-Kitsuregawa. '09, '10]

## ユーザ推薦

友達の候補生成

[Gupta-Goel-Lin-Sharma-Wang-Zadeh. WWW'13]



はじめに

# Personalized PageRank追跡の既存手法

	$m$ 辺の無作為挿入の時間	スケーラビリティ 0.1秒以下 / 辺 誤差約 $10^{-9}$
Aggregation/Disaggregation [Chien et al. '04]	$O(m S  \log 1/\epsilon)$	68M 辺
Monte-Carlo [Bahmani et al. '10]	$O(m + \log m / \epsilon^2)$	68M 辺
Power method ナイーブな手法	$O(m^2 \log 1/\epsilon)$	11M 辺

# 本研究の貢献

成長するグラフにおける

**Personalized PageRank 追跡**のための

**高速 & 高精度** な手法を提案

	<b><math>m</math>辺の無作為挿入の時間</b>	<b>スケーラビリティ</b> 0.1秒以下 / 辺 誤差約 $10^{-9}$
<b>提案手法</b>	平均 $\downarrow$ 最大出次数 <b><math>O(m + \Delta \log m / \epsilon)</math></b>	<b>3,700M 辺</b>
Aggregation/Disaggregation [Chien et al. '04]	$O(m S  \log 1/\epsilon)$	68M 辺
Monte-Carlo [Bahmani et al. '10]	$O(m + \log m / \epsilon^2)$	68M 辺
Power method ナイーブな手法	$O(m^2 \log 1/\epsilon)$	11M 辺

# 予備知識

# Personalized PageRank の定義

[Brin-Page. Comput. Networks ISDN Syst.'98] [Jeh-Widom. WWW'03]

## ■ 線形方程式

次の解

$$x = \alpha P x + (1 - \alpha) b$$

バイアス

確率遷移行列

非ジャンプ確率 = 0.85

Random walk  
による解釈



定常分布

## ■ Random walkによるWeb閲覧のモデル化

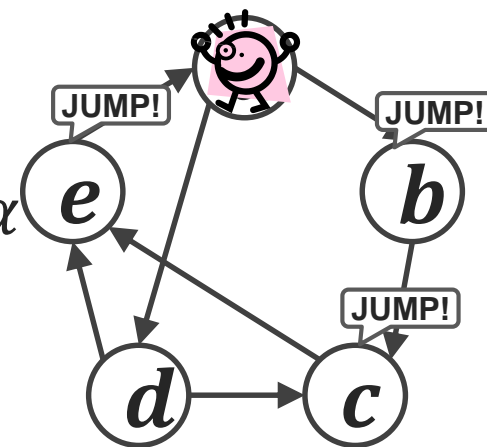
無作為に隣接頂点に移動

w.p.  $\alpha$

無作為に任意頂点にジャンプ  
分布  $b \in \mathbb{R}^n$

w.p.  $1 - \alpha$

$x_v = v$  を訪れる確率





# 静的グラフ上のPageRankの計算

- 線形方程式  $x = \alpha Px + (1 - \alpha)b$  を解く

Power method  $x^{(v)} = \alpha Px^{(v-1)} + (1 - \alpha)b$

- $v$  を訪れる割合  $x_v$  を見積もる

Monte-Carlo シミュレーション

# 動的グラフ上のPageRankの追跡

- Aggregation/disaggregation

[Chien-Dwork-Kumar-Simon-Sivakumar. Internet Math.'04]

変化のあった近傍にPower methodを適用

↑  
大きい☹

- Monte-Carlo ベース

[Bahmani-Chowdhury. VLDB'10]

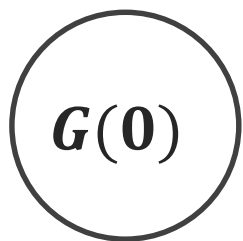
Random walkの軌跡を保持・更新

↑  
沢山必要☹

# 提案手法

# 問題設定

時刻0の入力 :  $G(0)$

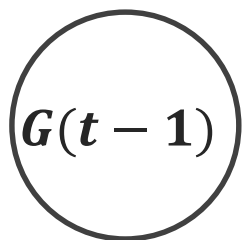


時刻0の問題 :

$G(0)$  の PageRank  $x(0)$  の近似計算

$$\|x(0) - x^*(0)\|_{\infty} < \epsilon$$

⋮

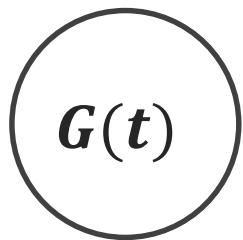


時刻  $t-1$  の問題 :

$G(t-1)$  の PageRank  $x(t-1)$  の近似計算

時刻  $t$  の入力 :

**追加 / 削除** される辺集合



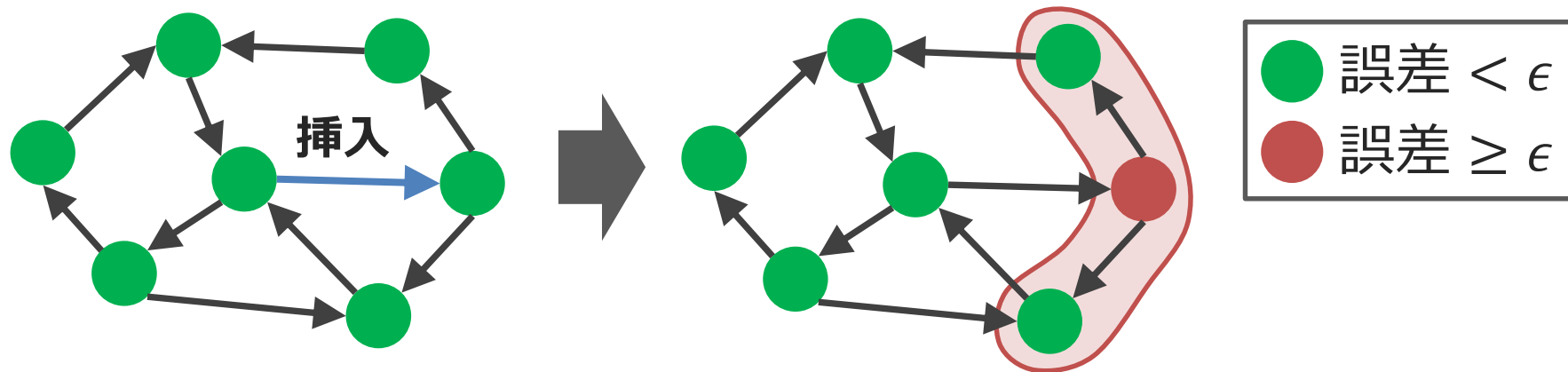
時刻  $t$  の問題 :

$G(t)$  の PageRank  $x(t)$  の近似計算

# そのアイデア

$$x(t) = \alpha P(t)x(t) + (1 - \alpha)b \text{ を解く}$$

1.  $x(t - 1)$  は  $x(t)$  の **良い** 初期近似解
2. 近似解を **局所的** に改善できないか？



- **Gauss-Southwell** 法 を採用 😊 [Southwell. '40, '46]

別名 *Local algorithm*

[Spielman-Teng. SIAM J. Comput.'13]

*Bookmark coloring algorithm*

[Berkhin. Internet Math.'06]

# Gauss-Southwell 法 [Southwell. '40, '46]

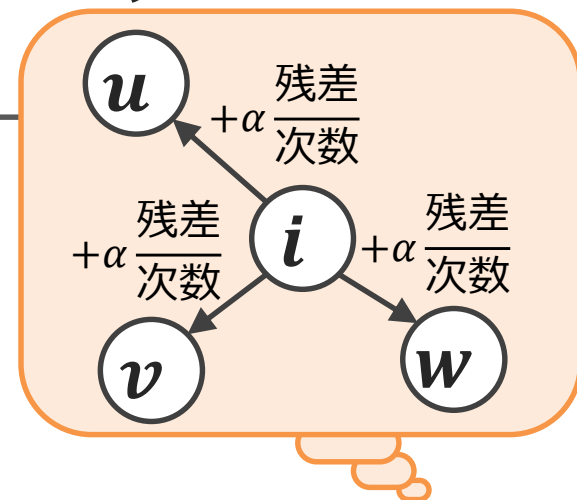
- $\nu^{\text{th}}$  近似解  $x^{(\nu)}$
- $\nu^{\text{th}}$  残差  $r^{(\nu)} = (1 - \alpha)b - (I - \alpha P)x^{(\nu)}$   
できるだけ0に近づける

$\nu = 1, 2, 3, \dots$

$i \leftarrow |r_i^{(\nu-1)}|$  が最大の頂点

**If**  $|r_i^{(\nu-1)}| < \epsilon$  **terminate**

$r_i^{(\nu)} = 0$  となるよう  $x^{(\nu-1)}$  と  $r^{(\nu-1)}$  を局所的に更新



近似保証:  $\|x^* - x^{(\nu)}\|_{\infty} \leq \frac{\epsilon}{1-\alpha}$

# Gauss-Southwell 法 [Southwell. '40, '46]

- $\nu^{\text{th}}$  近似解  $x^{(\nu)}$
- $\nu^{\text{th}}$  残差  $r^{(\nu)} = (1 - \alpha)b - (I - \alpha P)x^{(\nu)}$   
できるだけ0に近づける

$\nu = 1, 2, 3, \dots$

$i \leftarrow |r_i^{(\nu-1)}|$  が最大の頂点

**If**  $|r_i^{(\nu-1)}| < \epsilon$  **terminate**

$$x^{(\nu)} = x^{(\nu-1)} + r_i^{(\nu-1)} e_i$$

$$r^{(\nu)} = r^{(\nu-1)} - r_i^{(\nu-1)} e_i + \alpha r_i^{(\nu-1)} P e_i$$

$$\leq \frac{\|r^{(0)}\|}{(1-\alpha)\epsilon} \quad \square$$

$\square$  は  $\|r^{(\nu-1)}\|_1$  を  $(1 - \alpha)\epsilon$  以上減らす

# その概観

時刻  $t$ :

$$x(t)^{(0)} = x(t - 1)$$

$$r(t)^{(0)} = r(t - 1) + \alpha(P(t) - P(t - 1))x(t - 1)$$

Gauss-Southwell 法を適用



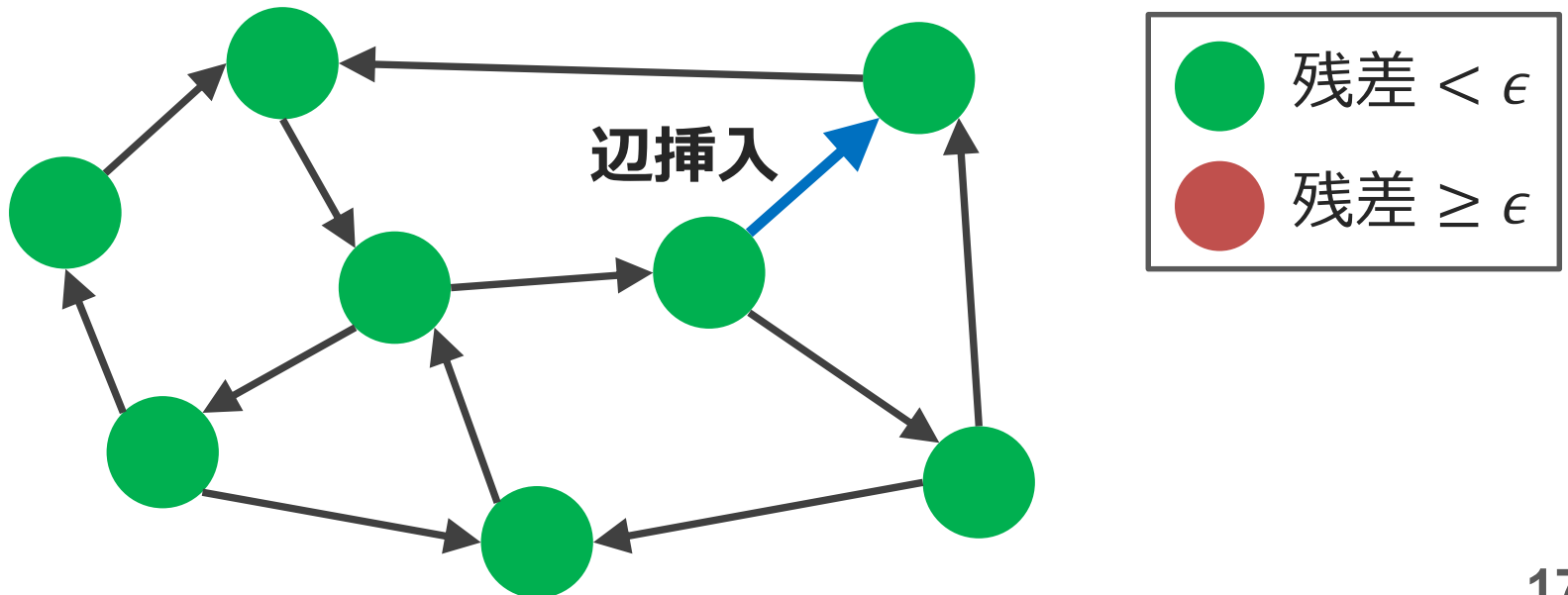
# 提案手法 挙動

時刻  $t$ :

$$x(t)^{(0)} = x(t - 1)$$

$$r(t)^{(0)} = r(t - 1) + \alpha(P(t) - P(t - 1))x(t - 1)$$

Gauss-Southwell 法を適用



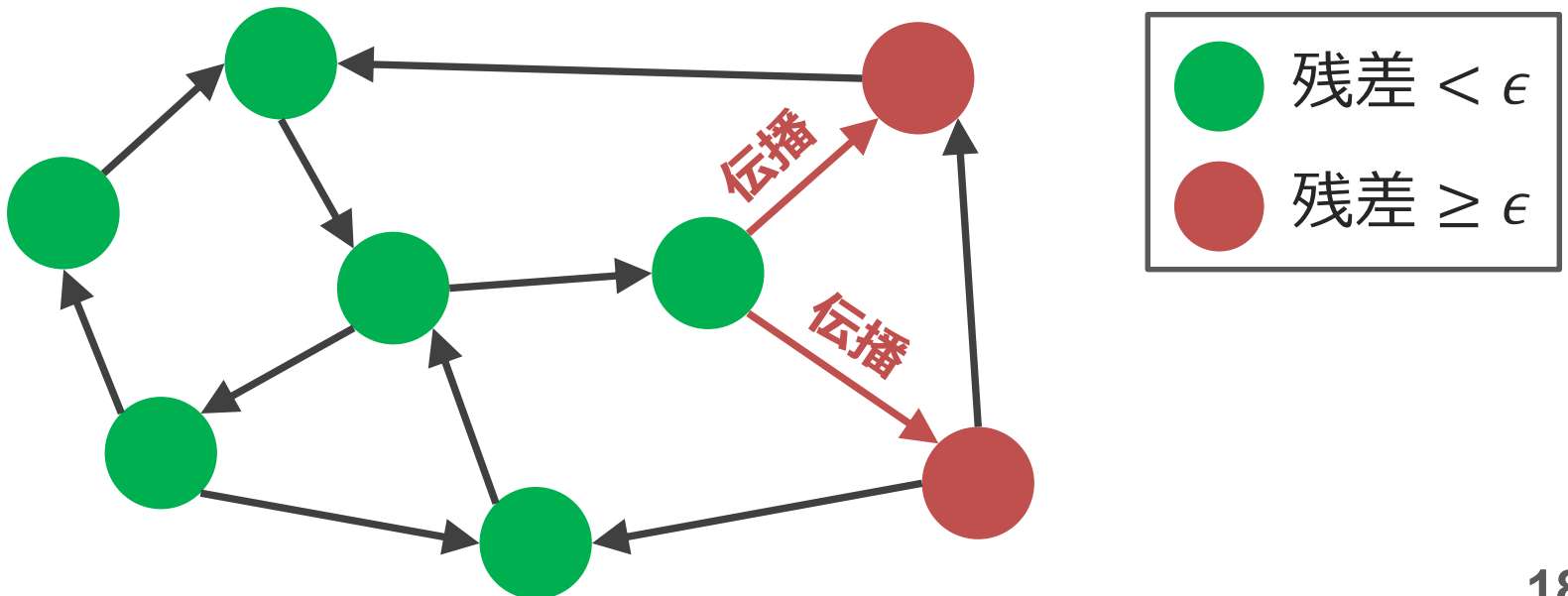
# 提案手法 挙動

時刻  $t$ :

$$x(t)^{(0)} = x(t - 1)$$

$$\rightarrow r(t)^{(0)} = r(t - 1) + \alpha(P(t) - P(t - 1))x(t - 1)$$

Gauss-Southwell 法を適用



# 提案手法 挙動

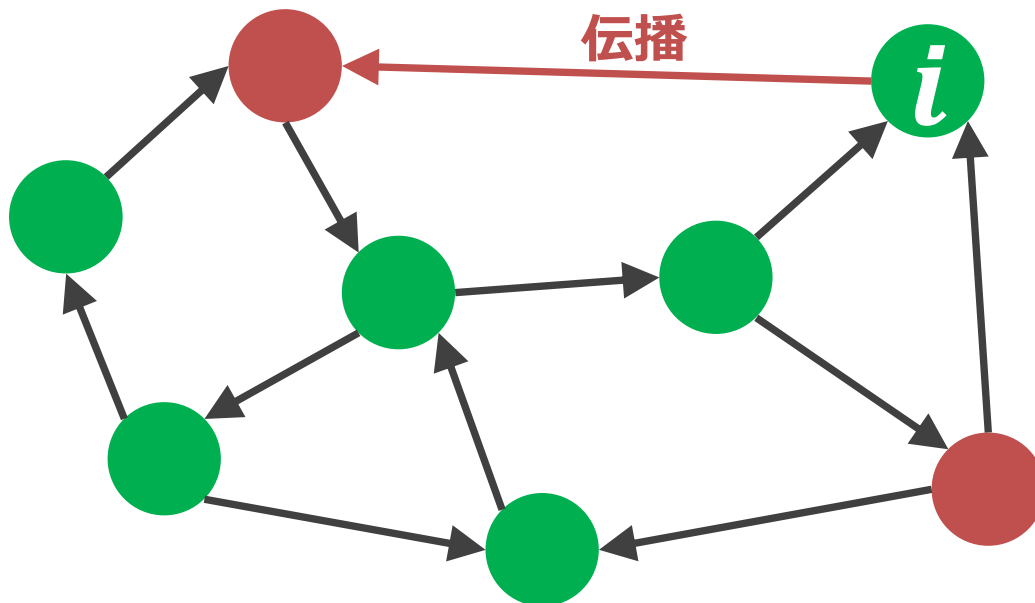
時刻  $t$ :

$$x(t)^{(0)} = x(t - 1)$$

$$r(t)^{(0)} = r(t - 1) + \alpha(P(t) - P(t - 1))x(t - 1)$$

➡ Gauss-Southwell 法を適用

$\nu = 1$



# 提案手法 挙動

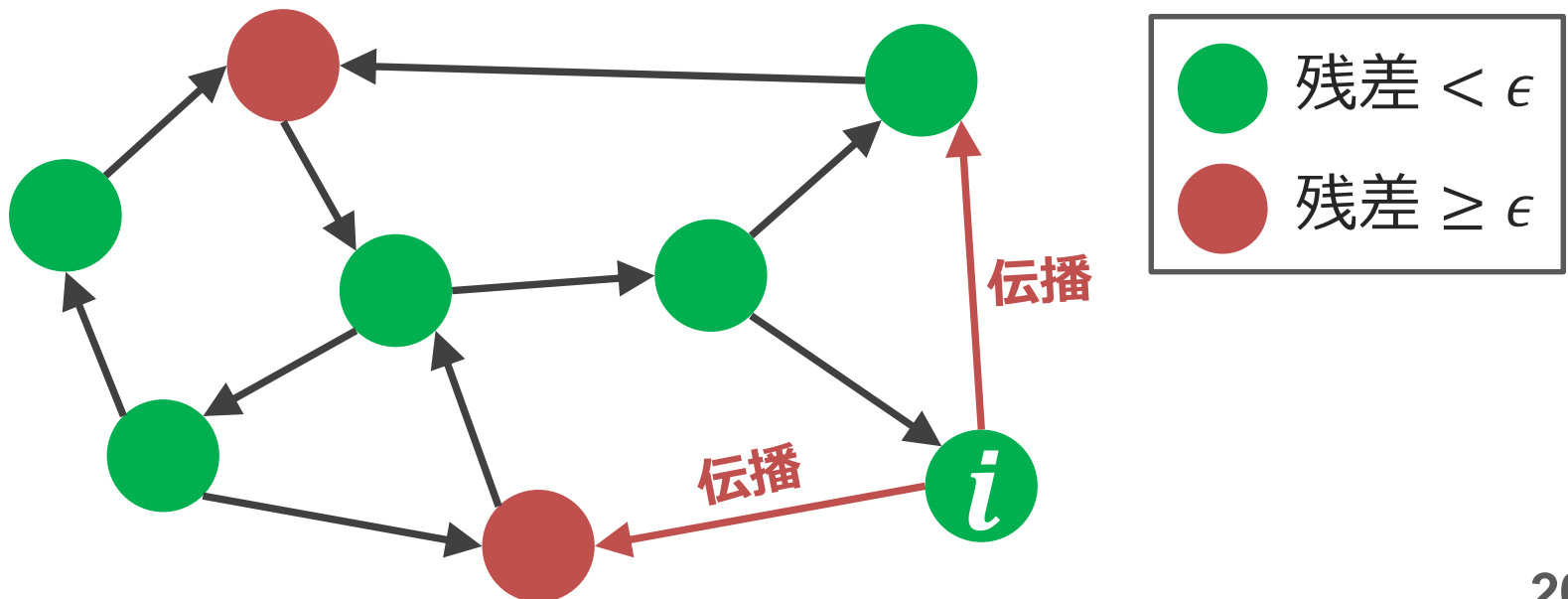
時刻  $t$ :

$$x(t)^{(0)} = x(t - 1)$$

$$r(t)^{(0)} = r(t - 1) + \alpha(P(t) - P(t - 1))x(t - 1)$$

➡ Gauss-Southwell 法を適用

$\nu = 2$



# 性能解析

時刻  $t$ :

$$x(t)^{(0)} = x(t - 1)$$

$$r(t)^{(0)} = r(t - 1) + \alpha(P(t) - P(t - 1))x(t - 1)$$

Gauss-Southwell 法を適用

は**効率的**に計算可

辺  $vw$  の追加／削除は  $O(d_v)$  時間

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 1 \\ 1/2 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1/2 & 0 \\ 1/2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/2 & 0 \end{bmatrix}$$

$P(t - 1)$



$$\begin{bmatrix} 0 & 0 & 0 & 0 & 1 \\ 1/2 & 0 & 0 & 1/3 & 0 \\ 0 & 1 & 0 & 1/3 & 0 \\ 1/2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/3 & 0 \end{bmatrix}$$

$P(t)$

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & -1/6 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1/6 & 0 \end{bmatrix}$$

$P(t) - P(t - 1)$

# 性能解析

時刻  $t$ :

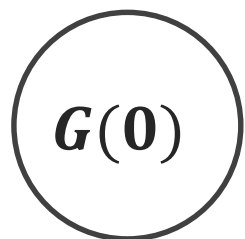
残差の増分  $\|\cdot\|_1$  はどの位？

$$x(t)^{(0)} = x(t-1)$$

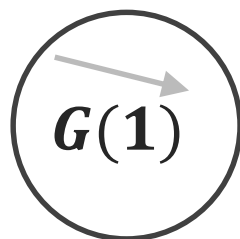
$$r(t)^{(0)} = r(t-1) + \alpha(P(t) - P(t-1))x(t-1)$$

**Gauss-Southwell 法を適用**

- 各時刻で単一辺が**無作為**に挿入



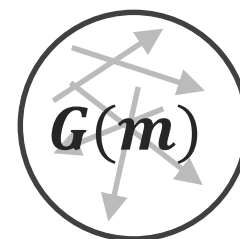
辺集合 =  $\emptyset$



1辺足す



...



$m$ 辺足し終えた

$$\mathbf{E}[\text{時刻 } t \text{ での残差の増分}] \leq 2\alpha/t$$

# 性能解析

## Gauss-Southwell 法を適用

### ■ 定理 1.

任意の変更に対する反復回数 は ならし  $\mathcal{O}(1/\epsilon)$

⇒ ならし時間は  $\mathcal{O}(\Delta/\epsilon)$

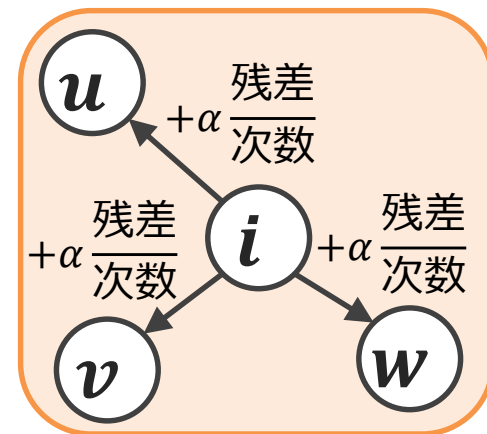
### ■ 定理 2.

$m$  辺を無作為・逐次的に挿入

期待総反復回数は  $\mathcal{O}(\log m / \epsilon)$

⇒ 期待総時間は  $\mathcal{O}(m + \Delta \log m / \epsilon)$

$\Delta =$  最大出次数

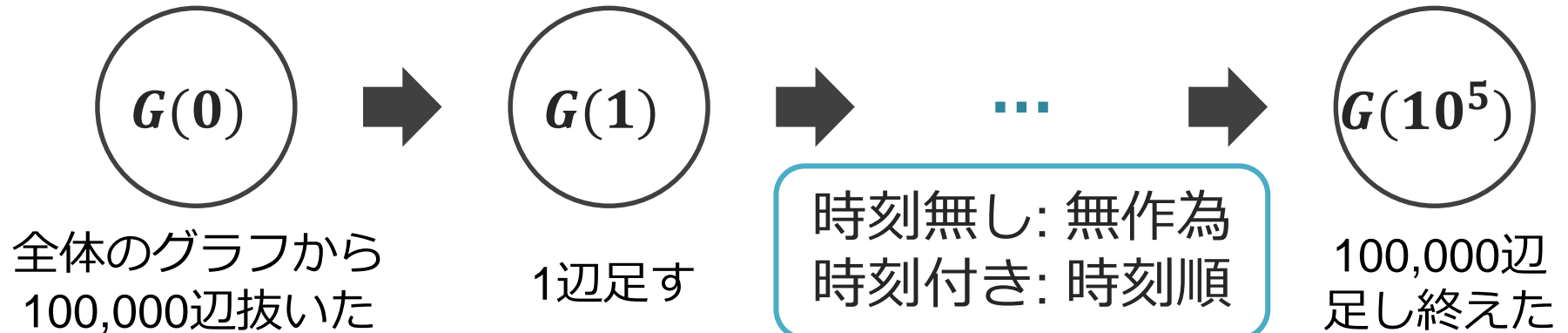


# 実験



# 実験

## 設定: 単一辺挿入の性能評価



### ■ パラメータ設定

- $\alpha = 0.85$
- $b = 100$ 要素が非ゼロのベクトル
- $\epsilon = 10^{-9}$

# 実験

## 性能評価: 単一辺挿入の平均時間 & 反復回数

データセット [出典]	頂点数 $ V $	辺数 $ E $	最大 出次数 $\Delta$	平均 時間	平均 反復回数
wiki-Talk [SNAP]	2M	5M	100,022	589.6 $\mu$ s	2.3
web-Google [SNAP]	1M	5M	3,444	7.2 $\mu$ s	22.6
as-Skitter [SNAP]	2M	11M	35,387	288.4 $\mu$ s	0.8
Flickr <sup>時刻</sup> [KONECT]	2M	33M	26,367	95.3 $\mu$ s	16.2
Wikipedia <sup>時刻</sup> [KONECT]	2M	40M	6,975	76.8 $\mu$ s	46.0
soc-LiveJournal1 [SNAP]	5M	68M	20,292	17.9 $\mu$ s	7.6
twitter-2010 [LAW]	42M	1,500M	2,997,469	29,382.8 $\mu$ s	0.7
uk-2007-05 [LAW]	105M	3,700M	15,402	2.3 $\mu$ s	0.0

[KONECT] The Koblenz Network Collection <http://konect.uni-koblenz.de/networks/>

[LAW] Laboratory for Web Algorithmics <http://law.di.unimi.it/datasets.php>

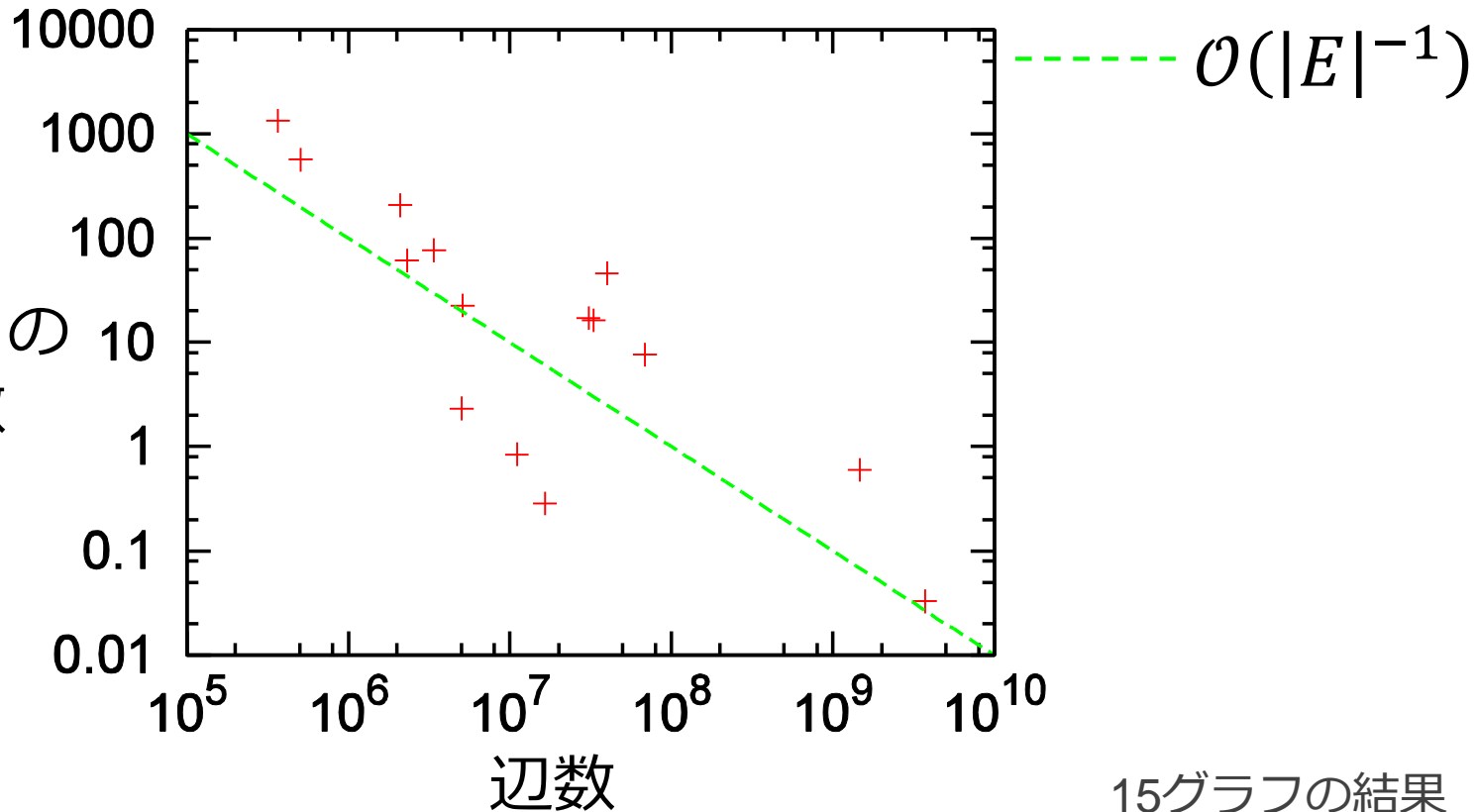
[SNAP] Stanford Large Network Dataset Collection <http://snap.stanford.edu/data/>

Environment: Intel Xeon E5-2690 2.90GHz CPU with 256GB memory

# 実験

## 性能評価: 辺数と反復回数の関係

辺挿入あたりの  
平均反復回数



辺が**多い**ほど少ない

# 実験

## 性能評価: 単一辺挿入の平均時間 & 反復回数

データセット [出典]	頂点数 $ V $	辺数 $ E $	最大 出次数 $\Delta$	平均 時間	平均 反復回数
wiki-Talk [SNAP]	2M	5M	100,022	589.6 $\mu$ s	2.3
web-Google [SNAP]	1M	5M	3,444	7.2 $\mu$ s	22.6
as-Skitter [SNAP]	2M	11M	35,387	288.4 $\mu$ s	0.8
<b>Flickr時刻</b> [KONECT]	2M	33M	26,367	95.3 $\mu$ s	16.2
<b>Wikipedia時刻</b> [KONECT]	2M	40M	6,975	76.8 $\mu$ s	46.0
soc-LiveJournal1 [SNAP]	5M	68M	20,292	17.9 $\mu$ s	7.6
twitter-2010 [LAW]	42M	1,500M	2,997,469	29,382.8 $\mu$ s	0.7
uk-2007-05 [LAW]	105M	3,700M	15,402	2.3 $\mu$ s	0.0

[KONECT] The Koblenz Network Collection <http://konect.uni-koblenz.de/networks/>

[LAW] Laboratory for Web Algorithmics <http://law.di.unimi.it/datasets.php>

[SNAP] Stanford Large Network Dataset Collection <http://snap.stanford.edu/data/>

Environment: Intel Xeon E5-2690 2.90GHz CPU with 256GB memory

# 実験

## 性能評価: 単一辺挿入の平均時間 & 反復回数

データセット [出典]	頂点数 $ V $	辺数 $ E $	最大 出次数 $\Delta$	平均 時間	平均 反復回数
wiki-Talk [SNAP]	2M	5M	100,022	589.6 $\mu$ s	2.3
web-Google [SNAP]	1M	5M	3,444	7.2 $\mu$ s	22.6
as-Skitter [SNAP]	2M	11M	35,387	288.4 $\mu$ s	0.8
Flickr <sup>時刻</sup> [KONECT]	2M	33M	26,367	95.3 $\mu$ s	16.2
Wikipedia <sup>時刻</sup> [KONECT]	2M	40M	6,975	76.8 $\mu$ s	46.0
soc-LiveJournal1 [SNAP]	5M	68M	20,292	17.9 $\mu$ s	7.6
<b>twitter-2010</b> [LAW]	42M	1,500M	<b>2,997,469</b>	<b>29,382.8</b> $\mu$ s	0.7
<b>uk-2007-05</b> [LAW]	105M	3,700M	<b>15,402</b>	<b>2.3</b> $\mu$ s	0.0

[KONECT] The Koblenz Network Collection <http://konect.uni-koblenz.de/networks/>

[LAW] Laboratory for Web Algorithmics <http://law.di.unimi.it/datasets.php>

[SNAP] Stanford Large Network Dataset Collection <http://snap.stanford.edu/data/>

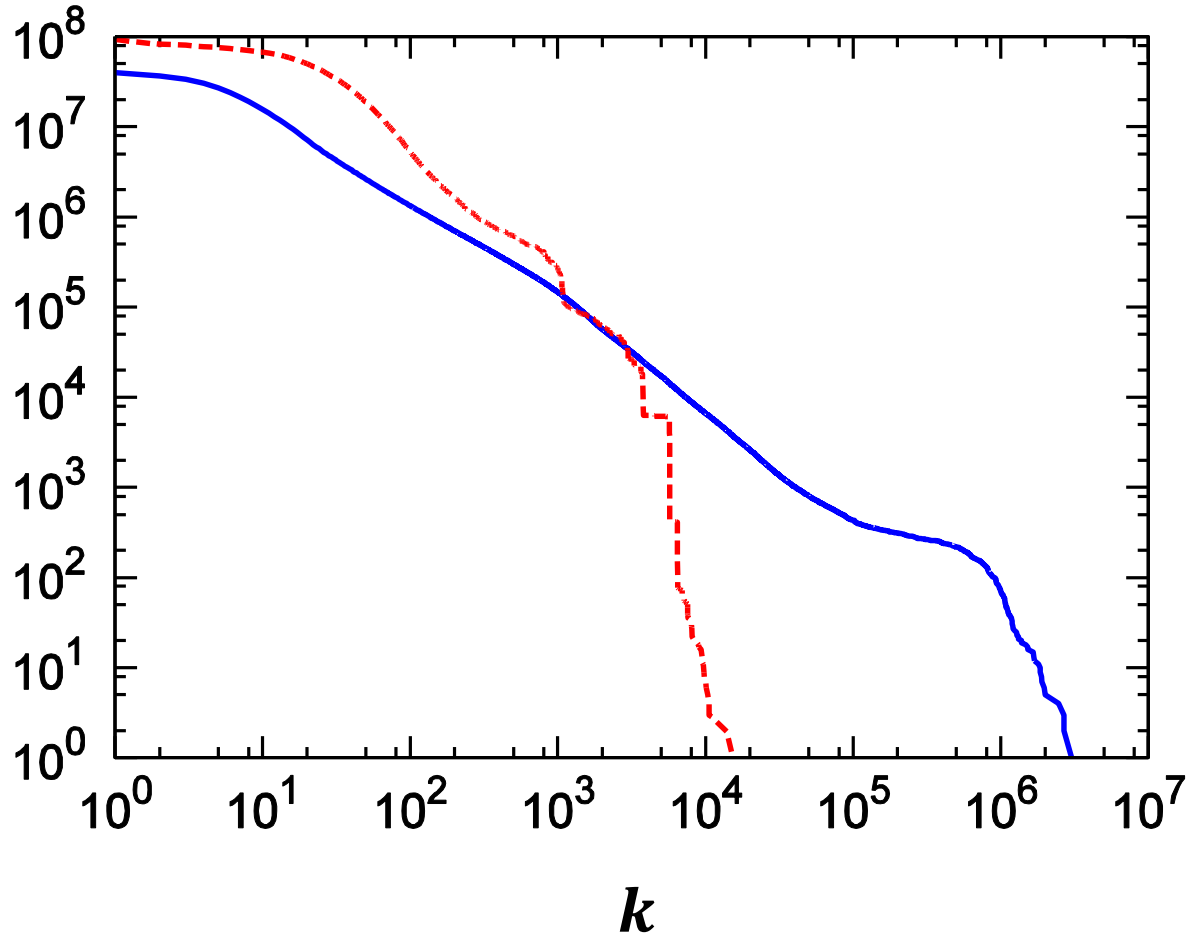
Environment: Intel Xeon E5-2690 2.90GHz CPU with 256GB memory

# 実験

## 次数分布の違い

- twitter-2010 ( $u, v$ )  $v$ が $u$ をフォロー
- - - uk-2007-05 ( $u, v$ ) ページ $u$ から $v$ へリンク

出次数  $\geq k$   
の頂点数



# 実験

## 既存手法との比較: 単一辺挿入の平均時間

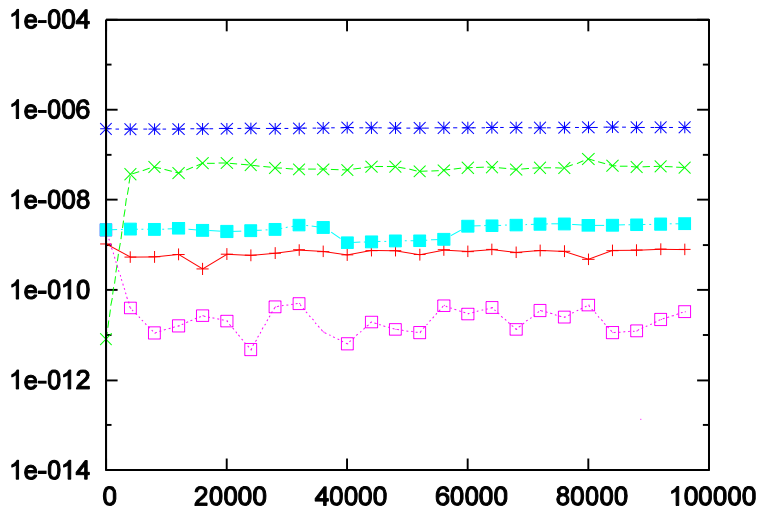
	web-Google [SNAP] $ V =1M$ $ E =5M$	Wikipedia [KONECT] $ V =2M$ $ E =40M$	uk-2007-05 [LAW] $ V =105M$ $ E =3,700M$
提案手法	7 $\mu$ s	77 $\mu$ s	2.3 $\mu$ s
Aggregation/Disaggregation [Chien et al. '04]	320 $\mu$ s	40,336 $\mu$ s	>100,000 $\mu$ s
Monte-Carlo [Bahmani et al. '10]	444 $\mu$ s	9,196 $\mu$ s	>100,000 $\mu$ s
Warm start power method	80,994 $\mu$ s	>100,000 $\mu$ s	>100,000 $\mu$ s
From scratch power method	>100,000 $\mu$ s	>100,000 $\mu$ s	>100,000 $\mu$ s

# 実験

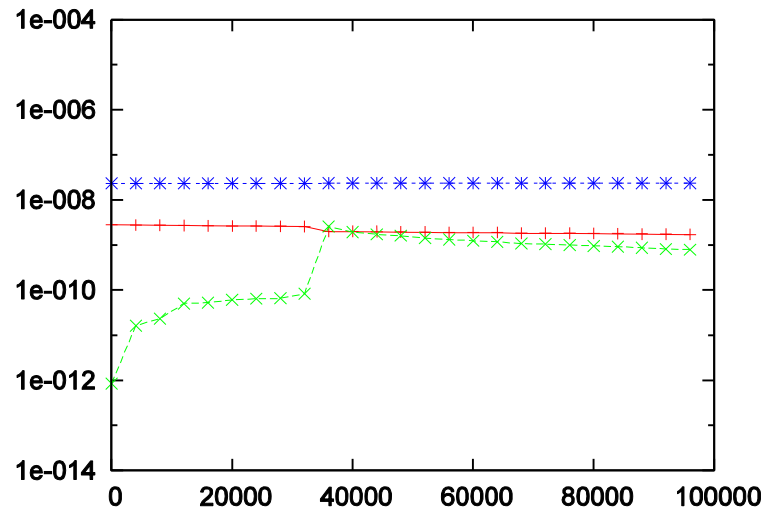
## 既存手法との比較: 精度

- +— 提案手法
- -x- - Aggregation/Disaggregation [Chien et al.'04]
- -\* - Monte-Carlo [Bahmani et al.'10]
- ...□... Warm start (power method)
- -■- From scratch (power method)

### 平均 $L_1$ 誤差の遷移



soc-Epinions1 [SNAP] |V|=76K |E|=509K



Wikipedia [KONECT] |V|=2M |E|=40M

愚直な手法に匹敵( $\sim 10^{-9}$ )



# まとめ

成長するネットワークにおける  
**Personalized PageRank追跡**のための  
**高速** & **高精度**な手法を提案

## 理論的

任意の変更にならし $\mathcal{O}(\Delta/\epsilon)$ 時間

$m$ 辺の無作為挿入に期待 $\mathcal{O}(m + \Delta \log m / \epsilon)$ 時間

$$\|x - x^*\|_{\infty} \leq \epsilon$$

## 実験的

**37億辺**をもつグラフへの単一辺挿入に**3 $\mu$ s**

**10<sup>-9</sup>**  $L_1$  誤差

# KDD'15は来週シドニー

KDD2015

CALL FOR

ATTENDING

PROGRAM

WORKSHOPS

TUTORIALS

KDD CUP

SPONSORSHIP

ORGANISERS

## Research Session RT17: Social and Graphs 3

Wednesday 1:00 pm–3:00 pm | Level 3 – Ballroom B

Chair: Tina Eliassi-Rad

### Edge-Weighted Personalized **PageRank**: Breaking A Decade-Old Performance Barrier

Wenlei Xie,Cornell University; David Bindel,Cornell University; Alan Demers,Cornell University; Johannes Gehrke,Cornell University

(Paper ID:117)

### SEISMIC: A Self-Exciting Point Process Model for Predicting Tweet Popularity

Qingyuan Zhao,Stanford University; Murat A.,Erdogdu; Stanford University Hera,Y.; He Stanford University,Anand; Rajaraman Stanford University,Jure; Leskovec Stanford Universit

(Paper ID:819)

### Beyond Triangles: A Distributed Framework for Estimating 3-profiles of Large Graphs

Ethan R.,Elenberg; The University of Texas Karthikeyan,Shanmugam; The University of Texas Michael,Borokhovich; The University of Texas Alexandros,G.; Dimakis The University of Texa

(Paper ID:896)

### Scalable Large Near-Clique Detection in Large-Scale Networks via Sampling

Michael Mitzenmacher,Harvard University; Jakub Pachocki,Carnegie Mellon University; Richard Peng,MIT; Charalampos Tsourakakis,Harvard University; Shen Chen Xu,Carnegie Mellon University

(Paper ID:720)

### Efficient **PageRank** Tracking in Evolving Networks

Naoto Ohsaka,The University of Tokyo; Takanori Maehara,Shizuoka University; Ken-ichi Kawarabayashi,National Institute of Informatics

(Paper ID:228)

### MASCOT: Memory-efficient and Accurate Sampling for Counting Local Triangles in Graph Streams

Yongsub Lim,KAIST; U Kang,KAIST

(Paper ID:163)

