

READING:

Online Topic-aware Influence Maximization Queries

Cigdem Aslay¹ Nicola Barbieri² Francesco Bonchi² Ricardo Baeza-Yates²

¹Web Research Group
University Pompeu Fabra
Barcelona, Spain
aslayci@acm.org

²Yahoo Labs
Barcelona, Spain
{barbieri,bonchi}@yahoo-inc.com
rbaeza@acm.org

from EDBT/ICDT 2014

SPEAKER: KAZUHIRO INABA

WHAT THE PAPER IS ABOUT?

- **Influence Maximization**

- ソーシャルグラフ上を口コミで情報が伝搬していく。できるだけ大勢に情報を広める発信源 “seed set” を求めたい

- **Topic-aware**

- 情報の伝わる確率が話題ごとに異なる

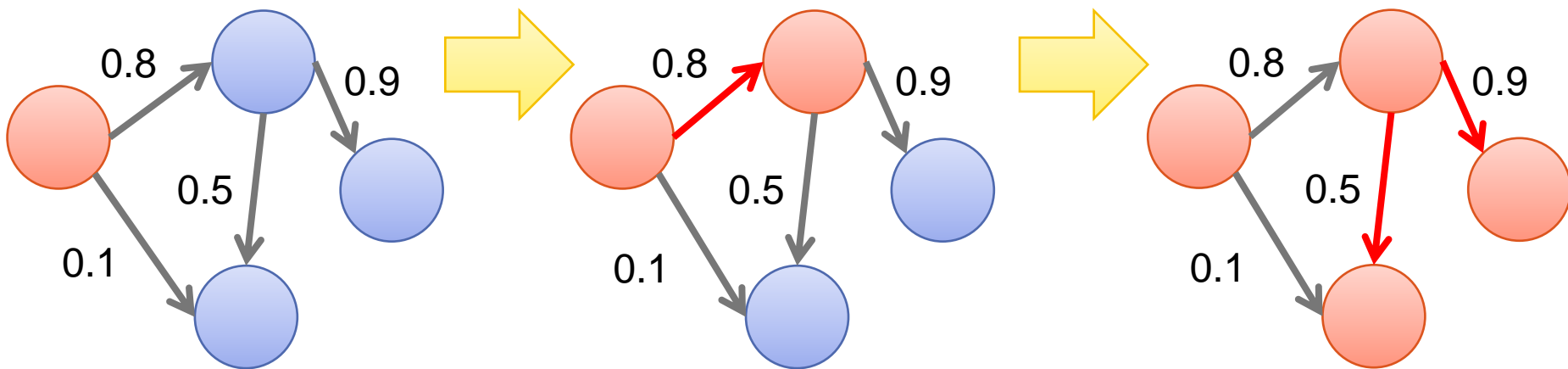
- **Online**

- 次々と違う話題の情報が飛んでくるので、その都度、よい seed set を、あらかじめ用意しておいた前計算を使って一瞬で求める

INDEPENDENT CASCADE (IC) MODEL

[Kempe, Kleinberg, and Tardos; KDD 2003]

有向辺 $u \rightarrow v$ それぞれに
確率 $p_{u,v}$ が割り当てられている



u が情報を最初に知った直後のターン
(1回だけ)に隣接する v に確率 $p_{u,v}$ で
情報が伝わる

TOPIC-AWARE IC (TIC) MODEL

- Z 個の *Topic* がある。
- 伝わる情報 (論文中では *Item*) は Z 次元ベクトル
 $\gamma = (\gamma^1, \gamma^2, \dots, \gamma^Z)$
where $0 \leq \gamma^k \leq 1, \sum \gamma^k = 1$
- 各辺に Z 個の確率 $(p^1_{u,v}, \dots, p^Z_{u,v})$ を割当て



確率 $\sum \gamma^k p^k_{u,v}$ で情報が伝わる

SOME NOTES ON

TOPIC-AWARE IC (TIC) MODEL

- Item を1つに固定すると IC-model に一致
- 実験では $Z = 10$. 映画レビューSNSで実験
 - 例: Topic = {1:ホラー, 2:恋愛, 3:SF}
 - $\gamma = (0.0, 0.8, 0.2)$ (少しSF色のあるロマンス)
 - $p_{u,v} = (0.9, 0.0, 0.4)$
(vさんはuさんの恋愛映画を見る目を全く信用していない)



確率 $\sum \gamma^k p^k_{u,v}$ で情報が伝わる

PROBLEM

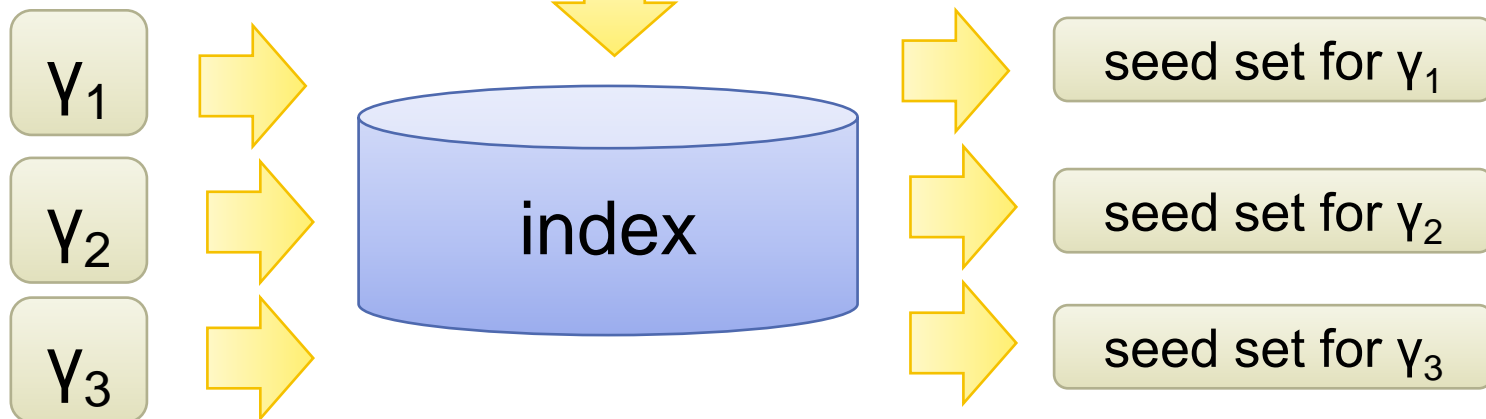
前もって1つ与えられる入力: グラフ (V, E, p)

あらかじめ
index を
構築しておくことで



Graph

実験では $|V| \sim 30k$, $|E| \sim 400k$



リアルタイムに与えられる Item y に対し
影響力の大きい seed set を即座に返したい

実験では $|\text{seed}| \leq 50$

PROBLEM?

- 「Item を1つに固定すると IC-model に一致」
 - Topic のことを忘れて、毎回 IC-model での影響最大化の既存アルゴリズムを実行すれば良いのでは？
 - 遅い。論文での比較対象は CELF++ [Goyal et al. WWW'11] で、数日かかる
 - [Ohsaka et al. AACL'14] [Tang et al. SIGMOD'14] などと比べるとどうでしょうか

OVERVIEW OF THE APPROACH

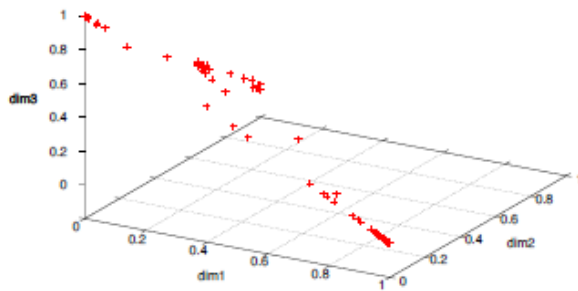
確率のついたグラフ (V, E, p) が与えられる

1. h 個の Item を代表点として選ぶ (実験では $h = 1000$)
2. 各代表点について、あらかじめ seed list を既存の IC-model 用オフラインアルゴリズムで計算しておく (実験では CELF++)

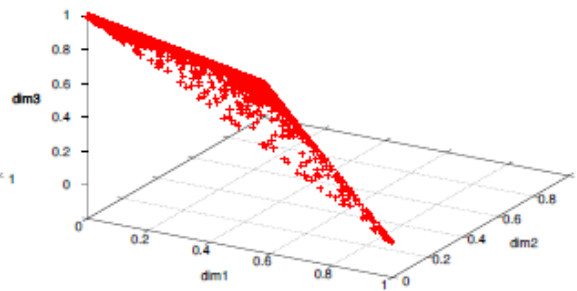
クエリ y が与えられる

3. y に近い代表点をいくつか (10個) 取り出す
4. それぞれの代表点での Seed List を巧く混ぜる

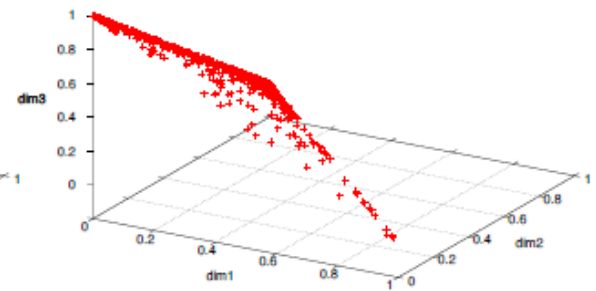
1. INDEX POINTS



(a) Distribution of TIC items



(b) Distribution of sampled items from $Dirichlet(\alpha)$



(c) Distribution of index items

元のデータセットに
含まれている Item
(レビューされた映画)



から学習した
Dirichlet分布で
Itemを大量に
ランダム生成



したものを
KL-divergence で
K-means++ で
h 個にクラスタリング

the hyper-parameters $\alpha = \{\alpha_1, \dots, \alpha_z\}$ which define the Dirichlet distribution that maximizes:

$$\prod_{i \in \mathcal{I}} P(\vec{\gamma}_i | \alpha) = \prod_{i \in \mathcal{I}} \frac{\Gamma(\sum_z \alpha_z)}{\prod_z \Gamma(\alpha_z)} \prod_z (\gamma_i^z)^{\alpha_z - 1}.$$

$$D_{KL}(P \| Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

$$D_{KL}(\vec{\gamma}_i \| \vec{\gamma}_q)$$

2. PRECOMPUTE

Index Point $H = \{\gamma_1, \dots, \gamma_h\}$ のそれぞれについて、
IC model での既存手法で $L (= 50)$ ノードの seed
list を計算



(p^1, \dots, p^Z)



$\sum \gamma_1^k p^k$



$\tau_1 = [V_{1,1}, V_{1,2}, \dots, V_{1,L}]$



$\sum \gamma_h^k p^k$



$\tau_2 = [V_{h,1}, V_{h,2}, \dots, V_{h,L}]$

(CELF++ を使用。

実験に使った述べ時間: 平均60時間 $\times h$)

OVERVIEW OF THE APPROACH

確率のついたグラフ (V, E, p) が与えられる

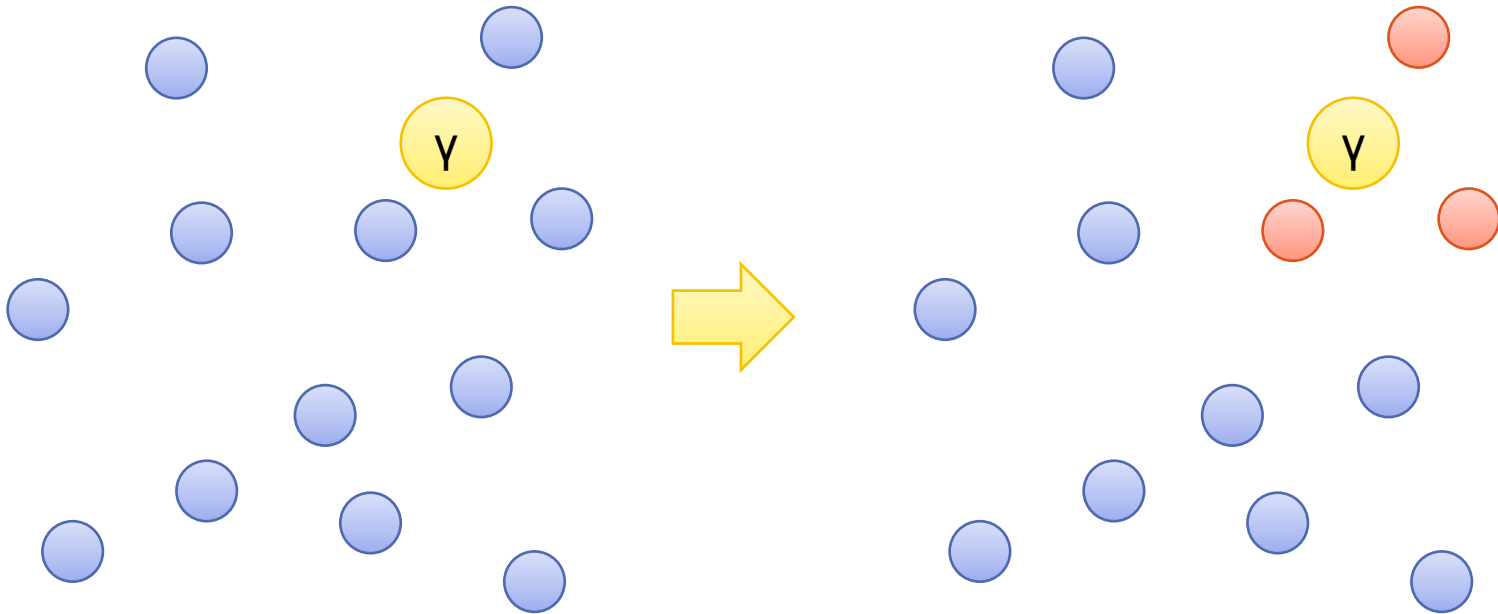
1. h 個の Item を代表点として選ぶ (実験では $h = 1000$)
2. 各代表点について、あらかじめ seed list を既存の IC-model 用オフラインアルゴリズムで計算しておく (実験では CELF++)

クエリ y が与えられる

3. y に近い代表点をいくつか (10個) 取り出す
4. それぞれの代表点での Seed List を巧く混ぜる

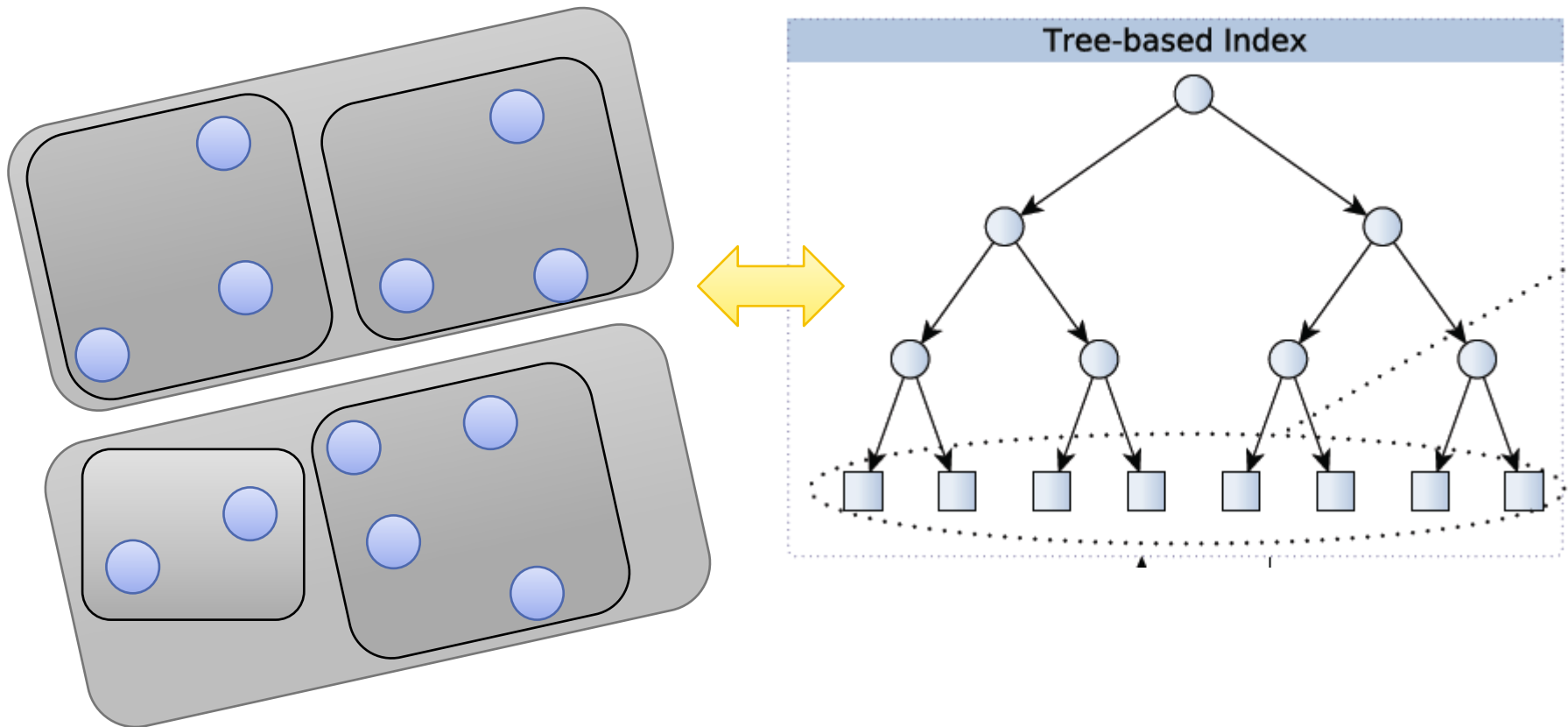
3. SEARCH IN INDEX

- 新しくアイテム y が与えられたときに、Index Point 集合 H から y に近い点をいくつか取り出す



3. SEARCH IN INDEX: HOW

Index集合を、KL-divergenceで定義されたballで繰り返し分割して、tree状に保持
("Bregman ball tree")



3. SEARCH IN INDEX: HEURISTICS

(How to search 10 neighbors from the tree index?)

- **Approximate K-NN:**
 - leafを5個訪れたら探索打ち切り
- **Anderson-Darling test:**
 - 含まれているIndex Pointに対しクエリ q を統計的に検定し、妥当ならそこで探索打ち切り
- **Selection:**
 - Rankを混ぜるときにKL-divergenceを正規化した重み付きで混ぜるので、重み0.5%
- **INFLEX (提案手法): 3つすべて**

3. SEARCH IN INDEX: HEURISTICS

(How to search 10 neighbors from the tree index?)

Seed list の質 (0 = CELF++と一致)

計算時間

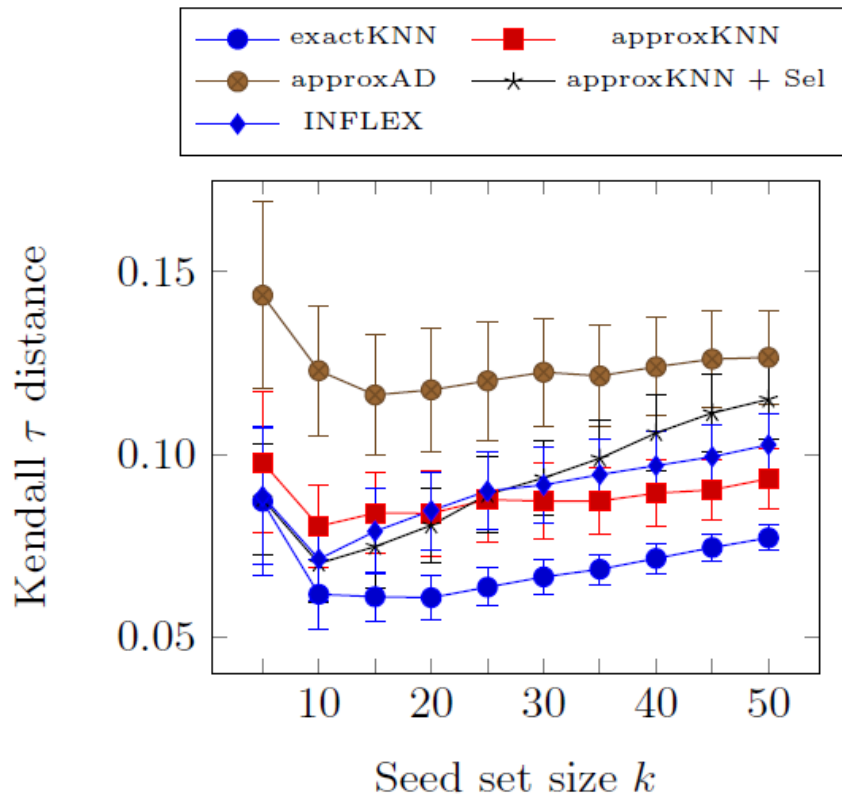


Figure 6: Accuracy comparison.

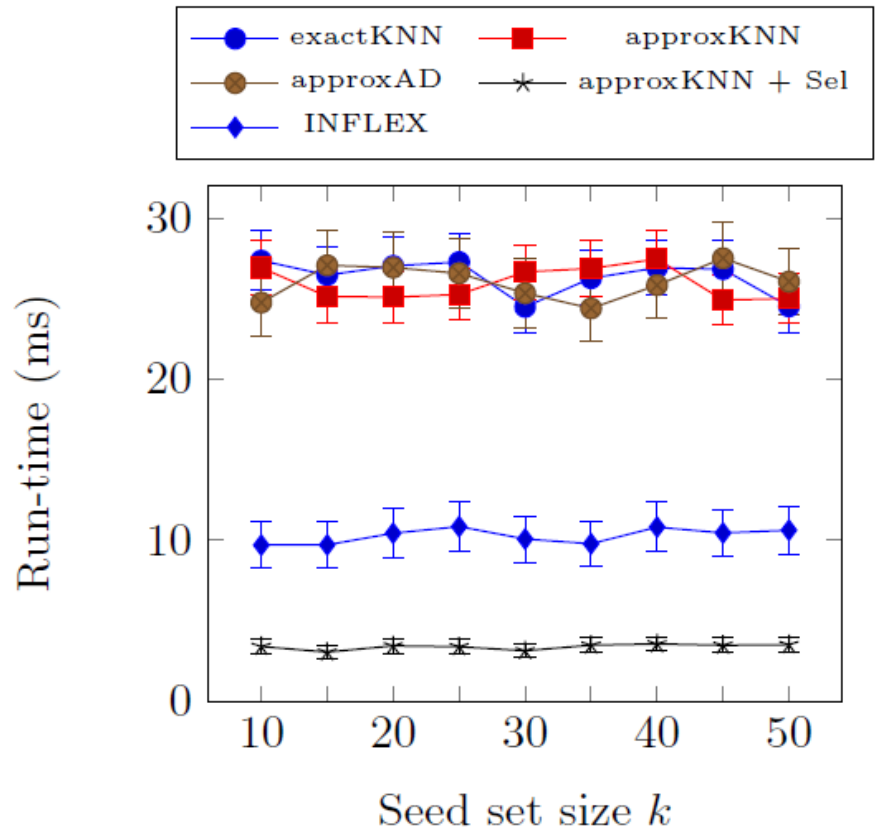


Figure 7: Run-time comparison.

4. RANK AGGREGATION

- h個の Index point から得られたノードのランキング ($0 \leq w \leq 1$ の重み付き) を混ぜる

- $(w_1, [v_{1,1}, v_{1,2}, \dots, v_{1,L}])$

- ...

- $(w_h, [v_{h,1}, v_{h,2}, \dots, v_{h,L}])$

$$W(\vec{\gamma}_i, \vec{\gamma}_q) = \frac{e^{KL_{max}} - e^{D_{KL}(\vec{\gamma}_i \| \vec{\gamma}_q)}}{1 - e^{-KL_{max}}}$$

4. RANK AGGREGATION

Weighted Copeland score

$$\text{Copeland}^w(v) = \sum_{v'} \sum_j \{w_j \mid \tau_j(v) < \tau_j(v')\}$$

where $\tau_j(v) :=$ Index Point j の seed list での v の順位
の高い順に並べる

Table 1: Kendall- τ distance between the seed sets produced by aggregation algorithms and the ground truth computed by standard offline influence maximization computation.

Seed Set size k	<i>Borda</i>	<i>Borda</i> ^w	<i>Copeland</i>	<i>Copeland</i> ^w
5	0.100	0.096	0.104	0.087
10	0.073	0.066	0.068	0.062
15	0.071	0.065	0.068	0.061
20	0.068	0.063	0.068	0.061
25	0.068	0.066	0.069	0.064
30	0.068	0.067	0.071	0.066
35	0.071	0.069	0.072	0.069
40	0.074	0.073	0.075	0.072
45	0.079	0.076	0.077	0.075
50	0.081	0.080	0.079	0.077

EXPERIMENTAL SETTINGS

- Flixster (映画レビューSNS) dataset
 - 30k users, 425k social links, 12k movies
 - “*User u rated Item m on Time t*” 形式のログ
 - 著者らが ICDM'12 で提案した手法で、Topic数 $Z (= 10)$ を決め打って各linkの p と各movieの γ を推測したデータを用いる。

FINAL RESULT

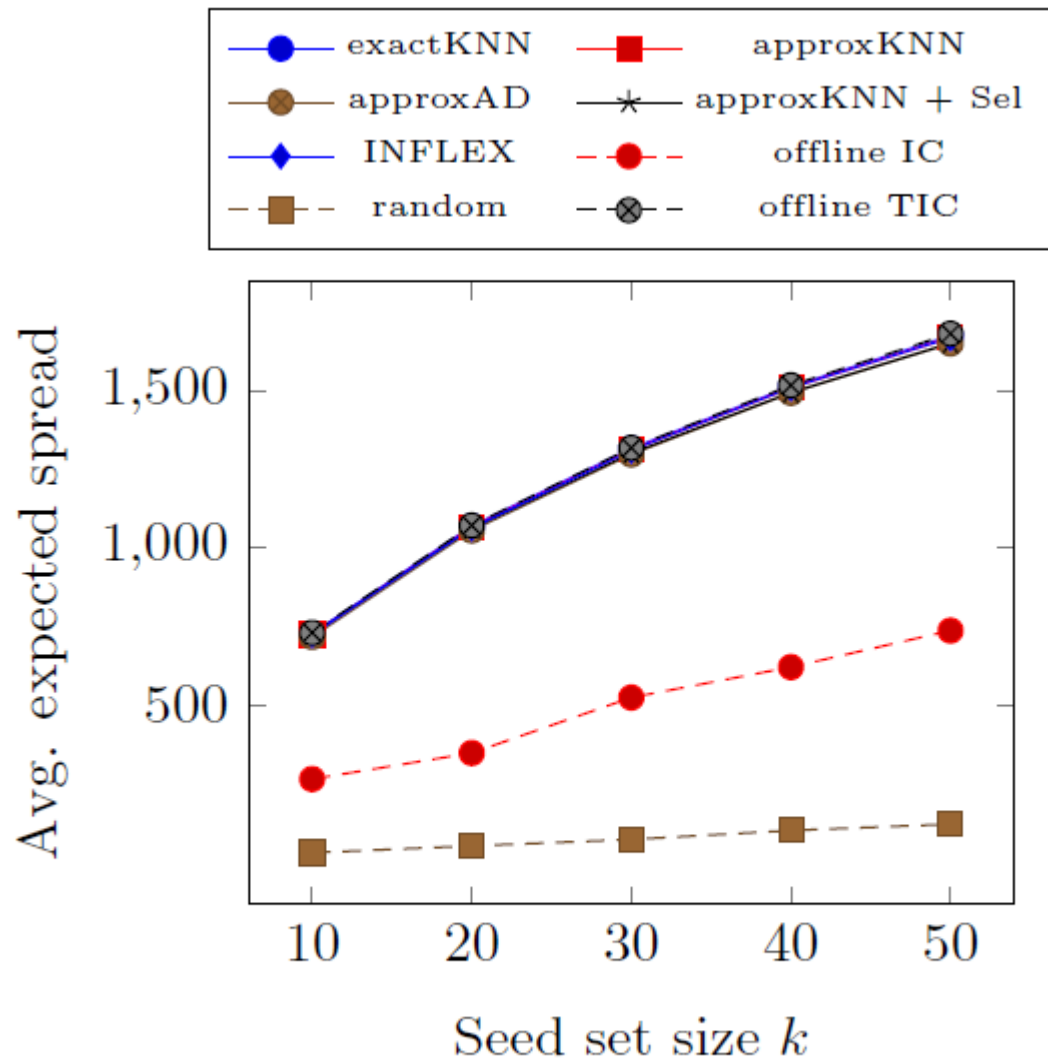


Figure 8: Expected spread comparison.