

SSDBM 2012

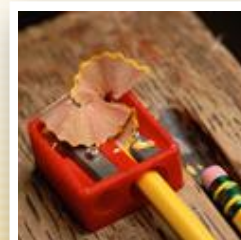
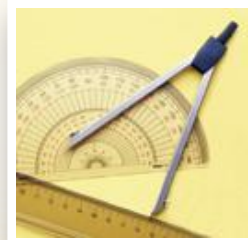


On the Efficiency of Estimating Penetrating Rank on Large Graphs

Weiren Yu¹, Jiajin Le², Xuemin Lin¹, Wenjie Zhang¹

¹ University of New South Wales & NICTA, Australia

² Donghua University, China



Contents



1. Introduction



2. Problem Definition



3. Optimization Techniques

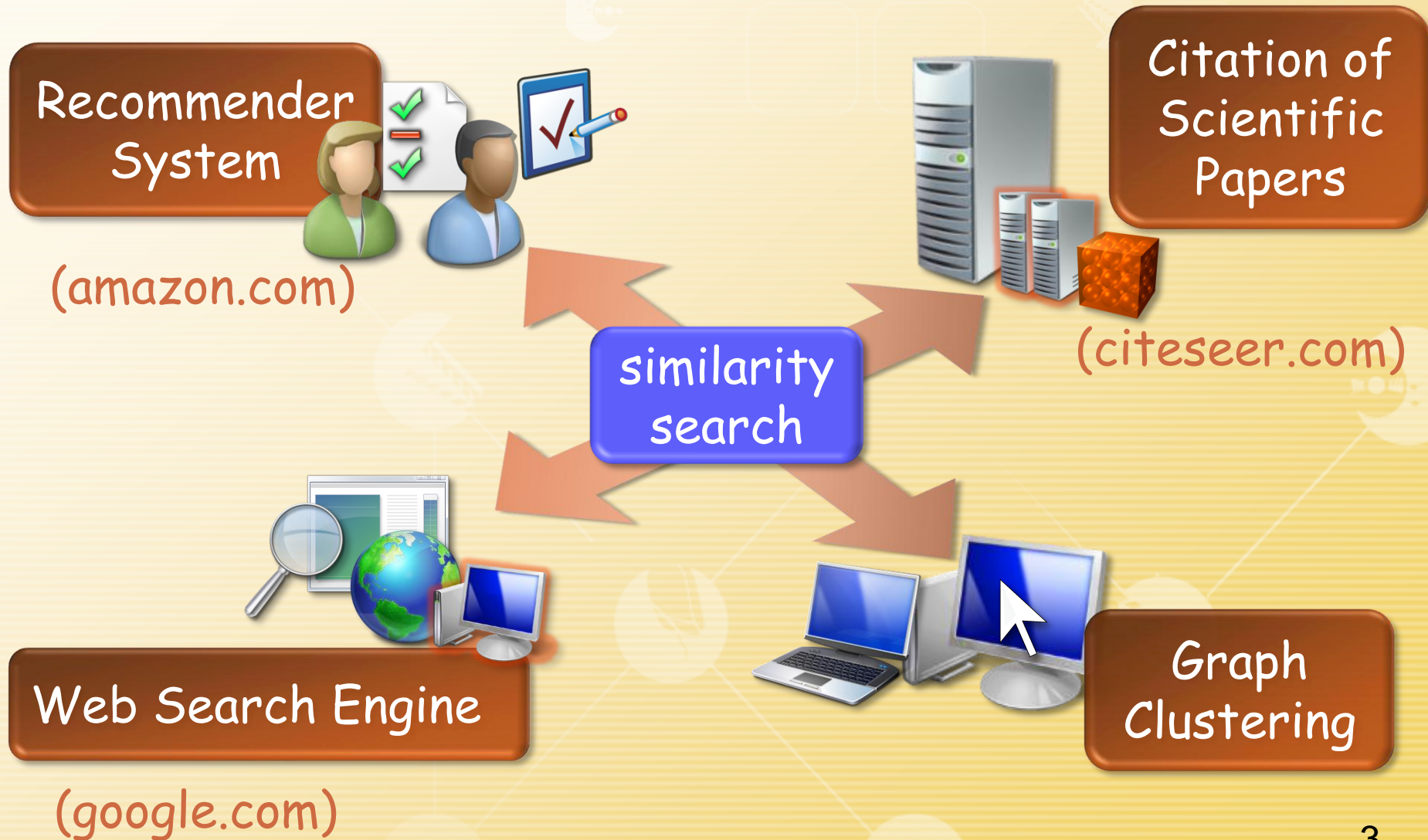


4. Experimental Results



1. Introduction

- ❖ Many applications require a measure of “similarity” between objects.



P-Rank : A New Link-based Similarity Measure

❖ Structural Similarity Measure

❖ PageRank [Page et. al, 1999]

❖ SimRank [Jeh and Widom, KDD 02]

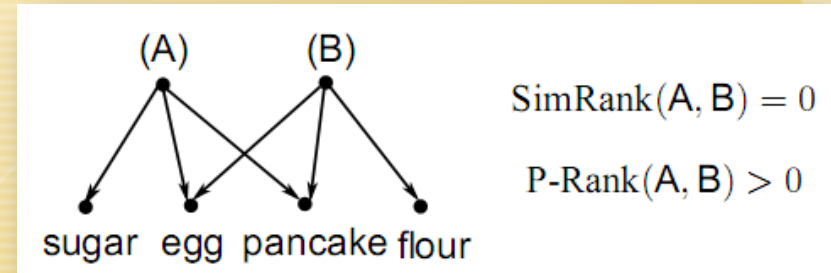
❖ P(enetrating)-Rank similarity

❖ A new promising structural measure [Zhao et. al. , CIKM 09]

❖ An extension of SimRank metrics

❖ Basic Philosophy

- ❖ Two entities are similar, if
 - (1) they are referenced by similar entities
 - (2) they reference similar entities



P-Rank Overview

❖ Features

- ❖ Avoiding “*limited information problem*” of SimRank
 - By taking account of both in- and out-links
- ❖ Defined recursively and is computed iteratively
- ❖ Applicable to any domain with object-to-object relationships

❖ Challenges

- ❖ Costly to compute P-Rank on large graphs
 - ❖ Naïve Iteration $O(Kn^4)$ [Zhao et. al. , CIKM 09]
 - ❖ Partial Sums Amortization $O(Kn^3)$ [Lizorkin et. al. , PVLDB 08]
- ❖ Hard to estimate the error for P-Rank approximation
 - ❖ Radius- and category-based Pruning Rule $O(Kd^2n^2)$
[Zhao et. al. , CIKM 09]

P-Rank Formulation

❖ Mathematical Formula

$$s(u, u) = 1;$$

$$s(u, v) = \underbrace{\frac{\lambda \cdot C_{\text{in}}}{|\mathcal{I}(u)| |\mathcal{I}(v)|} \sum_{i=1}^{|\mathcal{I}(u)|} \sum_{j=1}^{|\mathcal{I}(v)|} s(\mathcal{I}_i(u), \mathcal{I}_j(v))}_{\text{in-link part}} + \underbrace{\frac{(1 - \lambda) \cdot C_{\text{out}}}{|\mathcal{O}(u)| |\mathcal{O}(v)|} \sum_{i=1}^{|\mathcal{O}(u)|} \sum_{j=1}^{|\mathcal{O}(v)|} s(\mathcal{O}_i(u), \mathcal{O}_j(v))}_{\text{out-link part}}.$$

❖ Iterative Paradigm

$$\lim_{k \rightarrow \infty} s^{(k)}(u, v) = \sup_{k \geq 0} \{s^{(k)}(u, v)\} = s(u, v)$$

$$s^{(k+1)}(u, u) = 1.$$

$$s^{(k+1)}(u, v) = \frac{\lambda \cdot C_{\text{in}}}{|\mathcal{I}(u)| |\mathcal{I}(v)|} \sum_{i=1}^{|\mathcal{I}(u)|} \sum_{j=1}^{|\mathcal{I}(v)|} s^{(k)}(\mathcal{I}_i(u), \mathcal{I}_j(v)) + \frac{(1 - \lambda) \cdot C_{\text{out}}}{|\mathcal{O}(u)| |\mathcal{O}(v)|} \sum_{i=1}^{|\mathcal{O}(u)|} \sum_{j=1}^{|\mathcal{O}(v)|} s^{(k)}(\mathcal{O}_i(u), \mathcal{O}_j(v)).$$

Contributions

- ❖ Characterizing P-Rank as two forms
 - ❖ matrix inversion --- deterministic optimization
 - ❖ power series --- probabilistic computation
- ❖ Deterministic optimization (off-line)
 - ❖ eliminating neighborhood structure redundancy
 - ❖ quadratic-time with an error bound
- ❖ Probabilistic computation (on-line)
 - ❖ a sampling approach
 - ❖ linear-time with controlled accuracy

P-Rank Matrix Form

❖ Iterative Form

$$s(u, u) = 1;$$

$$s(u, v) = \underbrace{\frac{\lambda \cdot C_{\text{in}}}{|\mathcal{I}(u)| |\mathcal{I}(v)|} \sum_{i=1}^{|\mathcal{I}(u)|} \sum_{j=1}^{|\mathcal{I}(v)|} s(\mathcal{I}_i(u), \mathcal{I}_j(v))}_{\text{in-link part}} + \underbrace{\frac{(1 - \lambda) \cdot C_{\text{out}}}{|\mathcal{O}(u)| |\mathcal{O}(v)|} \sum_{i=1}^{|\mathcal{O}(u)|} \sum_{j=1}^{|\mathcal{O}(v)|} s(\mathcal{O}_i(u), \mathcal{O}_j(v))}_{\text{out-link part}}.$$

❖ Matrix Form

$$q_{i,j} \triangleq \begin{cases} a_{j,i} / \sum_{j=1}^n a_{j,i}, & \text{if } \mathcal{I}(i) \neq \emptyset; \\ 0, & \text{if } \mathcal{I}(i) = \emptyset. \end{cases}$$

$$p_{i,j} \triangleq \begin{cases} a_{i,j} / \sum_{j=1}^n a_{i,j}, & \text{if } \mathcal{O}(i) \neq \emptyset; \\ 0, & \text{if } \mathcal{O}(i) = \emptyset. \end{cases}$$

$$\mathbf{S} = \lambda C_{\text{in}} \cdot \mathbf{Q} \cdot \mathbf{S} \cdot \mathbf{Q}^T + (1 - \lambda) C_{\text{out}} \cdot \mathbf{P} \cdot \mathbf{S} \cdot \mathbf{P}^T + (1 - \lambda C_{\text{in}} - (1 - \lambda) C_{\text{out}}) \cdot \mathbf{I}_n,$$

$$\mathbf{S} = \lambda \cdot C_{\text{in}} \cdot \mathbf{Q} \cdot \mathbf{S} \cdot \mathbf{Q}^T + (1 - \lambda) \cdot C_{\text{out}} \cdot \mathbf{P} \cdot \mathbf{S} \cdot \mathbf{P}^T + \mathbf{I}_n$$

P-Rank is a Linear Matrix Equation

❖ Key Observation

$$\text{vec}(\mathbf{A} \cdot \mathbf{X} \cdot \mathbf{B}) = (\mathbf{B}^T \otimes \mathbf{A}) \cdot \text{vec}(\mathbf{X})$$

$$\mathbf{S} = \lambda \cdot C_{\text{in}} \cdot \mathbf{Q} \cdot \mathbf{S} \cdot \mathbf{Q}^T + (1 - \lambda) \cdot C_{\text{out}} \cdot \mathbf{P} \cdot \mathbf{S} \cdot \mathbf{P}^T + \mathbf{I}_n$$



$$(\mathbf{I} - \mathbf{M})^{-1} \mathbf{b}$$



$$\mathbf{x} = \mathbf{M} \cdot \mathbf{x} + \mathbf{b}$$



$$\mathbf{x} = \sum_{i=0}^{\infty} \mathbf{M}^i \cdot \mathbf{b}$$

$$\mathbf{b} = \text{vec}(\mathbf{I}_n)$$

$$\mathbf{x} = \text{vec}(\mathbf{s})$$

$$\mathbf{M} = \lambda \cdot C_{\text{in}} \cdot (\mathbf{Q} \otimes \mathbf{Q}) + (1 - \lambda) \cdot C_{\text{out}} \cdot (\mathbf{P} \otimes \mathbf{P})$$

Two Representations of P-Rank Solution

❖ Matrix Inversion Form

$$\text{vec}(\mathbf{S}) = [\mathbf{I}_{n^2} - \lambda C_{in} (\mathbf{Q} \otimes \mathbf{Q}) - (1 - \lambda) C_{out} (\mathbf{P} \otimes \mathbf{P})]^{-1} \cdot \text{vec}(\mathbf{I}_n).$$

huge size !!!

❖ Power Series Form

$$\text{vec}(\mathbf{S}) = \sum_{i=0}^{\infty} [\lambda \cdot C_{in} \cdot (\mathbf{Q} \otimes \mathbf{Q}) + (1 - \lambda) \cdot C_{out} \cdot (\mathbf{P} \otimes \mathbf{P})]^i \cdot \text{vec}(\mathbf{I}_n).$$

P-Rank Deterministic Optimization

❖ Basic Idea

- ❖ Most real-world graphs are low-rank and sparse.

$$\text{vec}(\mathbf{S}) = [\mathbf{I}_{n^2} - \underbrace{\lambda C_{in}(\mathbf{Q} \otimes \mathbf{Q})}_{\mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T} - \underbrace{(1 - \lambda) C_{out}(\mathbf{P} \otimes \mathbf{P})}_{\mathbf{U}_2 \Sigma_2 \mathbf{V}_2^T}]^{-1} \cdot \text{vec}(\mathbf{I}_n).$$

$$\mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T$$

$$\mathbf{U}_2 \Sigma_2 \mathbf{V}_2^T$$

- ❖ Extending Woodbury matrix identity

$$\begin{aligned} & \left(\begin{array}{c} \overbrace{\square}^n \\ \mathbf{I} \end{array} - \begin{array}{c} \overbrace{\square}^{r \ll n} \\ \mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T \end{array} - \begin{array}{c} \overbrace{\square}^{r \ll n} \\ \mathbf{U}_2 \Sigma_2 \mathbf{V}_2^T \end{array} \right)^{-1} \\ &= \begin{array}{c} \square \\ \mathbf{I} \end{array} + \begin{array}{c} \overbrace{\square}^{2r} \\ [\mathbf{U}_1 \ \mathbf{U}_2] \end{array} \cdot \begin{array}{c} \square^{-1} \\ \Sigma^{-1} \end{array} \cdot \begin{array}{c} \overbrace{\square}^{r \ll n} \\ \begin{bmatrix} \mathbf{V}_1^T \\ \mathbf{V}_2^T \end{bmatrix} \end{array} \end{aligned}$$

time complexity

$O(n^3)$

↓

$O(n^2 r)$

P-Rank Deterministic Optimization

- ❖ P-Rank can be solved as follows.

$$vec(\mathbf{S}) = (\tilde{\mathbf{U}}_Q \tilde{\mathbf{U}}_P) \Sigma \begin{pmatrix} \tilde{\mathbf{V}}_Q^T \\ \tilde{\mathbf{V}}_P^T \end{pmatrix} vec(\mathbf{I}_n) + vec(\mathbf{I}_n)$$

$$\Sigma = \begin{pmatrix} \frac{1}{\lambda C_{in}} \tilde{\Sigma}_Q^{-1} - \tilde{\mathbf{V}}_Q^T \tilde{\mathbf{U}}_Q & -\tilde{\mathbf{V}}_Q^T \tilde{\mathbf{U}}_P \\ -\tilde{\mathbf{V}}_P^T \tilde{\mathbf{U}}_Q & \frac{1}{(1-\lambda) C_{out}} \tilde{\Sigma}_P^{-1} - \tilde{\mathbf{V}}_P^T \tilde{\mathbf{U}}_P \end{pmatrix}^{-1}$$

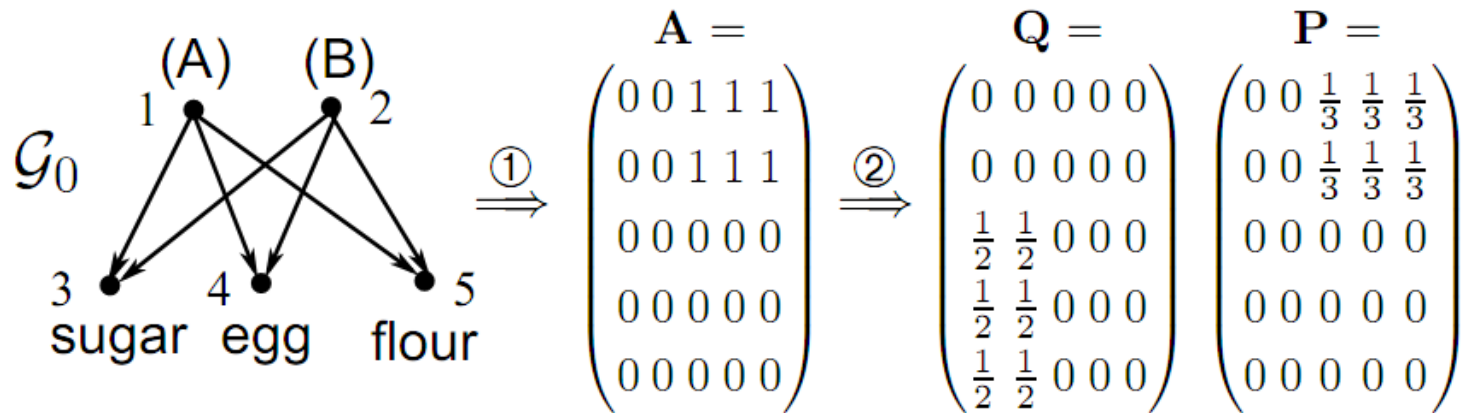
where a tilde denotes the self-Kronecker product of a matrix, e.g., $\tilde{\mathbf{U}}_Q = \mathbf{U}_Q \otimes \mathbf{U}_Q$

- ❖ Complexity

- ❖ $O(rn^2 + r^6)$ time, $O(r \cdot \max\{r^3, n\})$ space



Example



$$\xRightarrow{(3)} \quad \begin{pmatrix} U_Q = \\ 0 \\ 0 \\ -.577 \\ -.577 \\ -.577 \end{pmatrix} \quad \begin{matrix} \Sigma_Q = (1.225) & V_Q^T = (-.707 \quad -.707 \quad 0 \quad 0 \quad 0) \\ \hline V_P^T = (0 \quad 0 \quad -.577 \quad -.577 \quad -.577) & \Sigma_P = (.817) \end{matrix} \quad \begin{pmatrix} U_P = \\ 0 \\ 0 \\ -.577 \\ -.577 \\ -.577 \end{pmatrix}$$

$$\xRightarrow{(4)} \quad \begin{matrix} \Sigma_{11} = (3.33) \\ \Sigma_{12} = (-1) \\ \Sigma_{21} = (-1) \\ \Sigma_{22} = (5) \end{matrix} \quad \xRightarrow{(5)} \quad \begin{matrix} V_1 = (.383) \\ V_2 = (.277) \end{matrix} \quad \xRightarrow{(6)} \quad S = \begin{pmatrix} .569 & .069 & 0 & 0 & 0 \\ .069 & .569 & 0 & 0 & 0 \\ 0 & 0 & .564 & .064 & .064 \\ 0 & 0 & .064 & .564 & .064 \\ 0 & 0 & .064 & .064 & .564 \end{pmatrix}$$

P-Rank Deterministic Approximation

❖ P-Rank matrix inversion form

$$\text{vec}(\mathbf{S}) = [\underbrace{\mathbf{I}_{n^2} - \lambda C_{in} (\mathbf{Q} \otimes \mathbf{Q})}_{\mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T} - \underbrace{(1 - \lambda) C_{out} (\mathbf{P} \otimes \mathbf{P})}_{\mathbf{U}_2 \Sigma_2 \mathbf{V}_2^T}]^{-1} \cdot \text{vec}(\mathbf{I}_n).$$

❖ Reduced SVD for P-Rank Approximation

$$\begin{aligned} \begin{matrix} n \\ \left\{ \begin{array}{c} \mathbf{Q} \end{array} \right\} \end{matrix} &= \begin{matrix} r (\ll n) \\ \left\{ \begin{array}{c} \mathbf{U}_Q \end{array} \right\} \end{matrix} \cdot \begin{matrix} r & n-r \\ \left\{ \begin{array}{c} \Sigma_Q \end{array} \right\} \end{matrix} \cdot \begin{matrix} \left\{ \begin{array}{c} \mathbf{V}_Q \end{array} \right\} \end{matrix} \begin{matrix} r \\ n-r \end{matrix} \\ &\approx \Rightarrow \|\mathbf{Q}_v - \mathbf{Q}\|_2 = \sigma_{v+1} \\ \begin{matrix} n \\ \left\{ \begin{array}{c} \mathbf{Q}_v \end{array} \right\} \end{matrix} &= \begin{matrix} v & n-r \\ \left\{ \begin{array}{c} \mathbf{U}'_Q \end{array} \right\} \end{matrix} \cdot \begin{matrix} v & n-r \\ \left\{ \begin{array}{c} \Sigma'_Q \end{array} \right\} \end{matrix} \cdot \begin{matrix} \left\{ \begin{array}{c} \mathbf{V}'_Q \end{array} \right\} \end{matrix} \begin{matrix} v \\ n-r \end{matrix} \end{aligned}$$

Rank r SVD $O(rn^2)$ ($r \ll n$)

Low Rank v SVD $O(vn^2)$ ($v \ll r$)



P-Rank Deterministic Approximation

❖ Approximation Error

$$\epsilon_v \leq \frac{\lambda C_{in} \sigma_1 \sigma_{v+1} + (1-\lambda) C_{out} \bar{\sigma}_1 \bar{\sigma}_{v+1}}{1 - \lambda C_{in} - (1-\lambda) C_{out}} r$$

e.g., WIKI 0715 ($r = 15\text{K}$, $\sigma_1 = 1.12$, $\bar{\sigma}_1 = 1.08$)

Setting $C_{in} = 0.8$, $C_{out} = 0.6$, and $\lambda = 0.5$

$$\epsilon_v \leq \frac{0.5 \times 0.8 \times 1.12 + 0.5 \times 0.6 \times 1.08}{1 - 0.5 \times 0.8 - 0.5 \times 0.6} \times 10^{-7} \times 15\text{K} = 0.0039$$

❖ Complexity

Time: $O(vn^2 + v^6)$ with $v \leq r$

Space: $O(v \cdot \max\{v^3, n\})$

P-Rank Probabilistic Computation

❖ Key Observation

- ❖ P-Rank can be viewed as a geometric sum of random walks

$$\text{vec}(\mathbf{S}) = \sum_{i=0}^{\infty} [\lambda \cdot C_{in} \cdot (\mathbf{Q} \otimes \mathbf{Q}) + (1 - \lambda) \cdot C_{out} \cdot (\mathbf{P} \otimes \mathbf{P})]^i \cdot \text{vec}(\mathbf{I}_n)$$

- ❖ $s(u, v)$ represents how soon two surfers are expected to meet

❖ Main idea

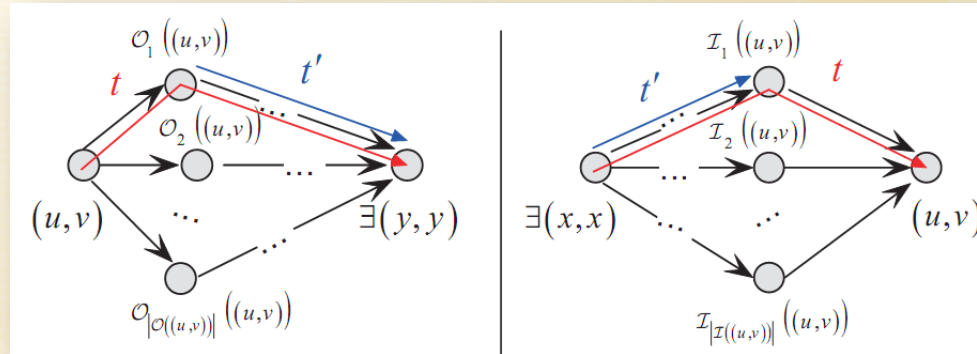
- ❖ utilize the first hitting time $\tau(u, v)$ of coalescing walks to estimate $s(u, v)$

$$s(u, v) = \mathbb{E}(\lambda \cdot C_{in}^{\tau_1(u, v)} + (1 - \lambda) \cdot C_{out}^{\tau_2(u, v)})$$

P-Rank Probabilistic Computation

❖ Random Surfer Model

❖ one-step path transformation $T : t' \rightarrow t$



❖ length

$$l(t) = l(t') + 1$$

❖ probability

$$p(T(t')) = \begin{cases} \frac{1}{|\mathcal{I}((u,v))|} \cdot p(t'), & t' : \exists(x,x) \rightarrow (u,v); \\ \frac{1}{|\mathcal{O}((u,v))|} \cdot p(t'), & t' : (u,v) \rightarrow \exists(y,y). \end{cases}$$

❖ Equivalence of Sampling approach

$$s(u,v) = \lambda \cdot \sum_{t: \exists(x,x) \rightarrow (u,v)} p(t) \cdot C_{\text{in}}^{l(t)} + (1 - \lambda) \cdot \sum_{t: (u,v) \rightarrow \exists(y,y)} p(t) \cdot C_{\text{out}}^{l(t)}.$$

$$= \frac{\lambda \cdot C_{\text{in}}}{|\mathcal{I}(u)| \cdot |\mathcal{I}(v)|} \cdot \sum_{i=1}^{|\mathcal{I}(u)|} \sum_{j=1}^{|\mathcal{I}(v)|} s(\mathcal{I}_i(u), \mathcal{I}_j(v)) + \frac{(1 - \lambda) \cdot C_{\text{out}}}{|\mathcal{O}(u)| \cdot |\mathcal{O}(v)|} \cdot \sum_{i=1}^{|\mathcal{O}(u)|} \sum_{j=1}^{|\mathcal{O}(v)|} s(\mathcal{O}_i(u), \mathcal{O}_j(v)).$$

P-Rank Probabilistic Computation

❖ Complexity

❖ Time $O(N \cdot n)$

❖ Space $O(n + N)$

where N : sample size, n : # of vertices

❖ Sample Size

❖ $N \geq -2 \lceil (\sigma/\epsilon)^2 \log \alpha \rceil$ suffices to ensure that

$$\Pr(|s_N - s| \geq \epsilon) < \alpha$$

❖ In practice, $N \ll n$.

e.g., on DBLP (98-07) For $n = 10K$, $\epsilon = 0.15\sigma$, $\alpha = 0.05$, we have

$$N \geq -2 \lceil 0.15^{-2} \log(0.05) \rceil = 267.$$

P-Rank Probabilistic Computation

❖ Error Bound

Let $Err \triangleq \sup_{N \geq 1} \Pr(|\hat{s}_N - s| \geq \epsilon)$

❖ upper bound - by Bernstein's Theorem

$$Err \leq \exp(-N\epsilon^2/(2\sigma^2))$$

❖ lower bound - by Central Limit Theorem

$$Err \geq \Pr\left(\left|\frac{1}{\sqrt{N}} \sum_{i=1}^N \left(\frac{\hat{s}_N^{(i)} - s}{\sigma}\right)\right| \geq \frac{\epsilon\sqrt{N}}{\sigma}\right) = 2 - 2\Phi\left(\frac{\epsilon\sqrt{N}}{\sigma}\right)$$

❖ Relative Order Preserving

If $s(u, v) > s(u, w) + \epsilon$, then

$$\Pr(\hat{s}_N(u, v) - \hat{s}_N(u, w) > \epsilon) \leq \exp(-N\epsilon^2/2)$$

Experiment

❖ Datasets

- ❖ Synthetic data (RAND 0.5M-3.5M)
- ❖ Real data (AMZN, DBLP, WIKI)

	0505	0601
$ \mathcal{V} $	410K	402K
$ \mathcal{E} $	3,356K	3,387K

Table 2: AMZN

	98-99	98-01	98-03	98-05	98-07
$ \mathcal{V} $	1,525	3,208	5,307	7,984	10,682
$ \mathcal{E} $	5,929	13,441	24,762	39,399	54,844

Table 3: DBLP

	0715	0827	0919
$ \mathcal{V} $	3,088K	3,102K	3,116K
$ \mathcal{E} $	1,126K	1,134K	1,142K

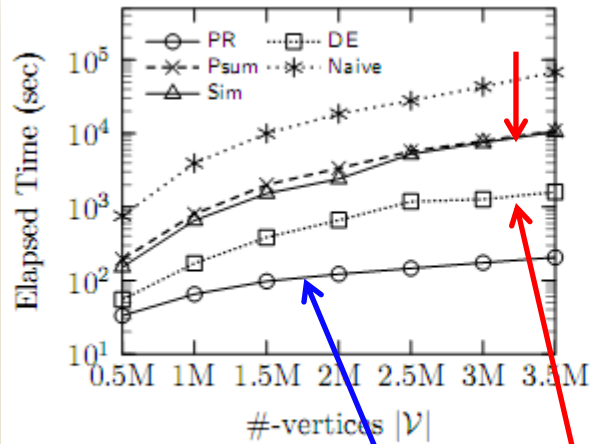
Table 4: WIKI

❖ Compared Algorithms

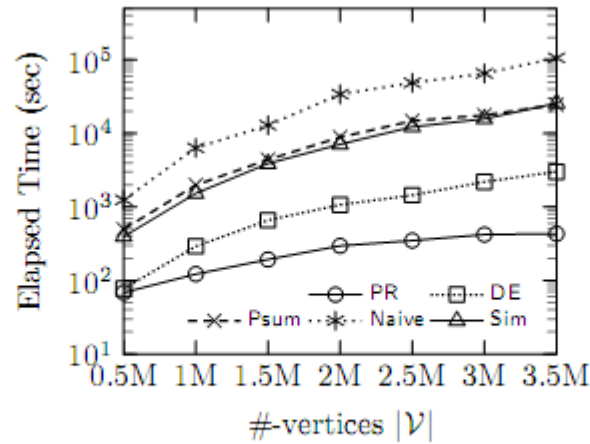
- ❖ DE P-Rank, PR P-Rank
- ❖ Naive, radius-based pruning iteration [Zhao et al, CIKM 2009]
- ❖ Psum, leveraging a partial sum function to compute P-Rank
- ❖ Sim, a SimRank algorithm, taking account of the evidence factor for incident vertices

Experiment (1)

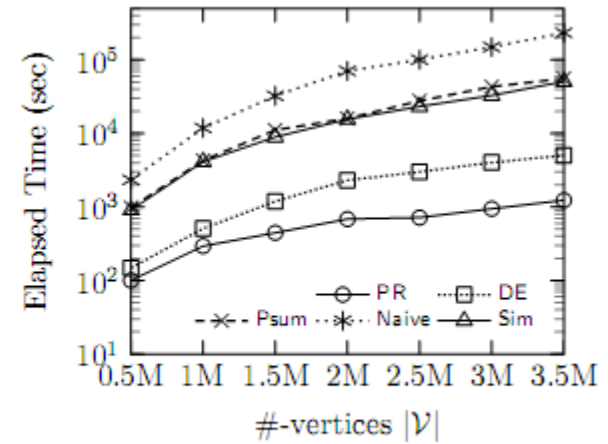
Scalability on Synthetic Datasets



(a) $|\mathcal{E}| = 2M$



(b) $|\mathcal{E}| = 4M$



(c) $|\mathcal{E}| = 6M$

DE is one-order-of-magnitude faster

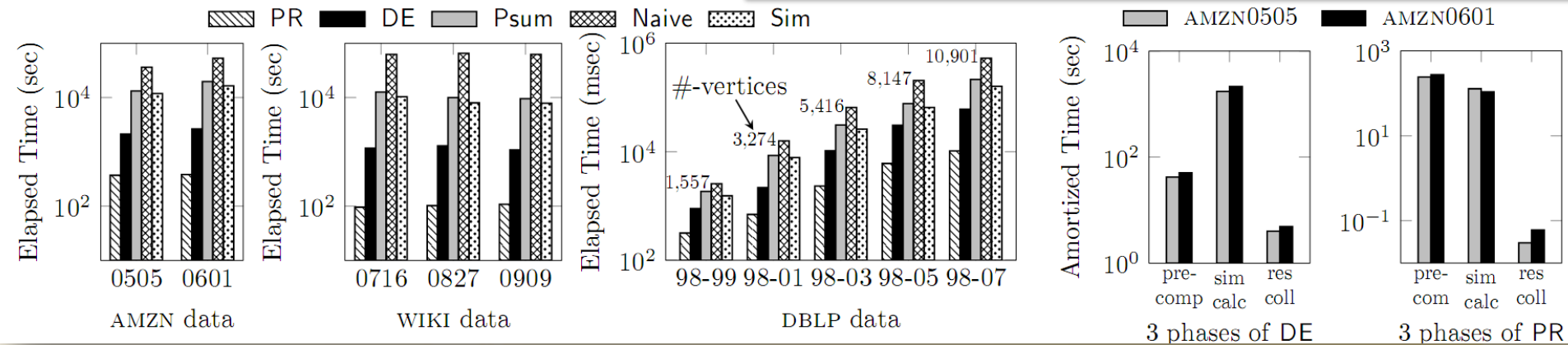
PR increases linearly with $|V|$



Experiment (2)

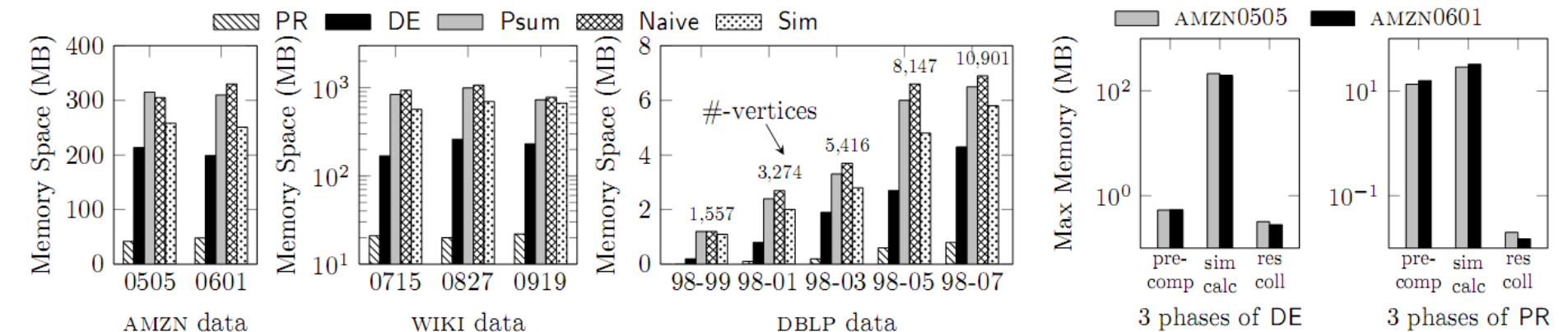
Computational Time on Real Datasets

PR outperforms the other approaches.



DE can cluster vertices with similar neighborhood.

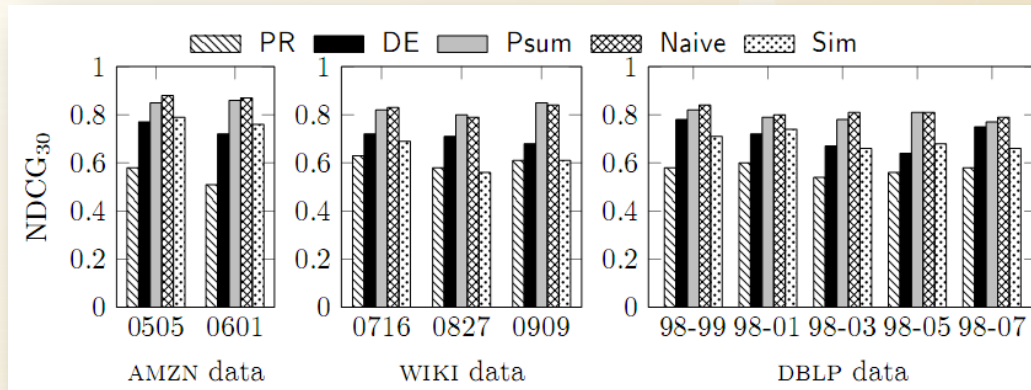
Memory Space on Real Datasets





Experiment (3)

Accuracy on Real Datasets



Rank	PR	DE	Naive
1	Shivnath Babu	Shivnath Babu	Shivnath Babu
2	Chris Olston	Yingwei Cui	Yingwei Cui
3	Jun Yang	Chris Olston	Chris Olston
4	Yingwei Cui	Jun Yang	Jun Yang
5	Rajeev Motwani	Arvind Arasu	Rajeev Motwani
6	Arvind Arasu	Rajeev Motwani	Arvind Arasu
7	David J. DeWitt	Alon Y. Halevy	Utkarsh Srivastava
8	Glen Jeh	Anish Das Sarma	David J. DeWitt
9	Utkarsh Srivastava	Omar Benjelloun	Omar Benjelloun
10	Omar Benjelloun	David J. DeWitt	Alon Y. Halevy

$$NDCG_p = \frac{1}{IDCG_p} \sum_{i=1}^p (2^{\text{rank}_i} - 1) / (\log_2(1 + i))$$

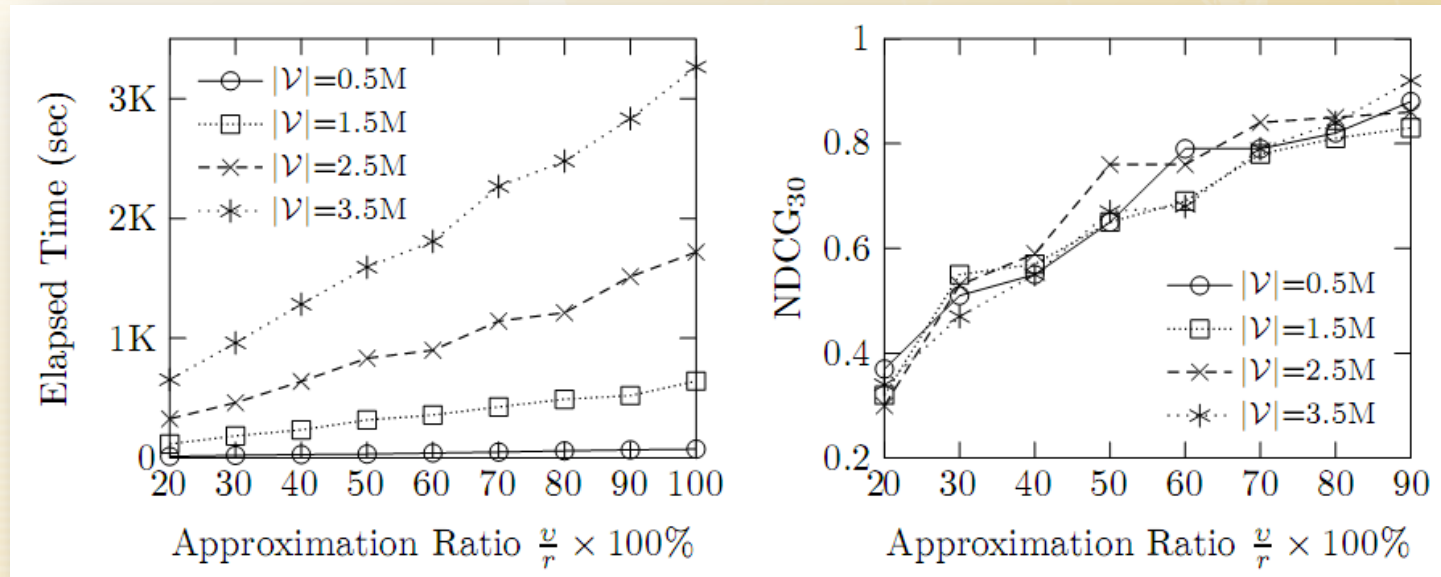
Top-10 Co-authors of Jennifer Widom on DBLP

DE achieves higher accuracy than PR.

The accuracy of PR is not that good because a few FPTs are neglected with certain probability by sampling.

Experiment (4)

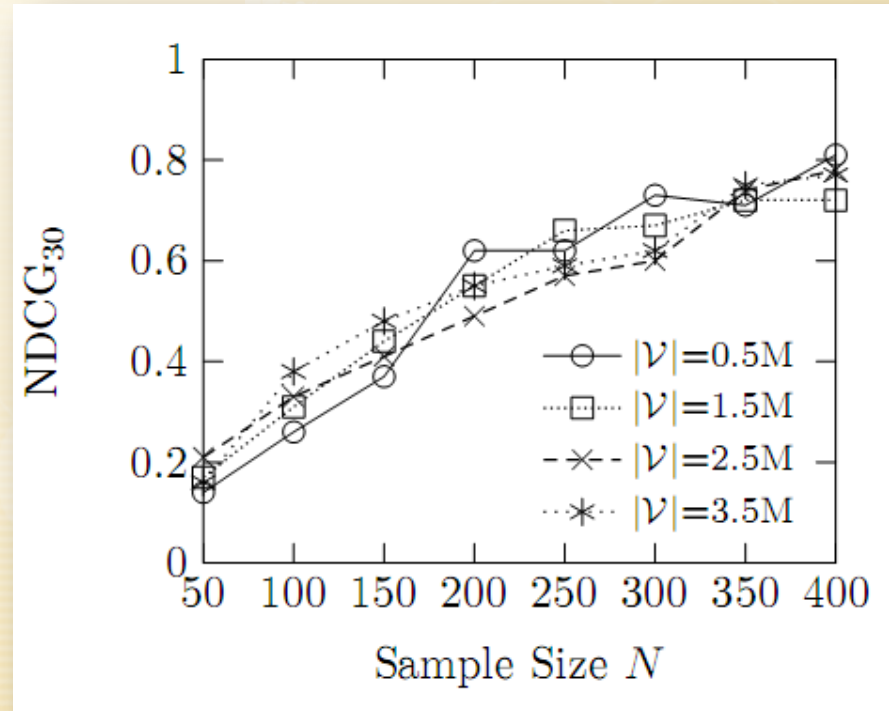
Effects of u for DE



Adding u induces smaller errors,
but increases the time up to rank r .

Experiment (5)

Effects of N for PR



Adding samples of FPTs reduces errors

When $N > 300$, higher accuracy could be expected ($NDCG_{30} > 0.6$)

Conclusions

- ❖ Two matrix forms are investigated to characterize P-Rank.
- ❖ Using matrix inversion form, we propose DE P-Rank to reduce the time from cubic to quadratic.
- ❖ By leveraging reduced SVD, the error estimate is obtained for P-Rank approximation.
- ❖ Using power series form, we present PR P-Rank to speed up the computation of P-Rank in linear time with controlled accuracy.



Thank You !