



Cross-domain Link Prediction and Recommendation

Jie Tang

Department of Computer Science and Technology
Tsinghua University

Networked World

facebook

- **1.26 billion** users
- **700 billion** minutes/month



- **280 million** users
- **80% of users** are 80-90's

twitter



- **555 million** users
- **.5 billion** tweets/day



- **560 million** users
- **influencing** our daily life

amazon.com

- **79 million** users per month
- **9.65 billion** items/year



- **500 million** users
- **35 billion** on 11/11



- **800 million** users
- **~50% revenue** from network life



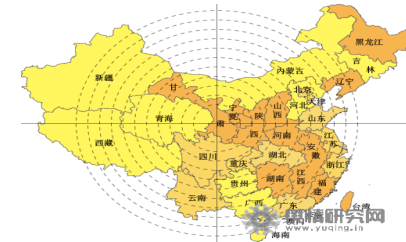
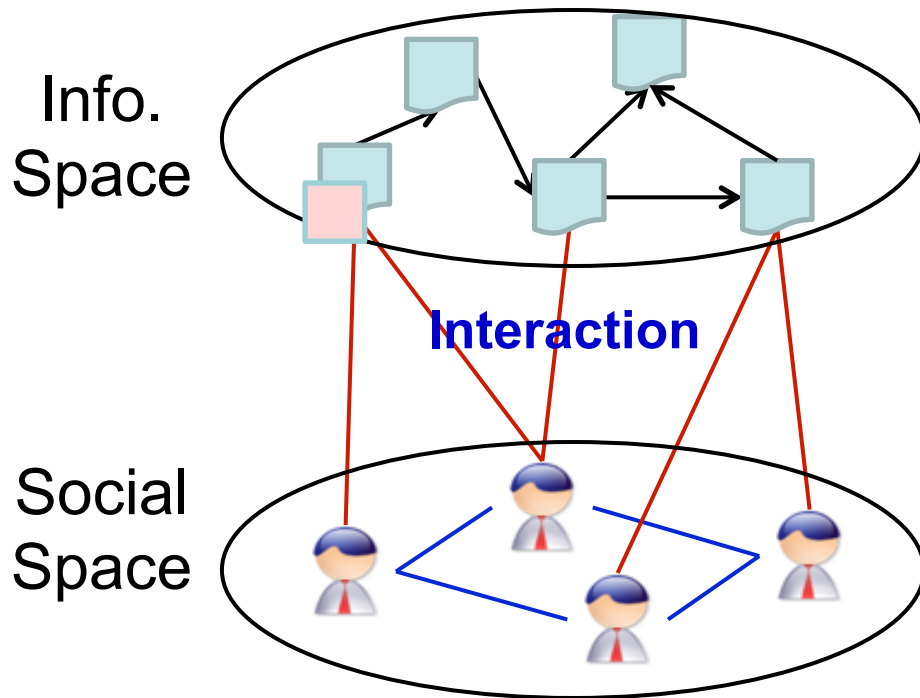
Challenge: Big Social Data

- We generate 2.5×10^{18} byte *big data* per day.
- Big social data:
 - 90% of the data was generated in the past 2 yrs
 - Mining in single data center → mining deep knowledge from multiple data sources

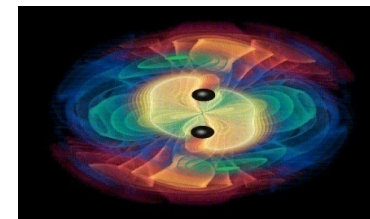
Social Networks



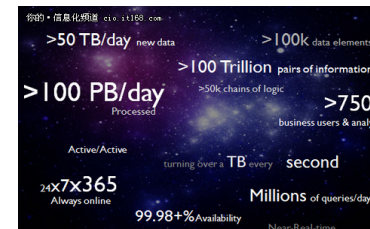
Info. Space vs. Social Space



Opinion Mining



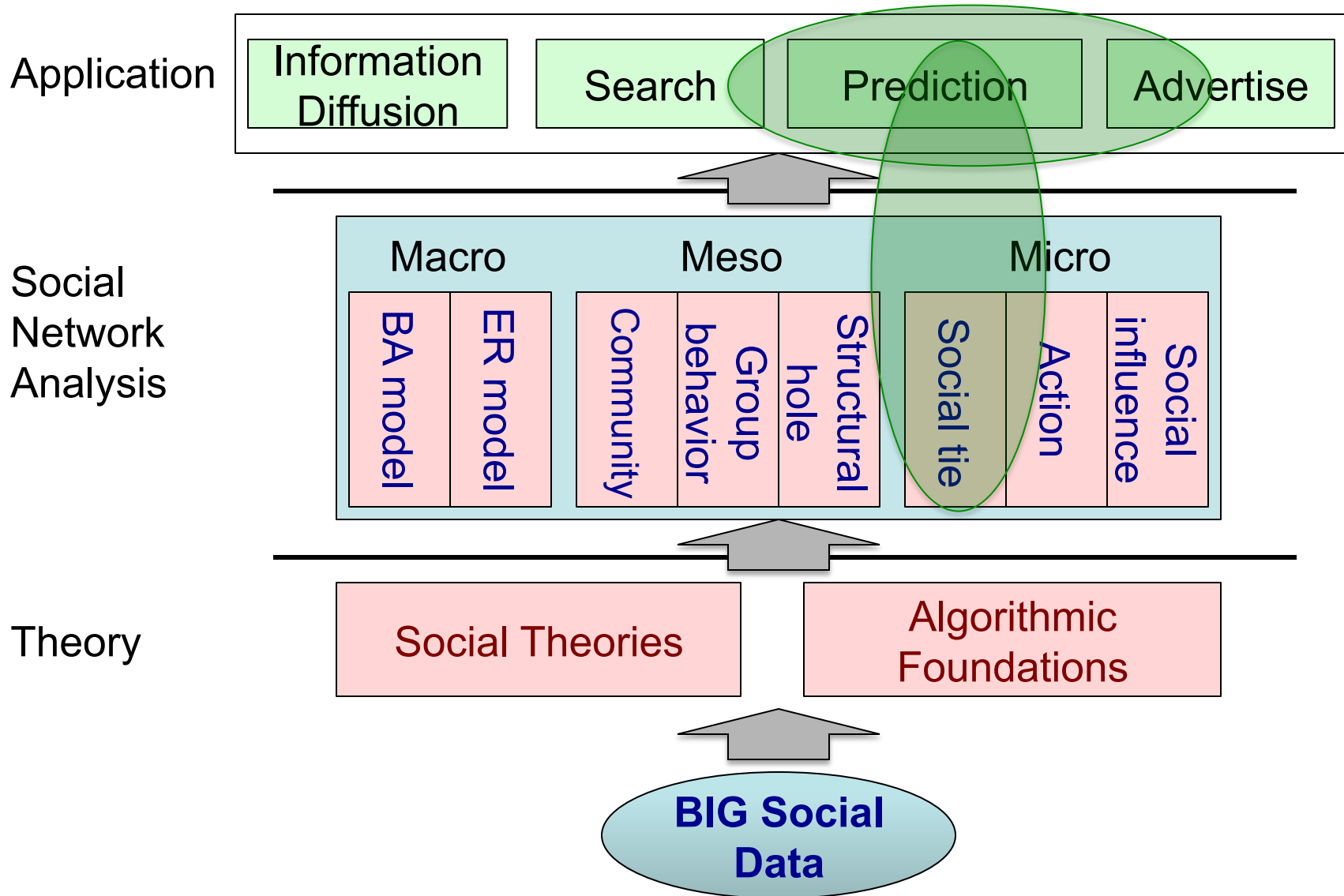
Innovation diffusion



Business intelligence

Understanding the mechanisms of interaction dynamics

Core Research in Social Network



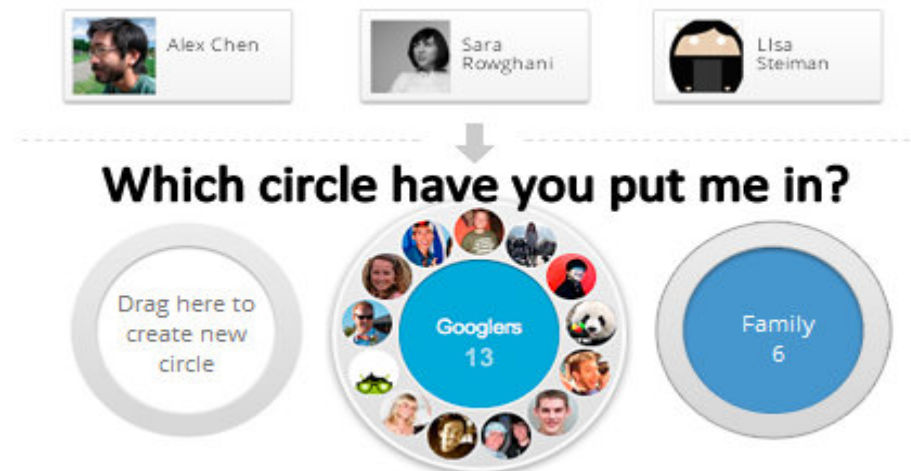
Part A:

Let us start with a simple case “inferring social ties in single network”

(KDD 2010, PKDD 2011 Best Runnerup)

Real social networks are complex...

- Nobody exists merely in one social network.
 - Public network vs. private network
 - Business network vs. family network
- However, existing networks (e.g., Facebook and Twitter) are trying to lump everyone into one big network
 - FB/QQ tries to solve this problem via **lists/groups**
 - however...
- Google circles



Even complex than we imaged!



- Only 16% of mobile phone users in Europe have created custom contact groups
 - *users do not* take the time to create it
 - *users do not* know how to circle their friends
- The Problem is that online social network are black **white**...

Example 2. From BW to Color

(PKDD'11, Best Paper Runnerup)



Enterprise email network

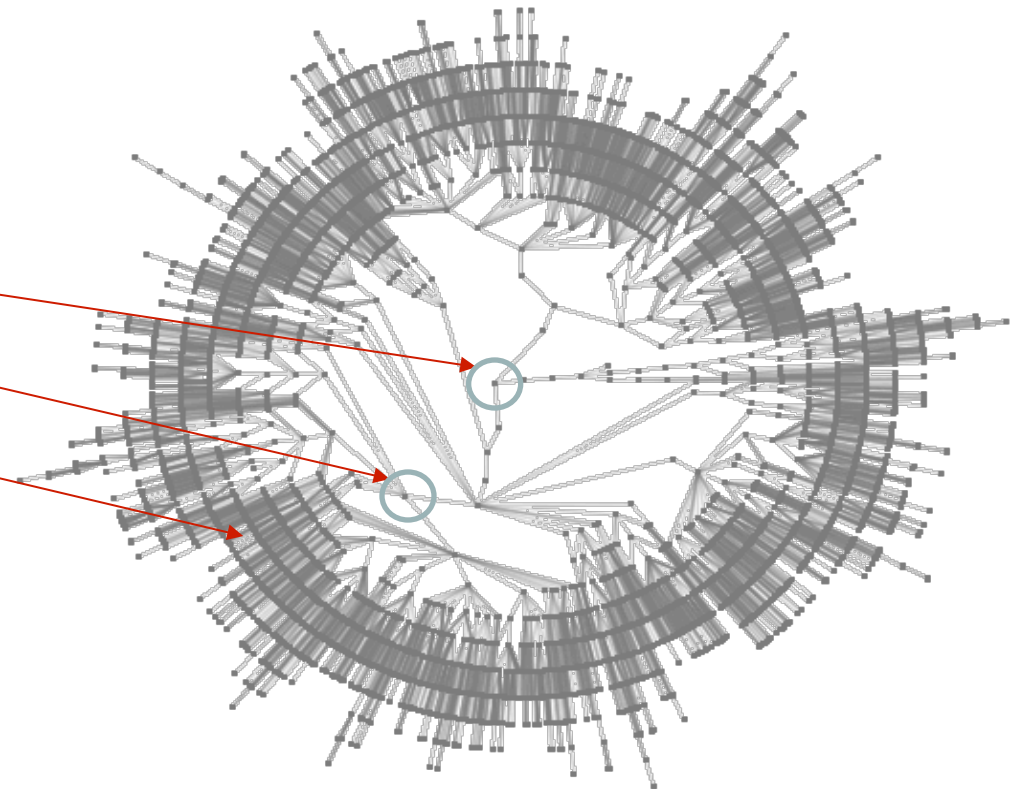
How to
infer



CEO

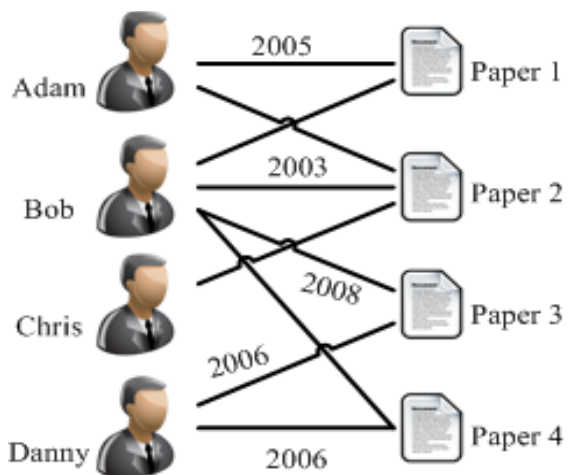
Manager

Employee

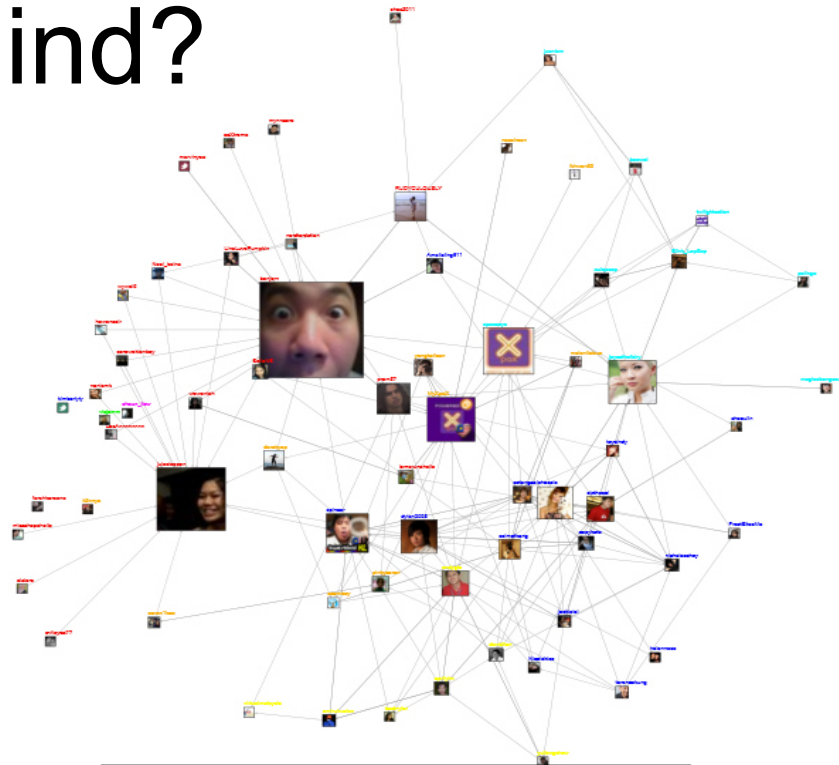


User interactions may form *implicit groups*

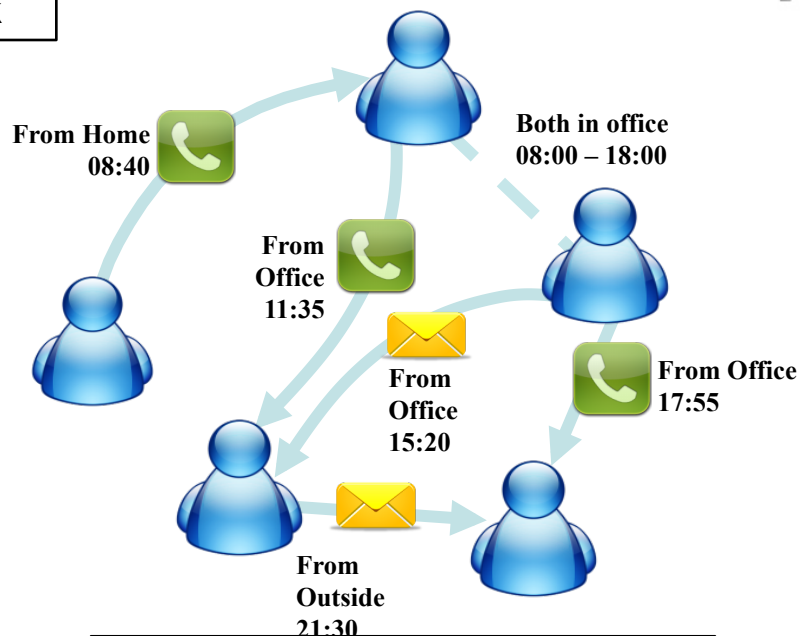
What is behind?



Publication network

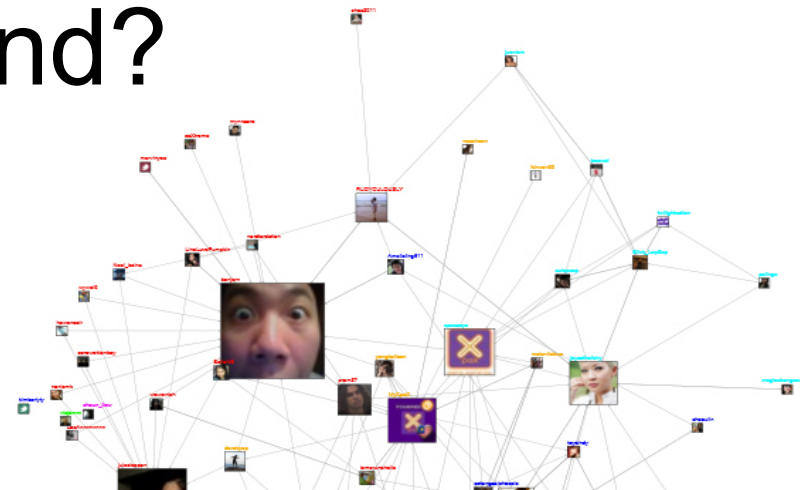
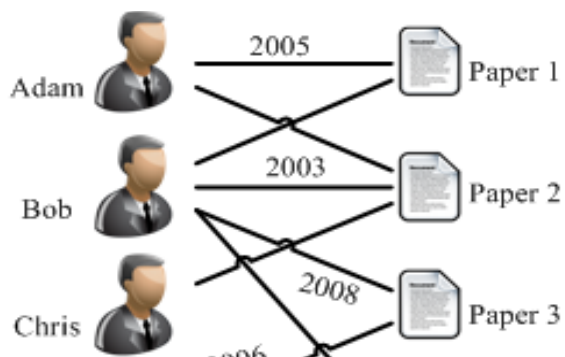


Twitter's following network



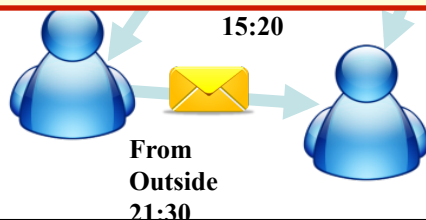
Mobile communication network

What is behind?



Questions:

- What are the **fundamental forces** behind?
- A **generalized framework** for inferring social ties?
- How to **connect** the different networks?



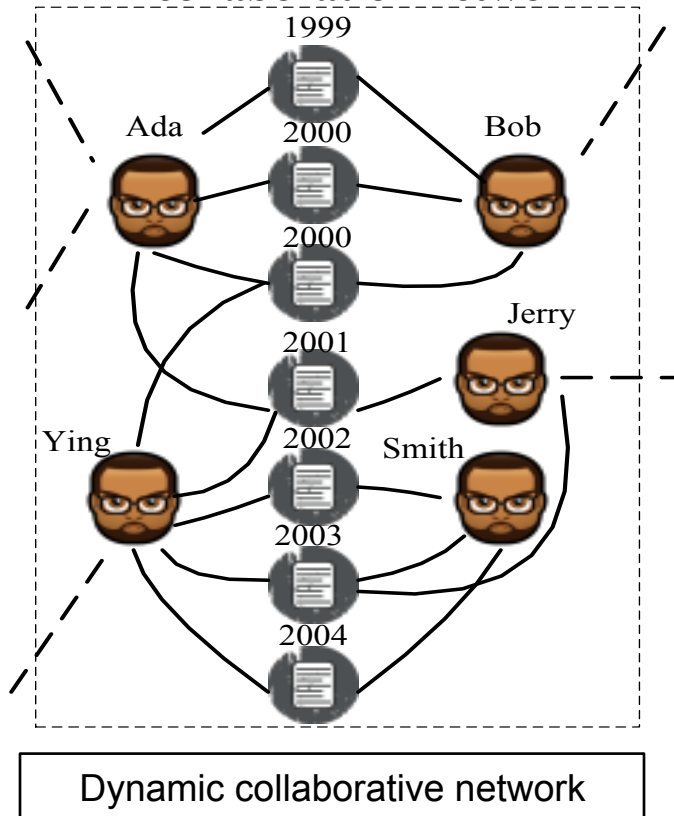
Mobile communication network

inferring social ties in single network

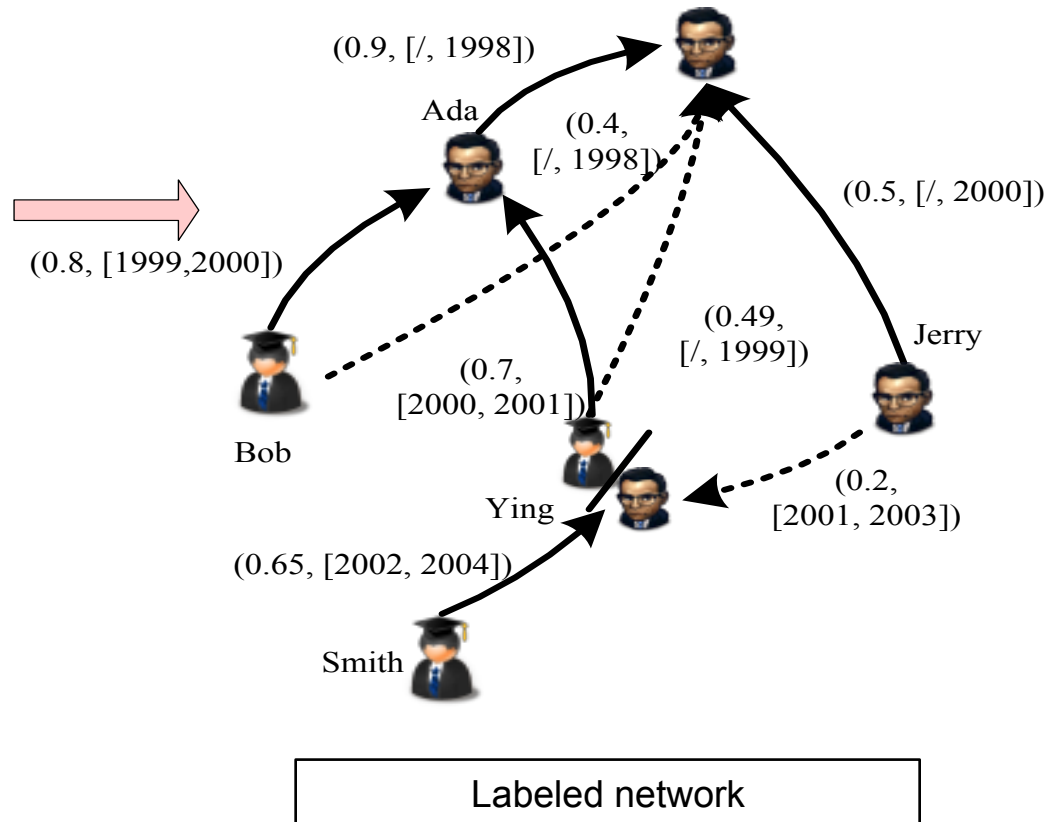
Learning Framework

Problem Analysis

Input: Temporal collaboration network

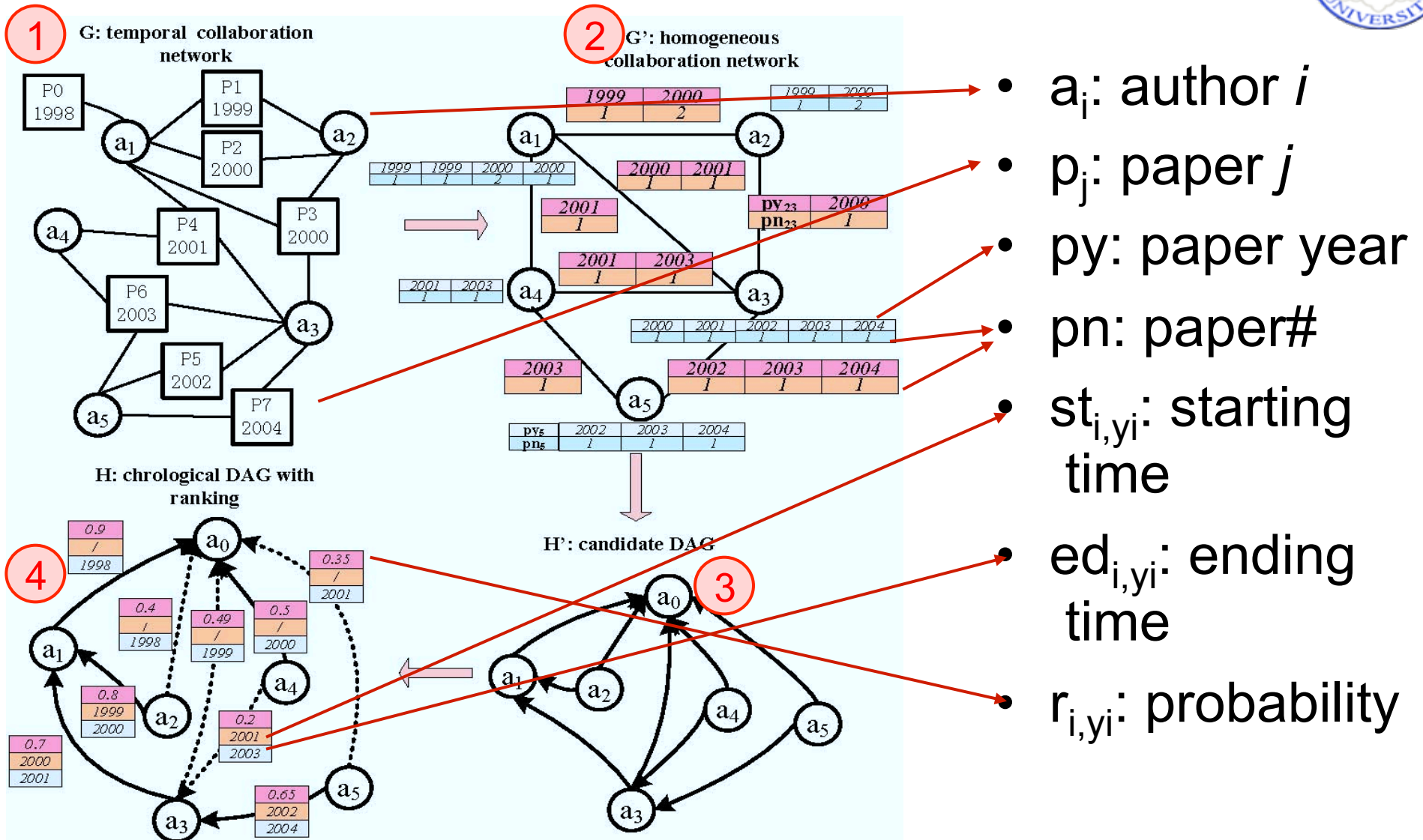


Output: Relationship analysis



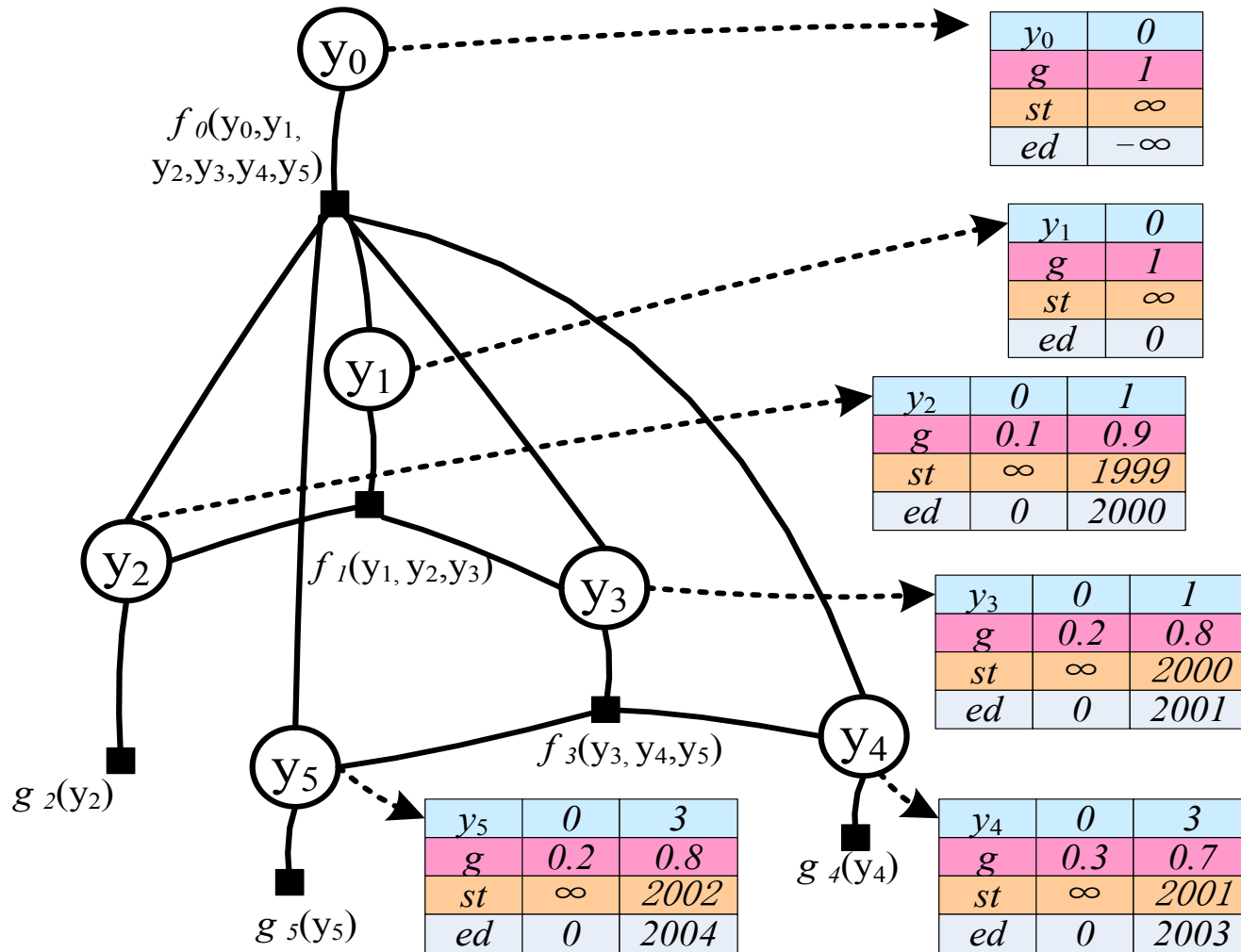
Output: potential types of relationships and their probabilities:
 (type, prob, [s_time, e_time])

Overall Framework



The problem is cast as, for each node, identifying which neighbor has the **highest probability** to be his/her advisor, i.e., $P(y_i=j | x_j, x_{\sim i}, \mathbf{y})$, where x_j and x_i are neighbors.

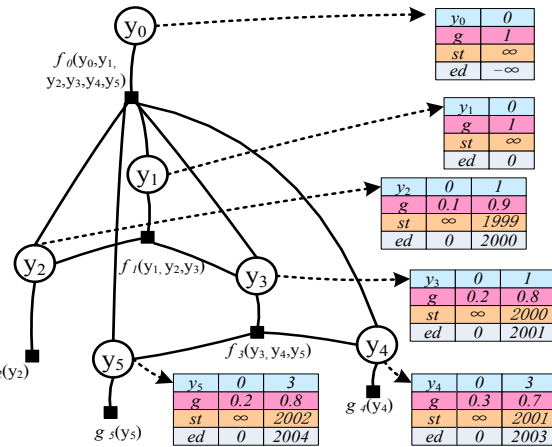
Time-constrained Probabilistic Factor Graph (TPFG)



- Hidden variable y_x : a_x 's advisor
- $st_{x,yx}$: starting time
 $ed_{x,yx}$: ending time
- $g(y_x, st_x, ed_x)$ is pairwise local feature
- $f_x(y_x, Z_x) = \max_{g(y_x, st_x, ed_x)}$ under time constraint
- Y_x : set of potential advisors of a_x

Maximum likelihood estimation

- A general likelihood objective func can be defined as



$$P(y_1, \dots, y_N) = \frac{1}{Z} \prod_{i=1}^N f_i(y_i | \{y_x | x \in Y_i^{-1}\})$$

with

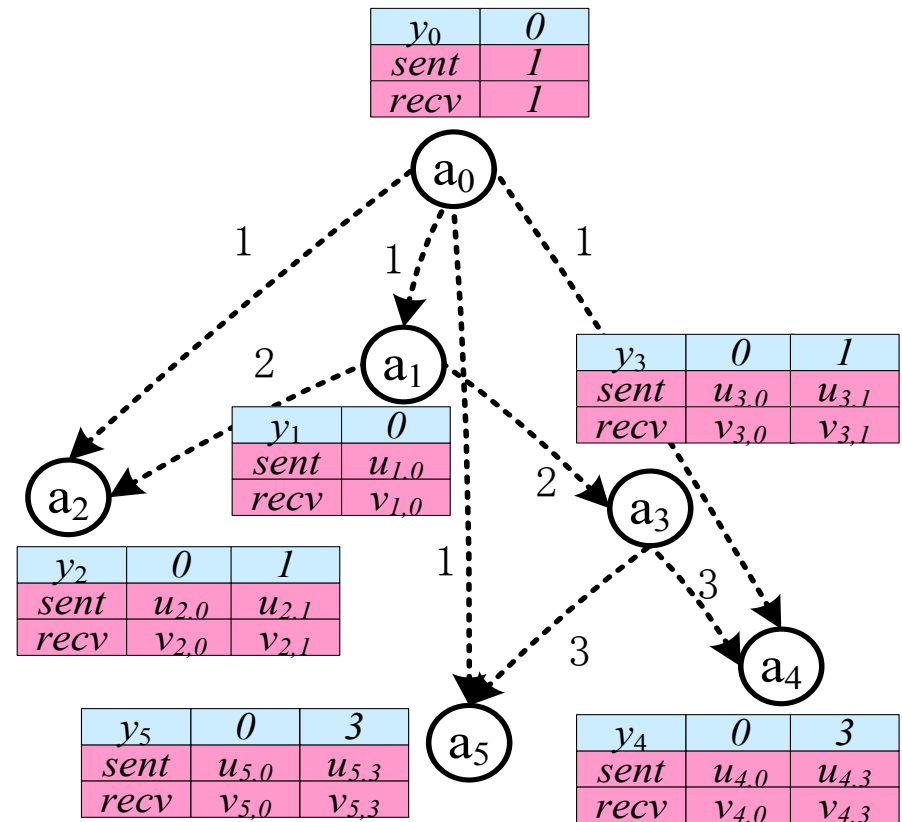
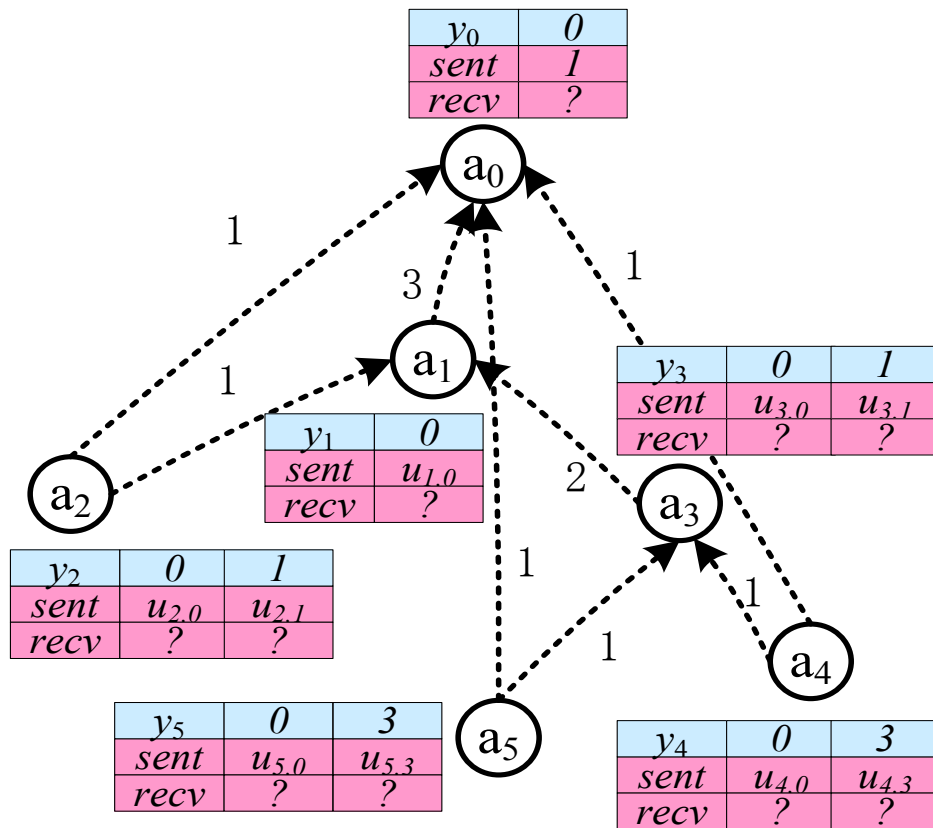
$$f_i(y_i | \{y_x | x \in Y_i^{-1}\}) = g(y_i, st_{ij}, ed_{ij}) \prod_{x \in Y_i^{-1}} \phi(y_x, ed_{ij}, st_{xi})$$

where $\Phi(\cdot)$ can be instantiated in different ways, e.g.,

$$\phi(y_x, ed_{ij}, st_{xi}) = \begin{cases} 1, & y_x \neq i \vee ed_{ij} < st_{xi} \\ 0, & y_x = i \wedge ed_{ij} \geq st_{xi} \end{cases}$$

Inference algorithm of TPFPG

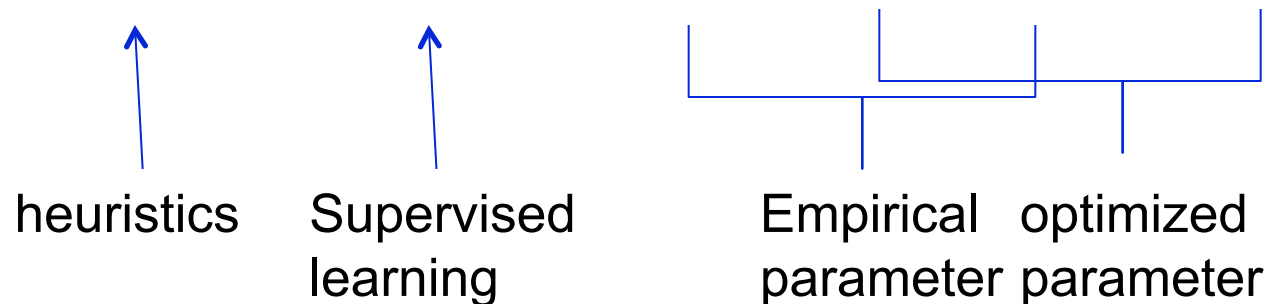
- $r_{ij} = \max P(y_1, \dots, y_{na} | y_i = j) = \exp(\text{sent}_{ij} + \text{recv}_{ij})$



Results of Model 1

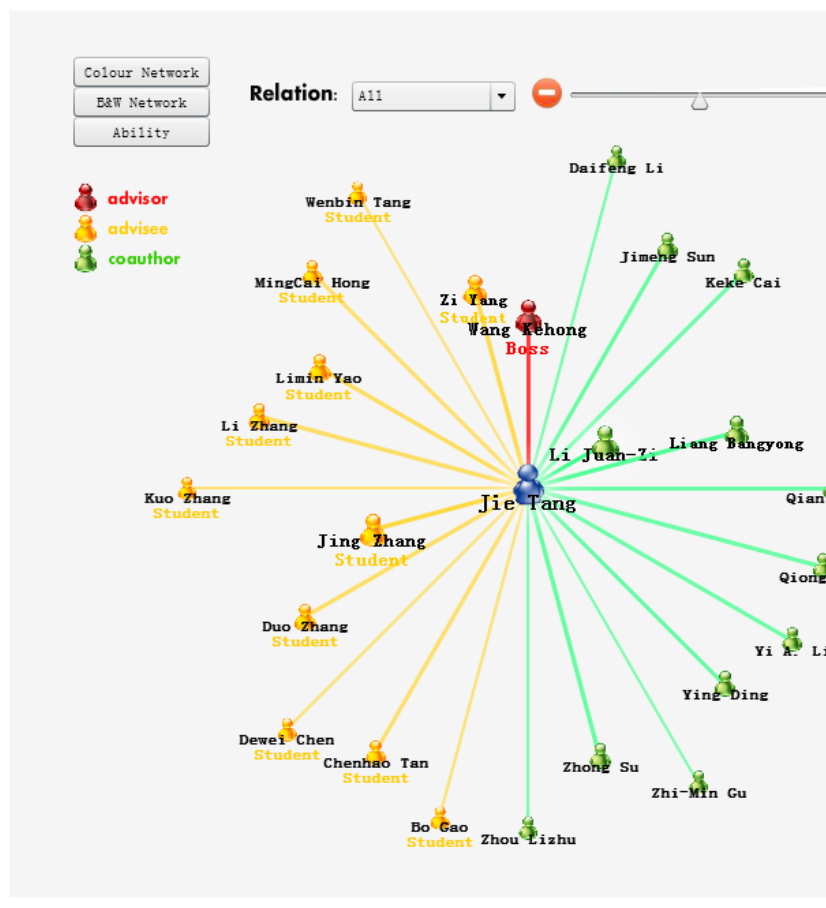
- DBLP data: 654, 628 authors, 1,076,946 publications, years provided.
- Ground truth: MathGenealogy Project; AI Genealogy Project; Faculty Homepage

Datasets	RULE	SVM	IndMAX		Model 1	
TEST1	69.9%	73.4%	75.2%	78.9%	80.2%	84.4%
TEST2	69.8%	74.6%	74.6%	79.0%	81.5%	84.3%
TEST3	80.6%	86.7%	83.1%	90.9%	88.8%	91.3%

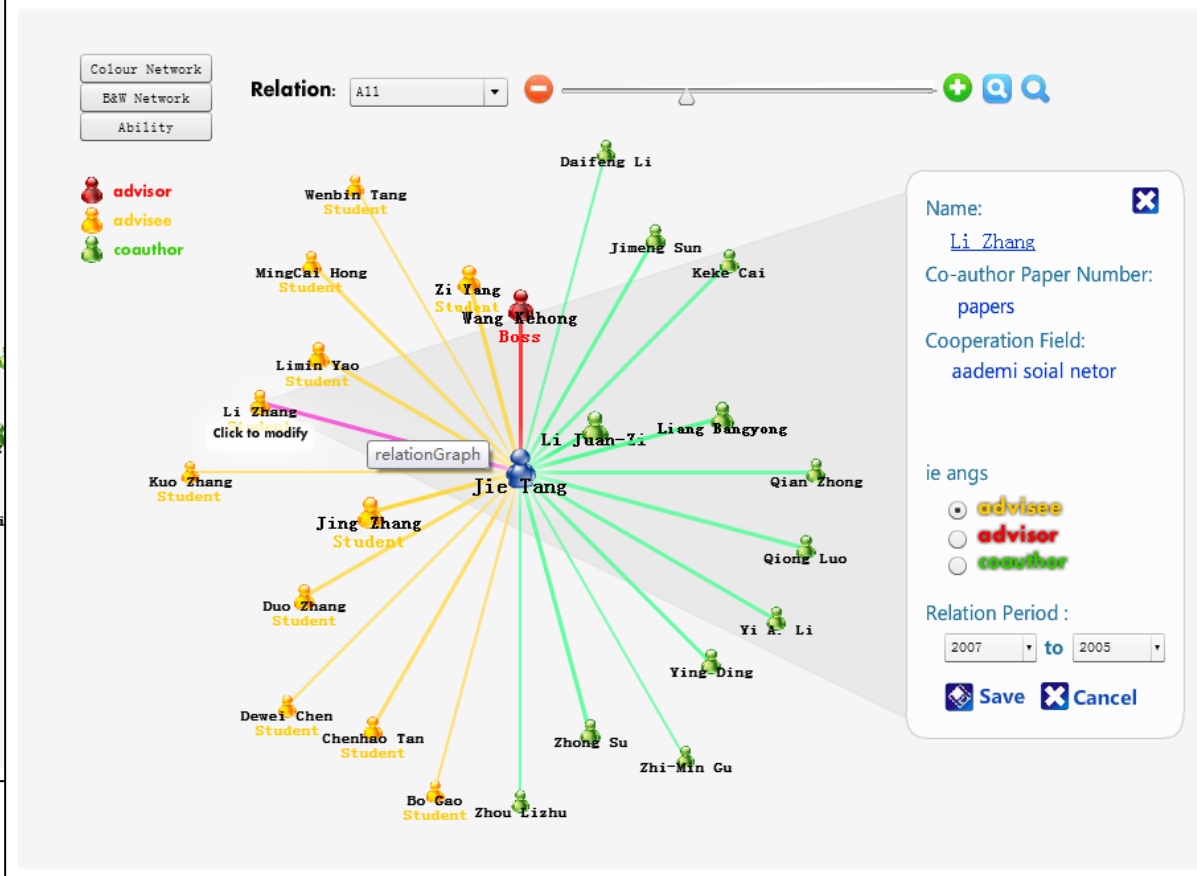


Results

Social Graph



Social Graph



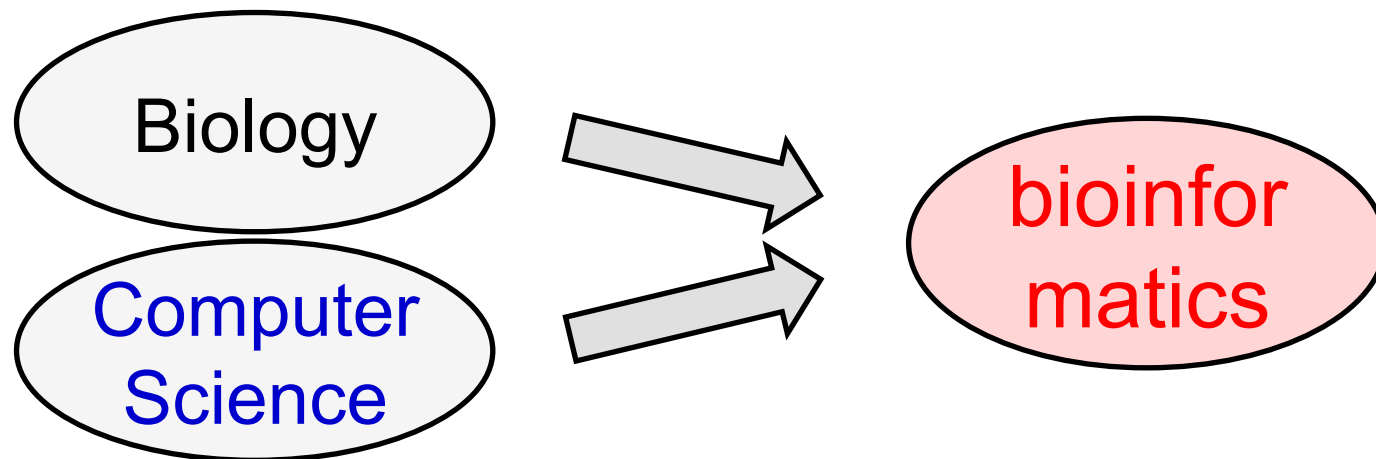
Part B:

Extend the problem to cross-domain
“cross-domain collaboration recommendation”

(KDD 2012, WSDM)

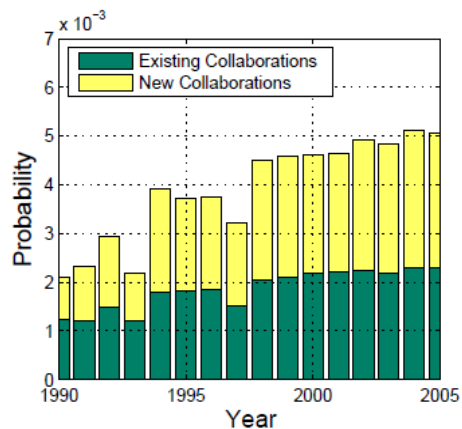
Cross-domain Collaboration

- Interdisciplinary collaborations have generated huge impact, for example,
 - 51 (>1/3) of the KDD 2012 papers are result of cross-domain collaborations between graph theory, visualization, economics, medical inf., DB, NLP, IR
 - Research field evolution

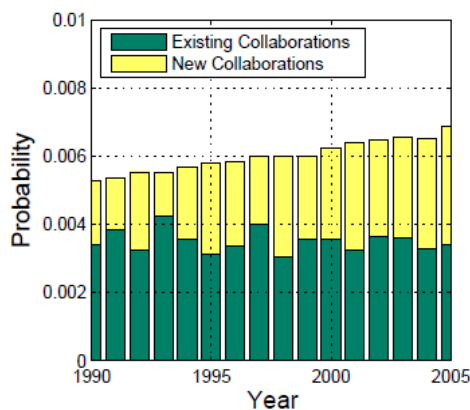


Cross-domain Collaboration (cont.)

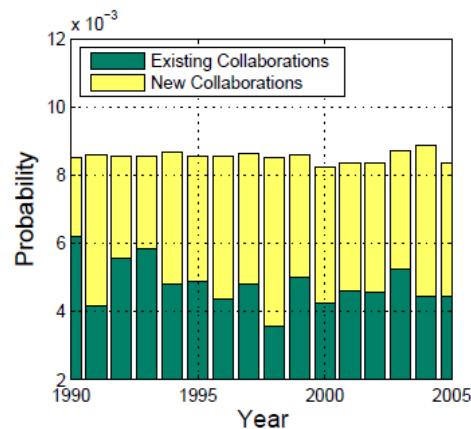
- Increasing trend of cross-domain collaborations



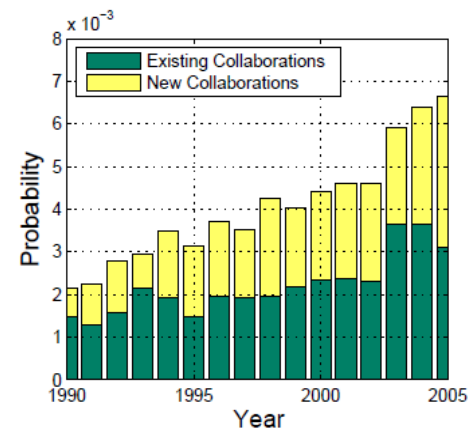
(a) DM - TH



(b) DM - MI



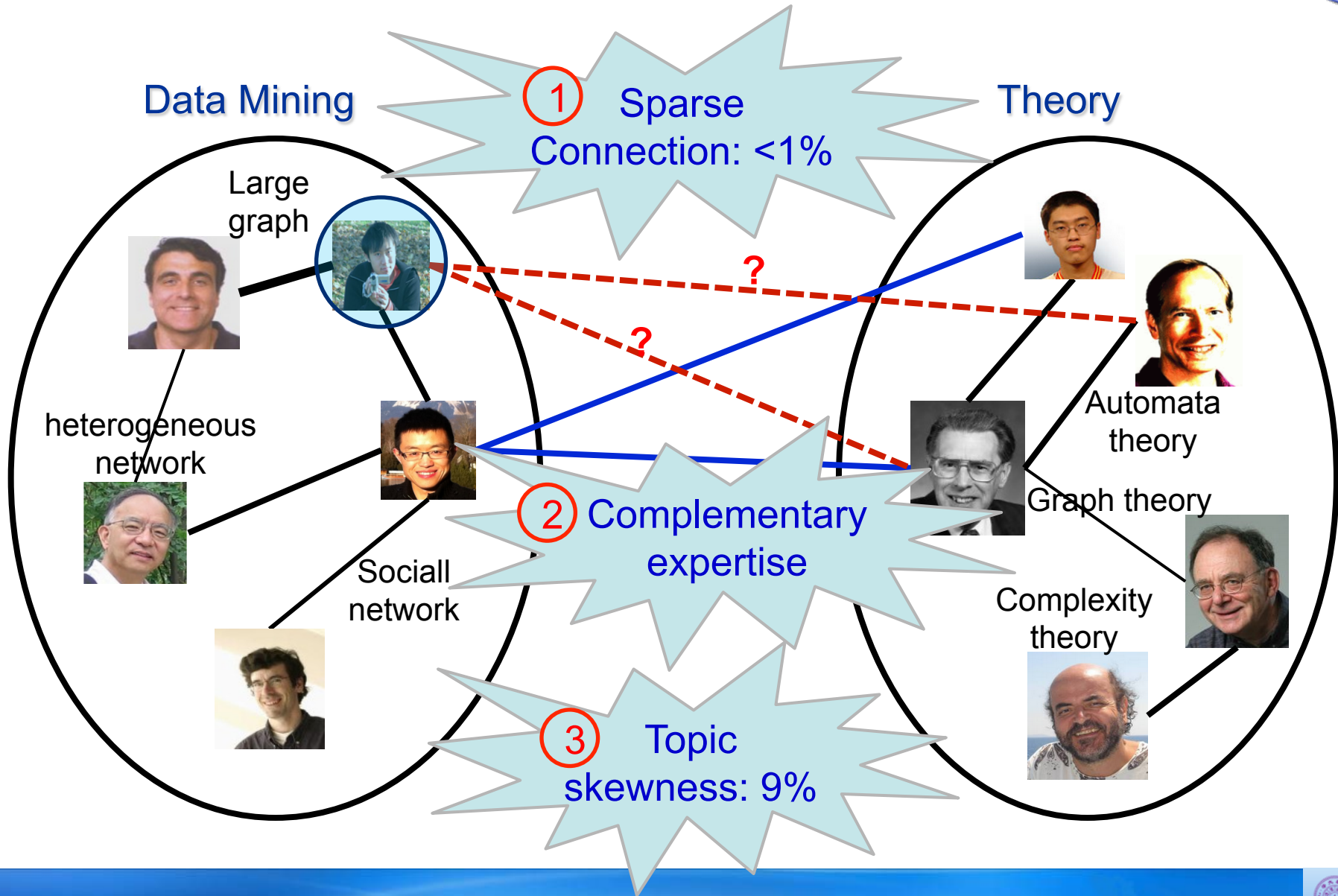
(c) DM - VIS



(d) MI - DB

Data Mining(DM), Medical Informatics(MI), Theory(TH), Visualization(VIS)

Challenges





Related Work-Collaboration recommendation

- Collaborative topic modeling for recommending papers
 - C. Wang and D.M. Blei. [2011]
- On social networks and collaborative recommendation
 - I. Konstas, V. Stathopoulos, and J. M. Jose. [2009]
- CollabSeer: a search engine for collaboration discovery
 - H.-H. Chen, L. Gou, X. Zhang, and C. L. Giles. [2007]
- Referral web: Combining social networks and collaborative filtering
 - H. Kautz, B. Selman, and M. Shah. [1997]
- Fab: content-based, collaborative recommendation
 - M. Balabanovi and Y. Shoham. [1997]



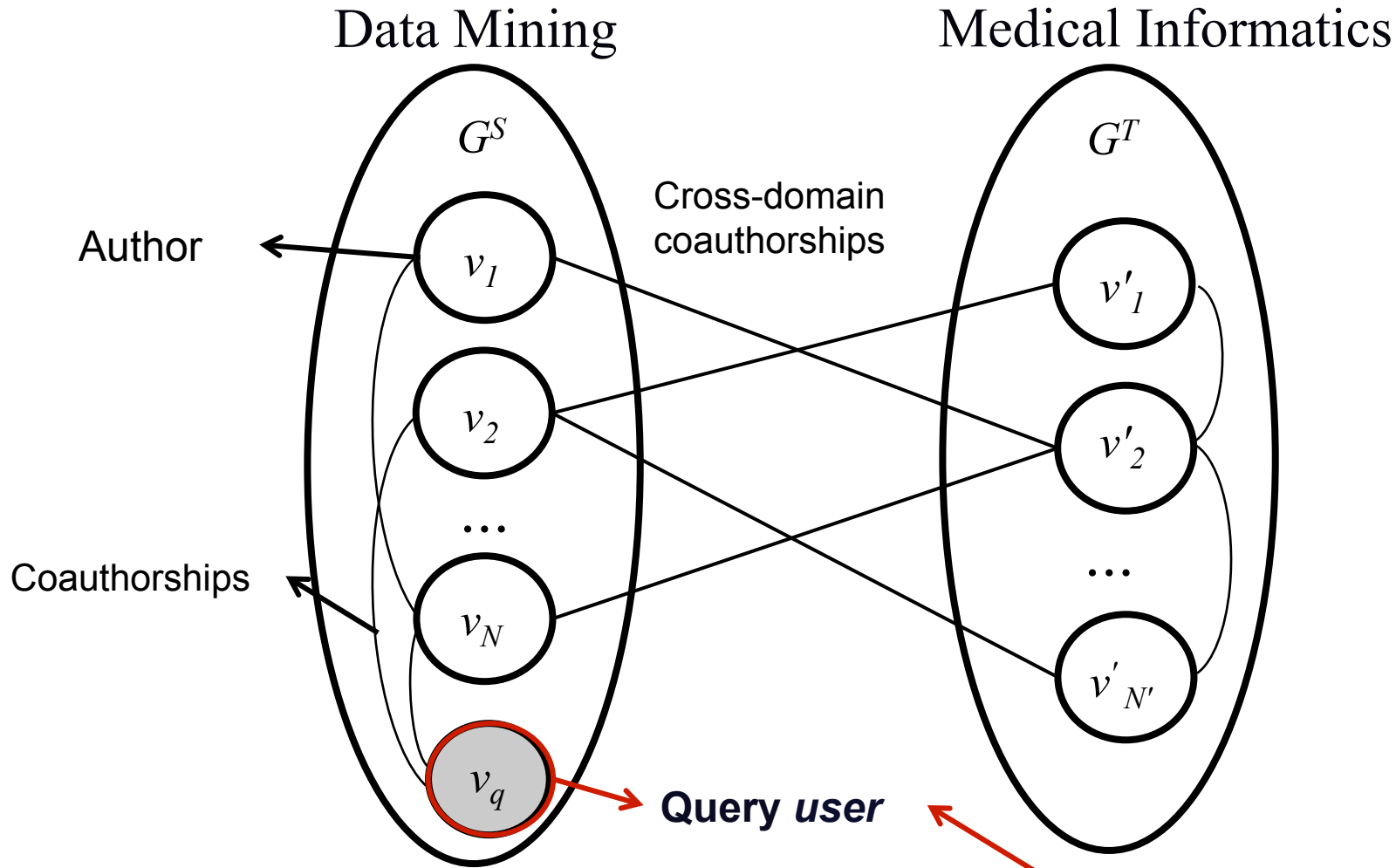
Related Work-Expert finding and matching

- Topic level expertise search over heterogeneous networks
 - J. Tang, J. Zhang, R. Jin, Z. Yang, K. Cai, L. Zhang, and Z. Su. [2011]
- Formal models for expert finding in enterprise corpora
 - K. Balog, L. Azzopardi, and M.de Rijke. [2006]
- Expertise modeling for matching papers with reviewers
 - D. Mimno and A. McCallum. [2007]
- On optimization of expertise matching with various constraints
 - W. Tang, J. Tang, T. Lei, C. Tan, B. Gao, and T. Li. [2012]

cross-domain collaboration recommendation

Approach Framework
—Cross-domain Topic Learning

Author Matching



$$\mathbf{r}^{(t+1)} = (1 - \tau)\mathbf{S} \cdot \mathbf{r}^{(t)} + \tau \mathbf{q}$$

Recall Random Walk

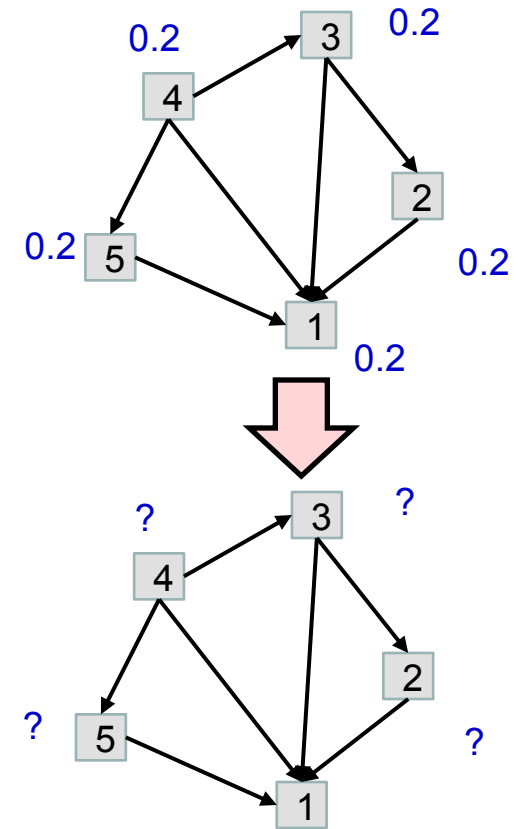
- Let us begin with PageRank^[1]

$$\mathbf{r} = (1 - \alpha)\mathbf{M} \cdot \mathbf{r} + \alpha\mathbf{U}$$

$$M_{ij} = \frac{1}{\text{outdeg}(v_i)}$$

$$U_i = \frac{1}{N}$$

$$\alpha = 0.15$$



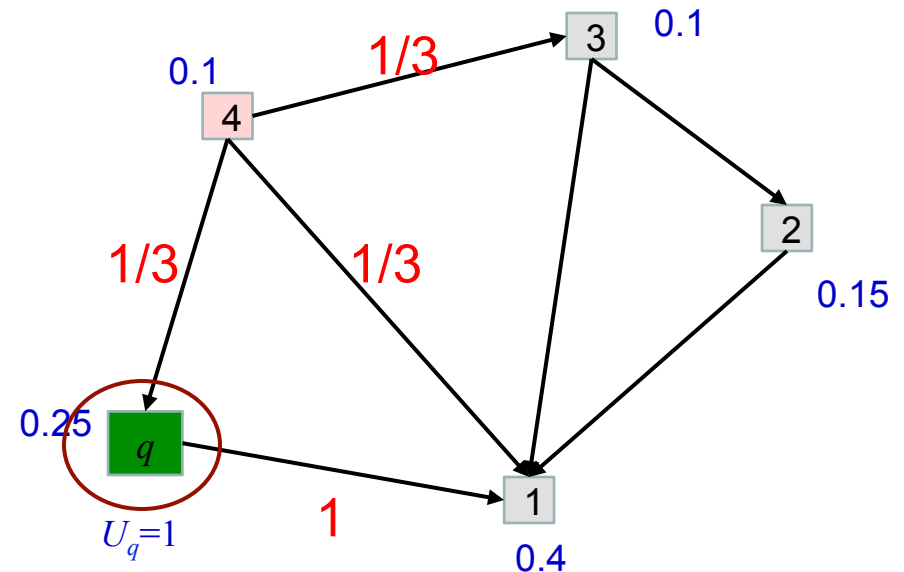
$$(0.2 + 0.2 \cdot 0.5 + 0.2 \cdot 1/3 + 0.2)0.85 + 0.15 \cdot 0.2$$

Random Walk with Restart^[1]

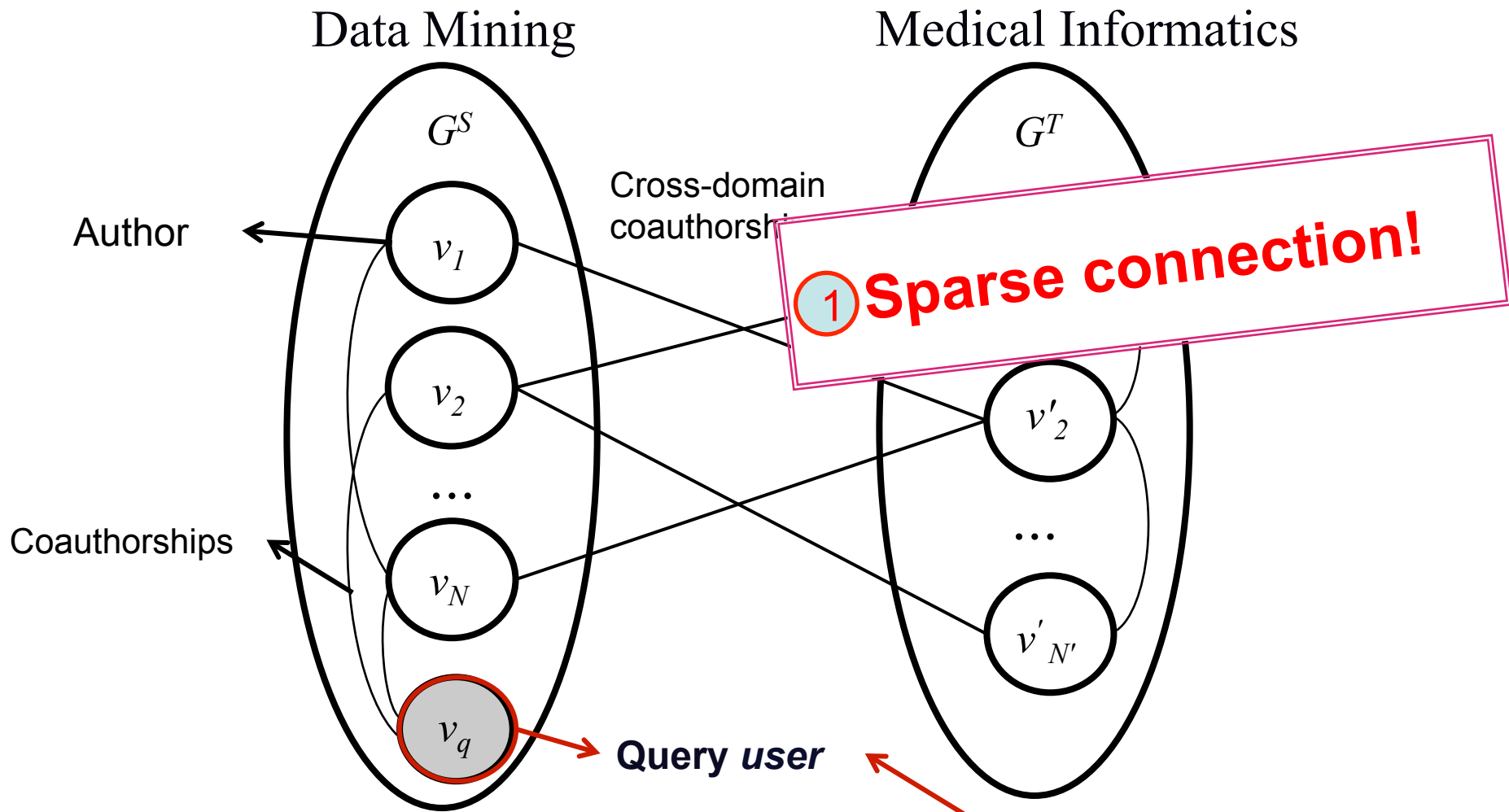
$$\mathbf{r}_q = (1 - \alpha)\mathbf{M} \cdot \mathbf{r}_q + \alpha\mathbf{U}$$

$$M_{ij} = \frac{1}{\text{outdeg}(v_i)}$$

$$U_i = \begin{cases} 1, & i = q \\ 0, & i \neq q \end{cases}$$

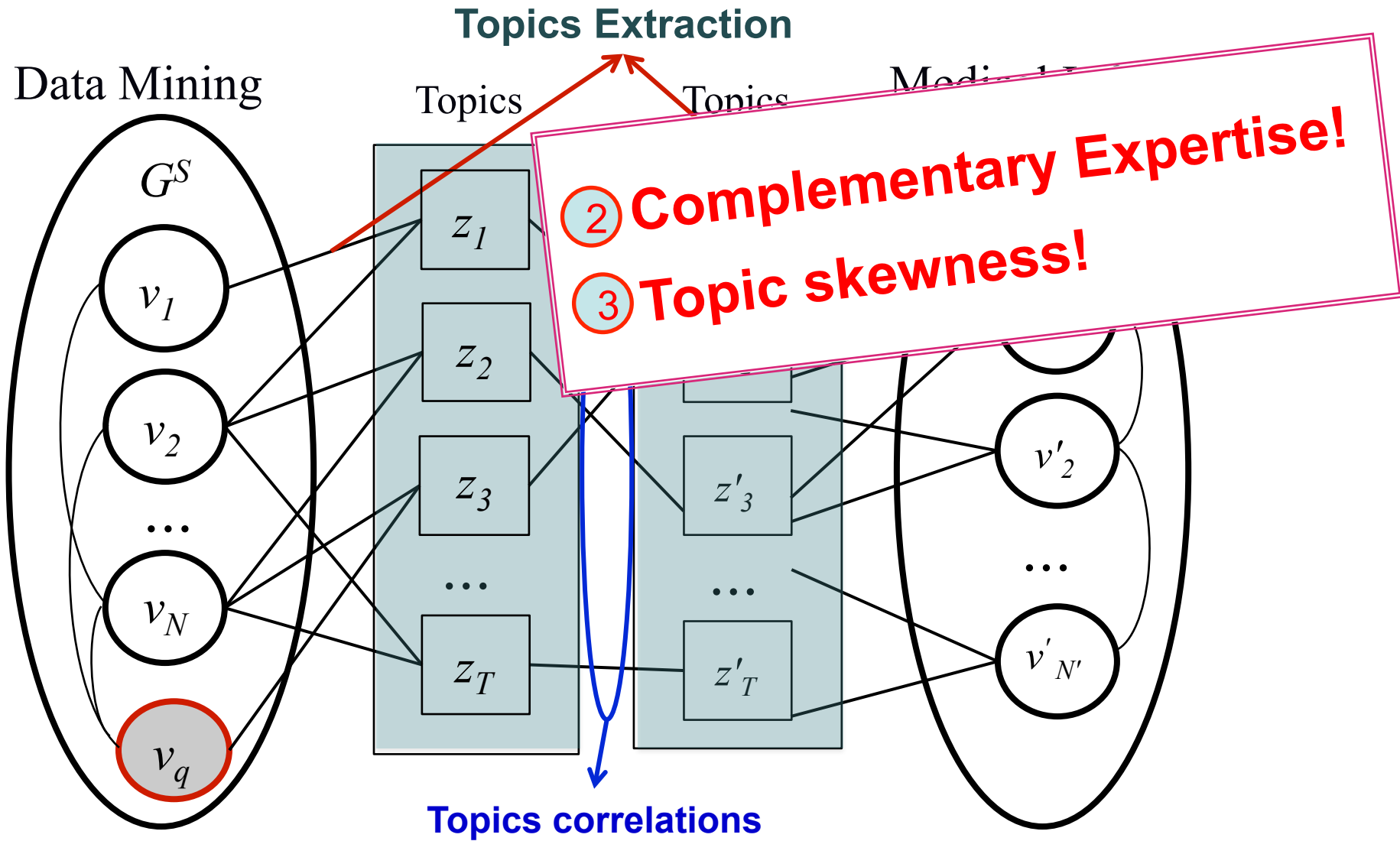


Author Matching



$$\mathbf{r}^{(t+1)} = (1 - \tau)\mathbf{S} \cdot \mathbf{r}^{(t)} + \tau \mathbf{q}$$

Topic Matching



Recall Topic Model



TOPIC 1



TOPIC 2

DOCUMENT 1: money¹ bank¹ bank¹ loan¹ river² stream²

bank¹ money¹
river² stream²
loan¹ bank¹

stream² bank²
river² bank² bank¹
money¹ loan¹
river² stream²
stream² river²

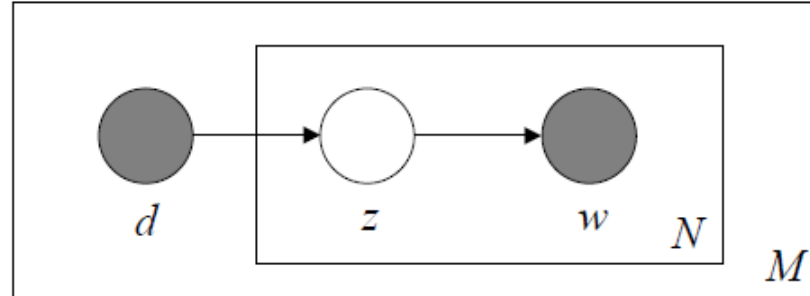
- Usage of a theme:
 - Summarize topics/subtopics
 - Navigate documents
 - Retrieve documents
 - Segment documents
 - All other tasks involving unigram language models

Mixture components

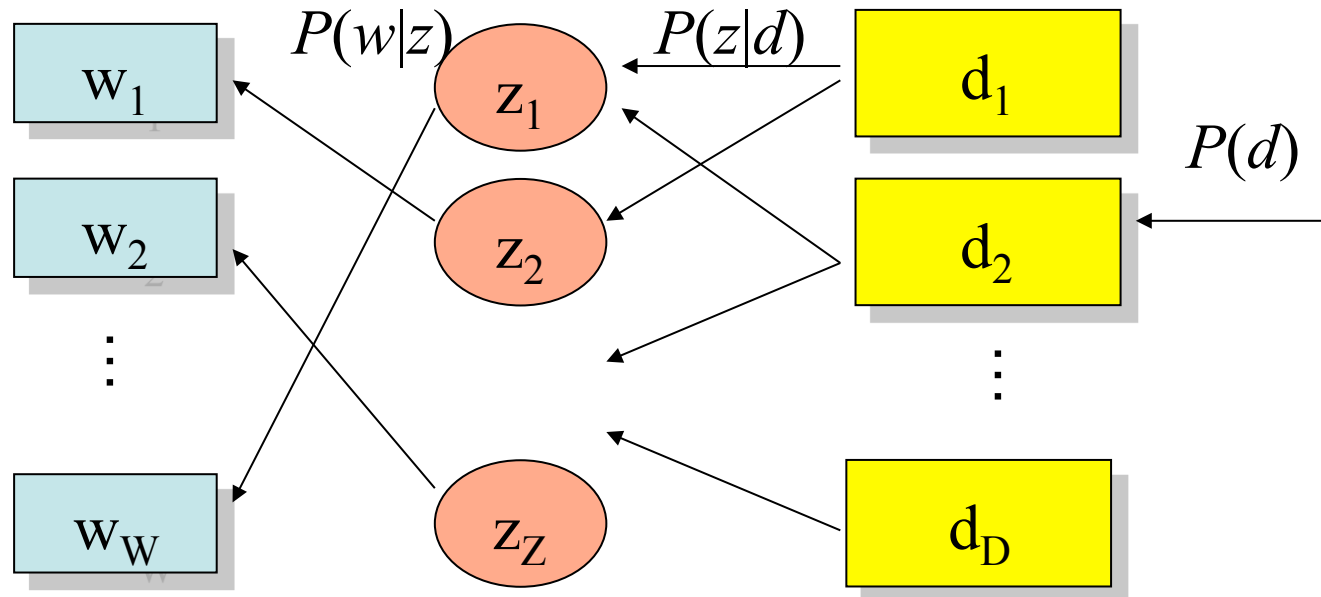
Mixture weights

Bayesian approach: use priors
 Mixture weights $\sim \text{Dirichlet}(\alpha)$
 Mixture components $\sim \text{Dirichlet}(\beta)$

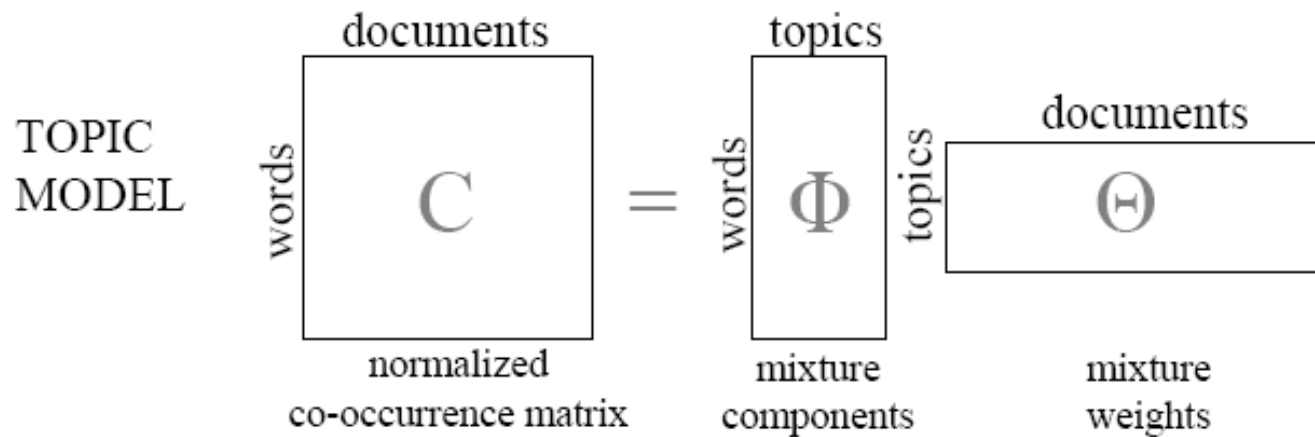
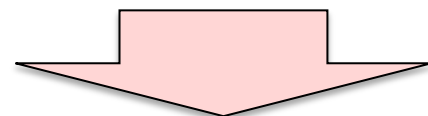
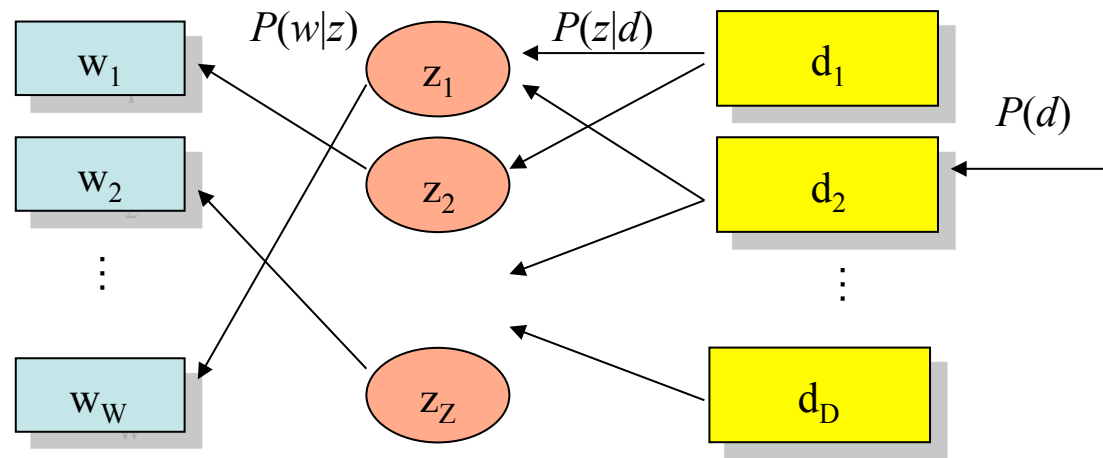
Topic Model



- A generative model for generating the co-occurrence of documents $d \in D = \{d_1, \dots, d_D\}$ and terms $w \in W = \{w_1, \dots, w_W\}$, which associates latent variable $z \in Z = \{z_1, \dots, z_Z\}$.
- The generative processing is:



Topic Model



Maximum-likelihood

- Definition

- We have a density function $P(x|\Theta)$ that is governed by the set of parameters Θ , e.g., P might be a set of Gaussians and Θ could be the means and covariances
- We also have a data set $X=\{x_1, \dots, x_N\}$, supposedly drawn from this distribution P , and assume these data vectors are i.i.d. with P .
- Then the log-likelihood function is:

$$L(\Theta | X) = \log p(X | \Theta) = \log \prod_i p(x_i | \Theta) = \sum_i \log p(x_i | \Theta)$$

- The log-likelihood is thought of as a function of the parameters Θ where the data X is fixed. Our goal is to find the Θ that maximizes L . That is

$$\Theta^* = \arg \max_{\Theta} L(\Theta | X)$$

Topic Model

- Following the likelihood principle, we determine $P(d)$, $P(z|d)$, and $P(w|d)$ by maximization of the likelihood

$P(d)$, $P(z|d)$, and $P(w|d)$

Unobserved data

co-occurrence times of d and w . Which is obtained according to the multi-distribution

$$L(\Theta | d, w, z) = \log \prod_d \prod_w P(d, w)^{n(d, w)}$$

Observed data

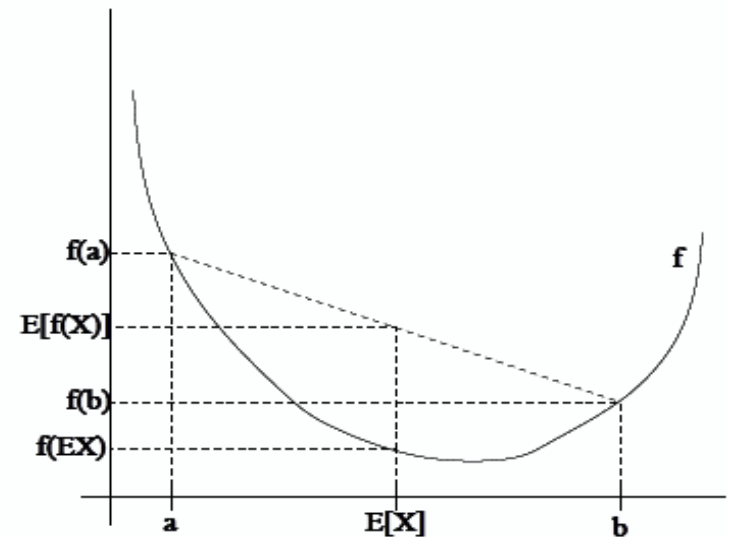
$$= \sum_{d \in D} \sum_{w \in W} n(d, w) \log P(d, w)$$

$$= \sum_{d \in D} \sum_{w \in W} n(d, w) \log \left(\sum_{z \in Z} P(w | z) P(d | z) P(z) \right)$$

Jensen's Inequality

- Recall that f is a **convex function** if $f''(x) \geq 0$, and f is strictly convex function if $f''(x) > 0$
- Let f be a convex function, and let X be a random variable, then:

$$E[f(X)] \geq f(EX)$$



- Moreover, if f is strictly convex, then $E[f(X)] = f(EX)$ holds true if and only if $X = E[X]$ with probability 1 (i.e., if **X is a constant**)

Basic EM Algorithm

- However, Optimizing the likelihood function is analytically intractable but when the likelihood function can be simplified by assuming the existence of and values for additional but missing (or hidden) parameters:

$$L(\Theta | X) = \sum_i \log p(x_i | \Theta) = \sum_i \log \sum_z p(x_i, z | \Theta)$$

- Maximizing $L(\Theta)$ explicitly might be difficult, and the strategy is to instead repeatedly construct a lower-bound on L (E-step), and then optimize that lower bound (M-step).
 - For each i , let Q_i be some distribution over z ($\sum_z Q_i(z)=1$, $Q_i(z) \geq 0$), then

$$\sum_i \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \Theta) = \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \Theta)}{Q_i(z^{(i)})} \geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \Theta)}{Q_i(z^{(i)})}$$

- The above derivation used Jensen's inequality. Specifically, $f(x) = \log x$ is a concave function, since $f''(x) = -1/x^2 < 0$

Parameter Estimation-Using EM

- According to Basic EM:

$$Q_i(z^{(i)}) = p(z^{(i)} | x^{(i)}; \Theta)$$

- Then we define

$$Q_i(z^{(i)}) = p(z | d, w)$$

- Thus according to Jensen's inequality

$$\begin{aligned} L(\Theta) &= \sum_{d \in D} \sum_{w \in W} n(d, w) \log \sum_{z \in Z} p(z | d, w) \frac{p(w | z) p(d | z) p(z)}{p(z | d, w)} \\ &\geq \sum_{d \in D} \sum_{w \in W} n(d, w) \sum_{z \in Z} p(z | d, w) \log \frac{p(w | z) p(d | z) p(z)}{p(z | d, w)} \end{aligned}$$

(1) Solve $P(w|z)$

- We introduce Lagrange multiplier λ with the constraint that $\sum_w P(w|z)=1$, and solve the following equation:

$$\frac{\partial}{\partial P(w|z)} \left\{ \sum_{d \in D} \sum_{w \in W} n(d, w) \sum_{z \in Z} p(z|d, w) \log \frac{p(w|z)p(d|z)p(z)}{p(z|d, w)} + \lambda \left[\sum_z P(w|z) - 1 \right] \right\} = 0$$

$$\therefore \frac{\sum_{d \in D} n(d, w) P(z|d, w)}{P(w|z)} + \lambda = 0,$$

$$\therefore P(w|z) = - \frac{\sum_{d \in D} n(d, w) P(z|d, w)}{\lambda},$$

$$\sum_w P(w|z) = 1,$$

$$\therefore \lambda = - \sum_{w \in W} \sum_{d \in D} n(d, w) P(z|d, w),$$

$$\therefore P(w|z) = \frac{\sum_{d \in D} n(d, w) P(z|d, w)}{\sum_{w \in W} \sum_{d \in D} n(d, w) P(z|d, w)}$$

The final update Equations

- E-step:

$$P(z | d, w) = \frac{P(w | z)P(d | z)P(z)}{\sum_{z \in Z} P(w | z)P(d | z)P(z)}$$

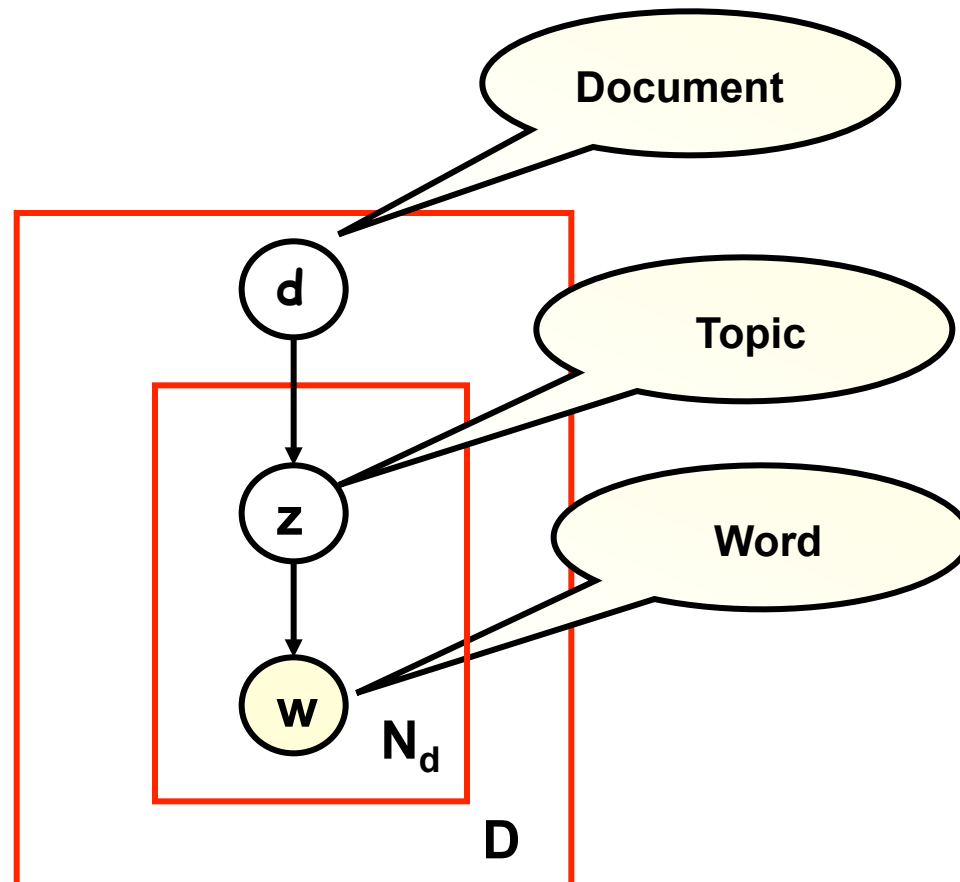
- M-step:

$$P(w | z) = \frac{\sum_{d \in D} n(d, w)P(z | d, w)}{\sum_{w \in W} \sum_{d \in D} n(d, w)P(z | d, w)}$$

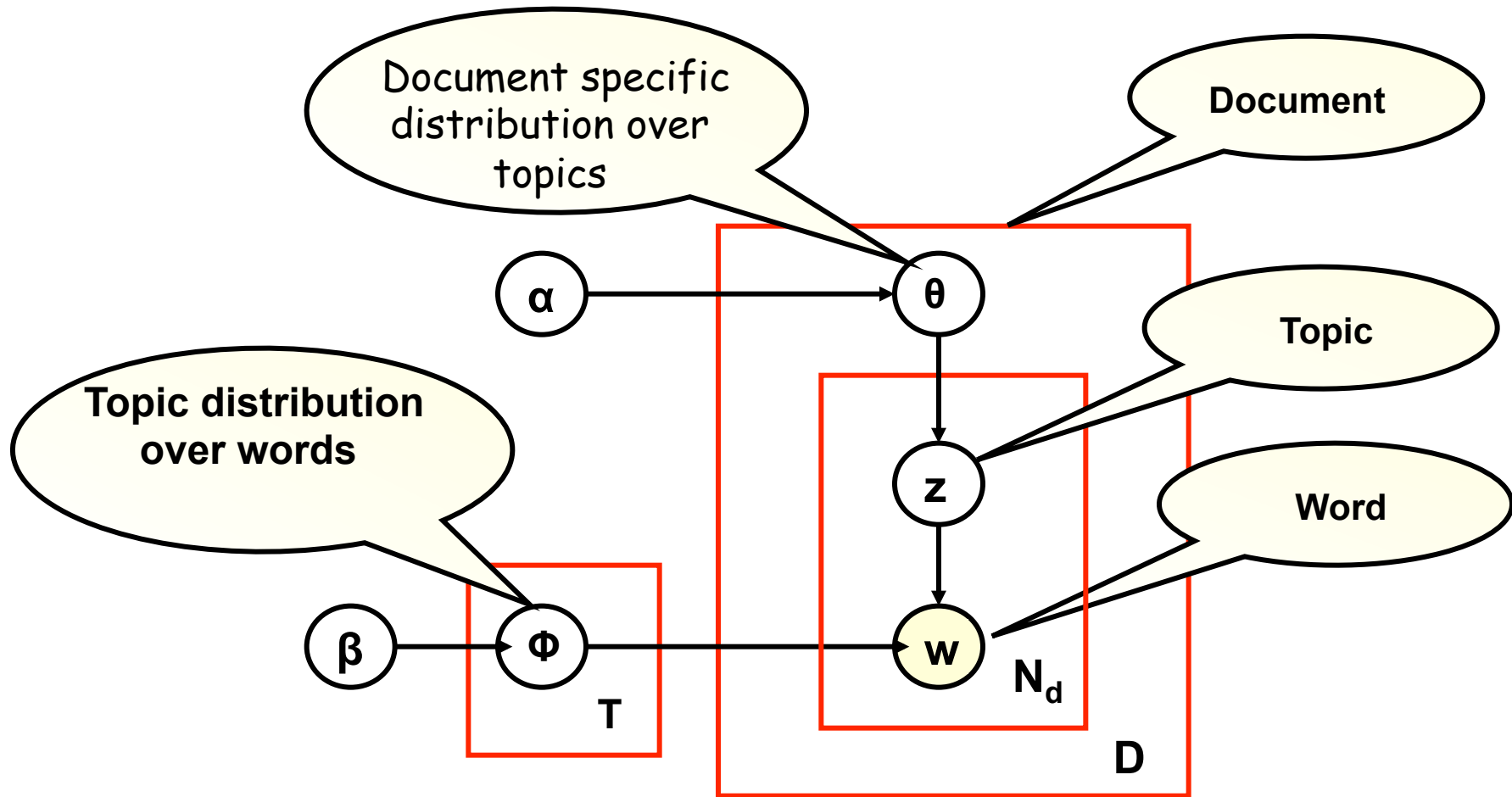
$$P(d | z) = \frac{\sum_{w \in W} n(d, w)P(z | d, w)}{\sum_{d \in D} \sum_{w \in W} n(d, w)P(z | d, w)}$$

$$P(z) = \frac{\sum_{d \in D} \sum_{w \in W} n(d, w)P(z | d, w)}{\sum_{w \in W} \sum_{d \in D} n(d, w)}$$

PLSI(SIGIR'99)

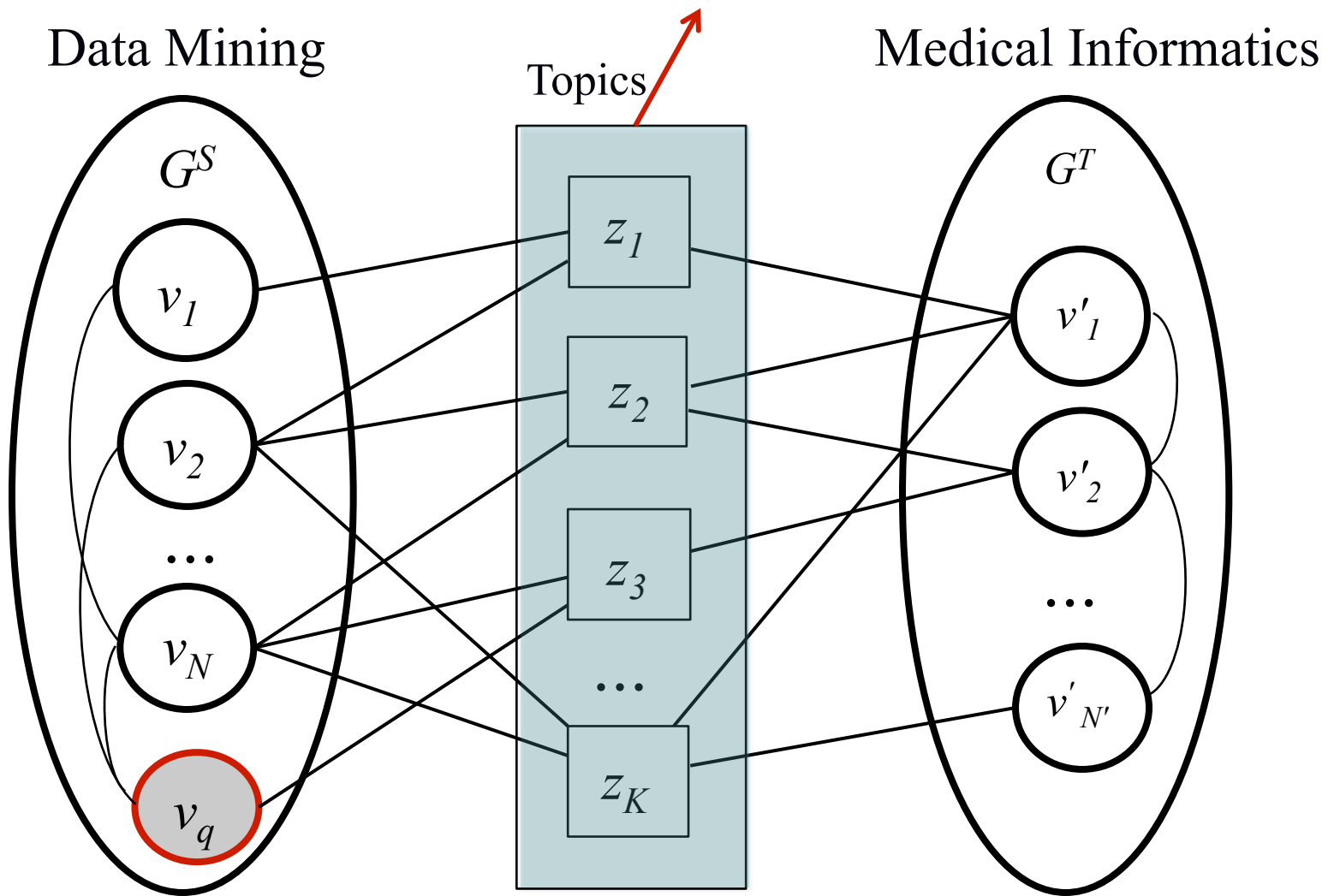


LDA (JMLR'03)



Cross-domain Topic Learning

Identify “cross-domain” Topics



Collaboration Topics Extraction



Step 1:

Input: a source domain G^S and a target domain G^T

Output: estimated parameters $\theta, \theta', \phi, \vartheta$, and λ

Initialize an ACT model in G^S by learning from documents written by authors only from G^S ;

Similarly, initialize an ACT model for target domain G^T ;

foreach collaborated document d **do**

foreach word $x_{di} \in d$ **do**

Toss a coin s_{di} according to $bernoulli(s_{di}) \sim beta(\gamma_t, \gamma)$, where $beta(\cdot)$ is a Beta distribution, and γ_t and γ are two parameters;

if $s_{di} = 0$ **then**

Randomly select a pair (v, v') from d 's authors, where v is an author from G^S and v' from G^T ;

Draw a topic $z_{di} \sim multi(\vartheta_{vv'})$ from the topic mixture $\vartheta_{vv'}$ specific to (v, v') ;

end

if $s_{di} = 1$ **then**

Randomly select a user v ;

Draw a topic $z_{di} \sim multi(\theta_v)$ from the topic model of user v ;

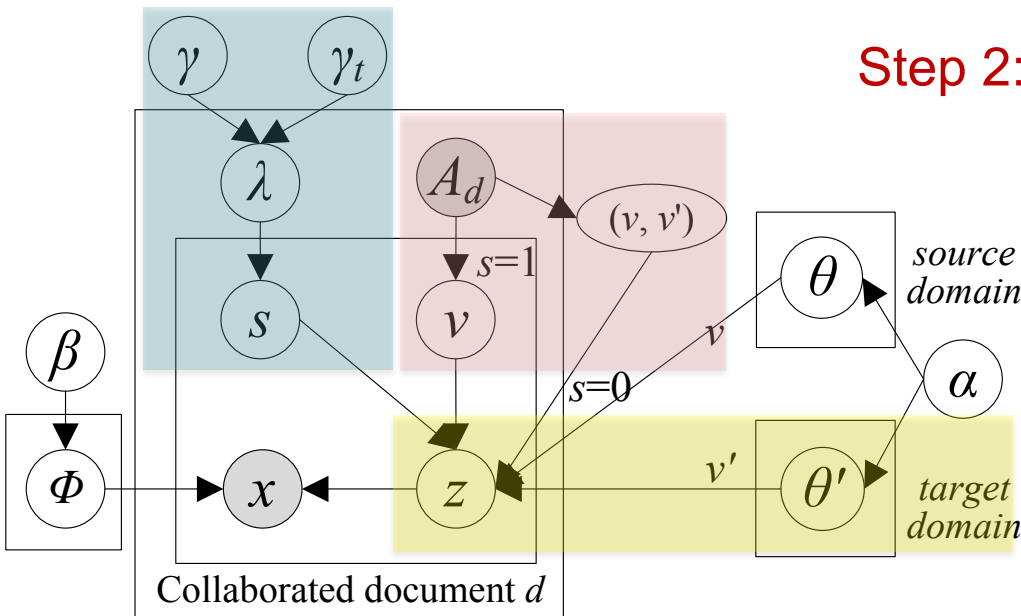
end

end

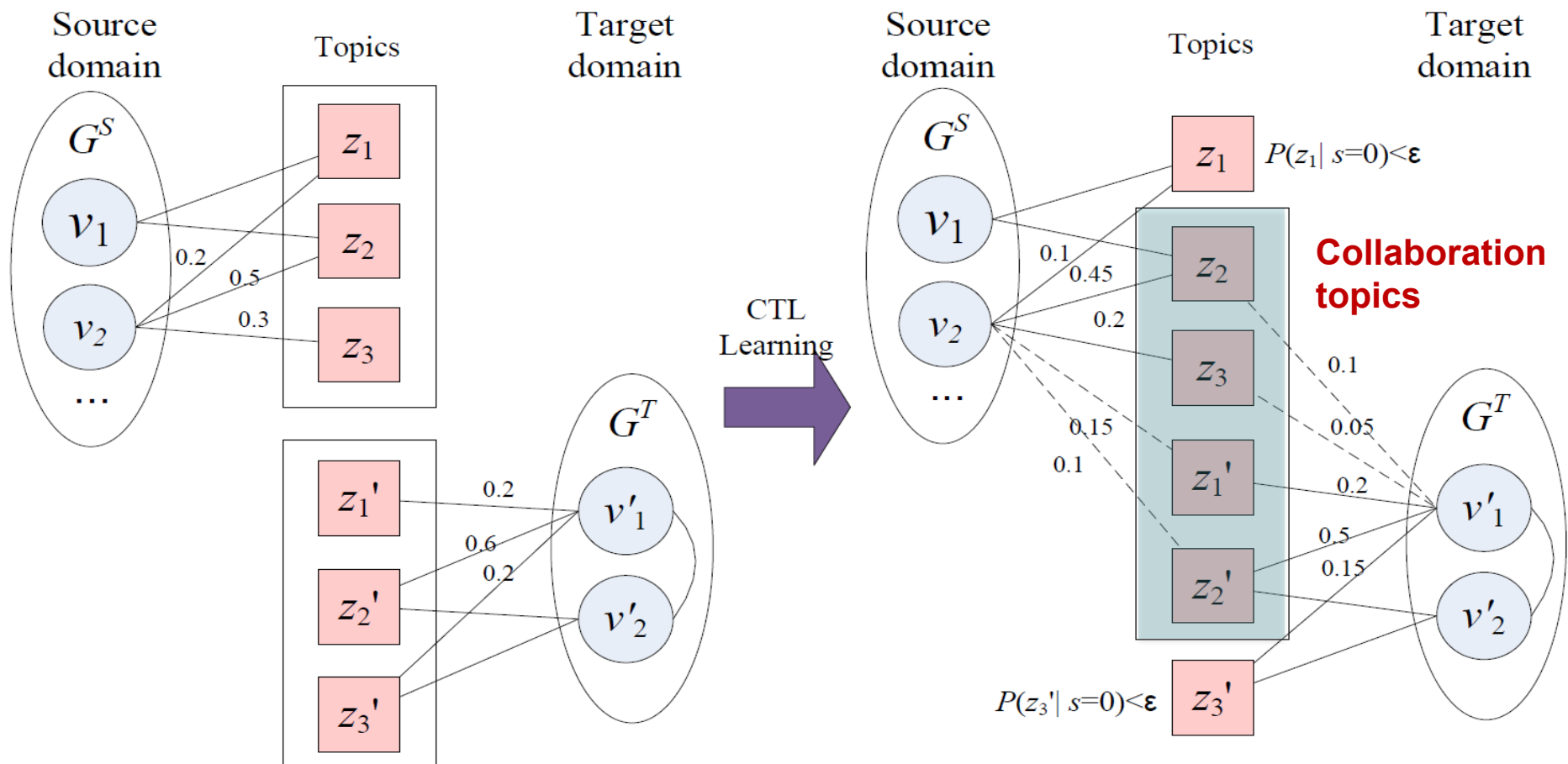
Draw a word $x_{di} \sim multi(\phi_{z_{di}})$ from z_{di} -specific word distribution;

end

Step 2:



Intuitive explanation of Step 2 in CTL



cross-domain collaboration recommendation

Experiments

Data Set and Baselines

- Arnetminer (available at <http://arnetminer.org/collaboration>)

Domain	Authors	Relationships	Source
Data Mining	6,282	22,862	KDD, SDM, ICDM, WSDM, PKDD
Medical Informatics	9,150	31,851	JAMIA, JBI, AIM, TMI, TITB
Theory	5,449	27,712	STOC, FOCS, SODA
Visualization	5,268	19,261	CVPR, ICCV, VAST, TVCG, IV
Database	7,590	37,592	SIGMOD, VLDB, ICDE

- **Baselines**
 - Content Similarity(Content)
 - Collaborative Filtering(CF)
 - Hybrid
 - Katz
 - Author Matching(Author), Topic Matching(Topic)

Performance Analysis

Training: collaboration before 2001 **Validation:** 2001-2005

Cross Domain	ALG	P@10	P@20	MAP	R@100	ARHR -10	ARHR -20
Data Mining(S) to Theory(T)	Content	10.3	10.2	10.9	31.4	4.9	2.1
	CF	15.6	13.3	23.1	26.2	4.9	2.8
	Hybrid	17.4	19.1	20.0	29.5	5.0	2.4
	Author	27.2	22.3	25.7	32.4	10.1	6.4
	Topic	28.0	26.0	32.4	33.5	13.4	7.1
	Katz	30.4	29.8	21.6	27.4	11.2	5.9
	CTL	37.7	36.4	40.6	35.6	14.3	7.5

Content Similarity(Content): based on similarity between authors' publications

Collaborative Filtering(CF): based on existing collaborations

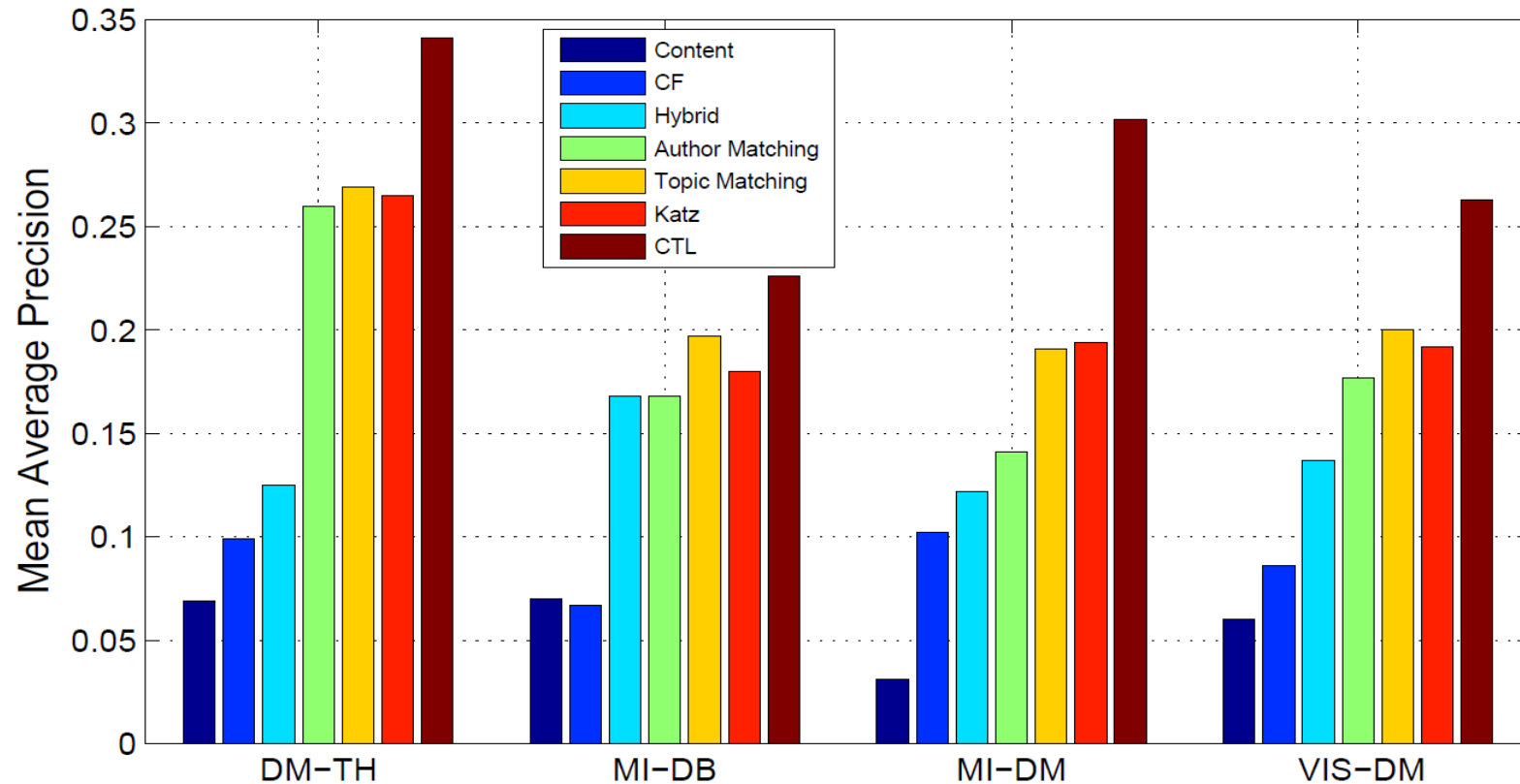
Hybrid: a linear combination of the scores obtained by the Content and the CF methods.

Katz: the best link predictor in link-prediction problem for social networks

Author Matching(Author): based on the random walk with restart on the collaboration graph

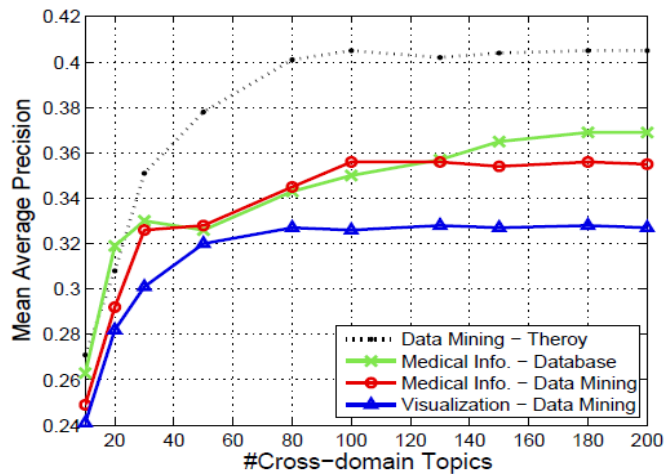
Topic Matching(Topic): combining the extracted topics into the random walking algorithm

Performance on New Collaboration Prediction

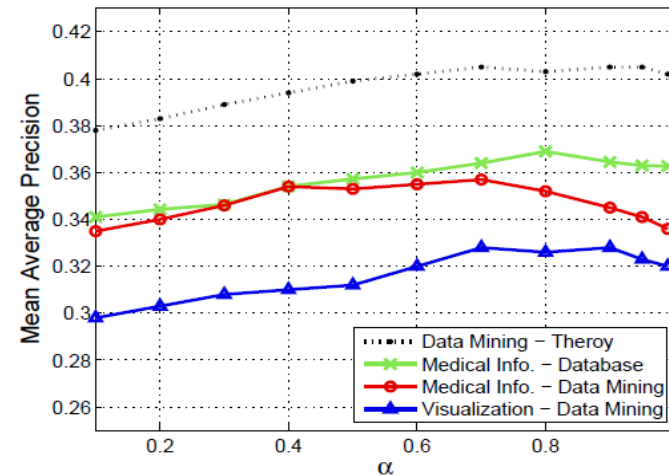


CTL can still maintain about 0.3 in terms of MAP which is significantly higher than baselines.

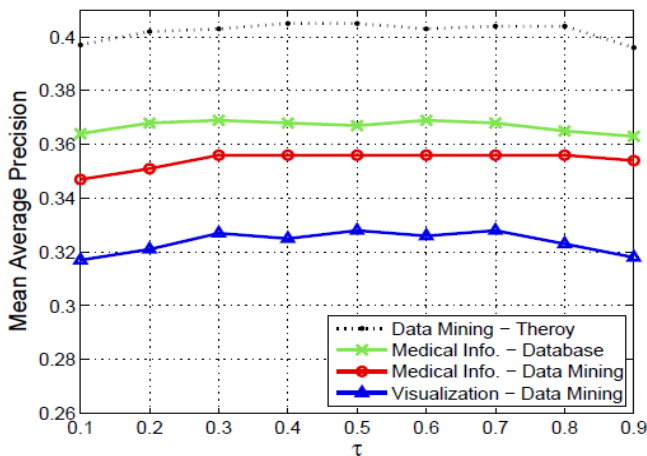
Parameter Analysis



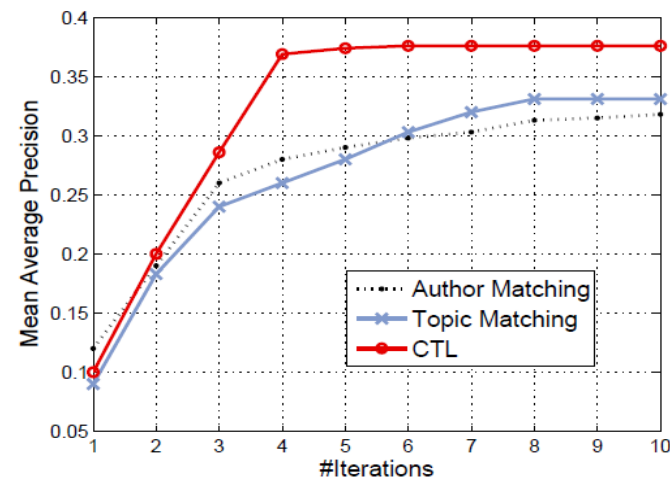
(a) number of topics T



(b) Hyperparameter α



(c) RWR parameter τ



(d) Convergence analysis

(a) varying the number of topics T

(c) varying the restart parameter τ in the random walk

(b) varying α parameter

(d) Convergence analysis

Prototype System

<http://arnetminer.org/collaborator>

Cross-Domain Collaboration Recommendation



Treemap: representing subtopic in the target domain

Yigong Shi
 Professor, School of Life Science, Tsinghua University
 H-index: 38, #Papers: 231, #Citations: 5788
 Cell Biology, Molecular Biology, Genetics & Genealogy

David J. Chen
 Professor, Rehabilitation Institute of Chicago
 H-index: 27, #Papers: 182, #Citations: 2424
 Cell Biology, Molecular Biology, Oncology

Pascale Cossart
 Huazhong Agriculture University
 H-index: 41, #Papers: 208, #Citations: 5794
 Molecular Biology, Microbiology, Cell Biology

Shoichiro Tsukita
 Professor, Department of Cell Biology, Kyoto University
 H-index: 49, #Papers: 174, #Citations: 7467
 Cell Biology, Biochemistry, Molecular Biology

Fred Chang
 Professor of Microbiology & Immunology M.D., Ph.D., University of California at San Francisco
 H-index: 22, #Papers: 115, #Citations: 1519
 Cell Biology, Molecular Biology, Pharmacology

Yang Shi
 Professor, Department of Pathology, Harvard Medical School
 H-index: 38, #Papers: 323, #Citations: 5646
 Molecular Biology, Cell Biology, Biochemistry

Robert W. Doms
 University of Pennsylvania
 H-index: 46, #Papers: 179, #Citations: 6700
 Virology, Biochemistry, Cell Biology

Recommend collaborators for Jimeng Sun from "Cell Biology"

Yigong Shi
 Professor, School of Life Science, Tsinghua University
 H-index: 38, #Papers: 231, #Citations: 5788
 Cell Biology, Molecular Biology, Genetics & Genealogy

David J. Chen
 Professor, Rehabilitation Institute of Chicago
 H-index: 27, #Papers: 182, #Citations: 2424
 Cell Biology, Molecular Biology, Oncology

Pascale Cossart
 Huazhong Agriculture University
 H-index: 41, #Papers: 208, #Citations: 5794
 Molecular Biology, Microbiology, Cell Biology

Shoichiro Tsukita
 Professor, Department of Cell Biology, Kyoto University
 H-index: 49, #Papers: 174, #Citations: 7467
 Cell Biology, Biochemistry, Molecular Biology

Fred Chang
 Professor of Microbiology & Immunology M.D., Ph.D., University of California at San Francisco
 H-index: 22, #Papers: 115, #Citations: 1519
 Cell Biology, Molecular Biology, Pharmacology

Yang Shi
 Professor, Department of Pathology, Harvard Medical School
 H-index: 38, #Papers: 323, #Citations: 5646
 Molecular Biology, Cell Biology, Biochemistry

Robert W. Doms
 University of Pennsylvania
 H-index: 46, #Papers: 179, #Citations: 6700
 Virology, Biochemistry, Cell Biology

Yigong Shi's Publications
 1 - 5 of 231 publications

Molecular mechanisms of caspase regulation during apoptosis
 Authors: Stefan J. Riedl, Yigong Shi.
 JConf: Nature Reviews Molecular Cell Biology
 Published Year: 2004 CitedBy 233

C. elegans mitochondrial factor WAH1 promotes phosphatidylserine externalization in apoptotic cells through phospholipid scramblase SCRM-1
 Authors: Xiaochen Wang, Jin Wang, Keiko Genyo-Ando, Lichuan Gu, Chun-Ling Sun, Chonglin Yang, Yong Shi, Tetsuo Kobayashi, Yigong Shi, Shohei Mizani, Xiao-Song Xie, Ding Xue.
 JConf: Nature Cell Biology
 Published Year: 2007 CitedBy 23

Multiple Apoptotic Caspase Cascades Are Required in Nonapoptotic Roles for Drosophila Spermatid Individualization
 Authors: Jun R. Huh, Stephanie Y. Vernooy, Hong Yu, Nieng Yan, Yigong Shi, Ming Guo, Bruce A. Hay.
 JConf: Plos Biology
 Published Year: 2004 CitedBy 49

Transforming Growth Factor-Mediated Transcriptional Repression of c-myc Is Dependent on Direct Binding of Smad3 to a Novel Repressive Smad Binding Element
 Authors: Joshua P. Frederick, Nicole T. Liberati, David S. Waddell, Yigong Shi, Xiao-Fan Wang.
 JConf: Molecular and Cellular Biology
 Published Year: 2004 CitedBy 42

MECHANISMS OF APOPTOSIS THROUGH STRUCTURAL BIOLOGY
 Authors: Nieng Yan, Yigong Shi.
 JConf: Annual Review of Cell and Developmental Biology
 Published Year: 2005 CitedBy 29

Recommend Collaborators & Their relevant publications

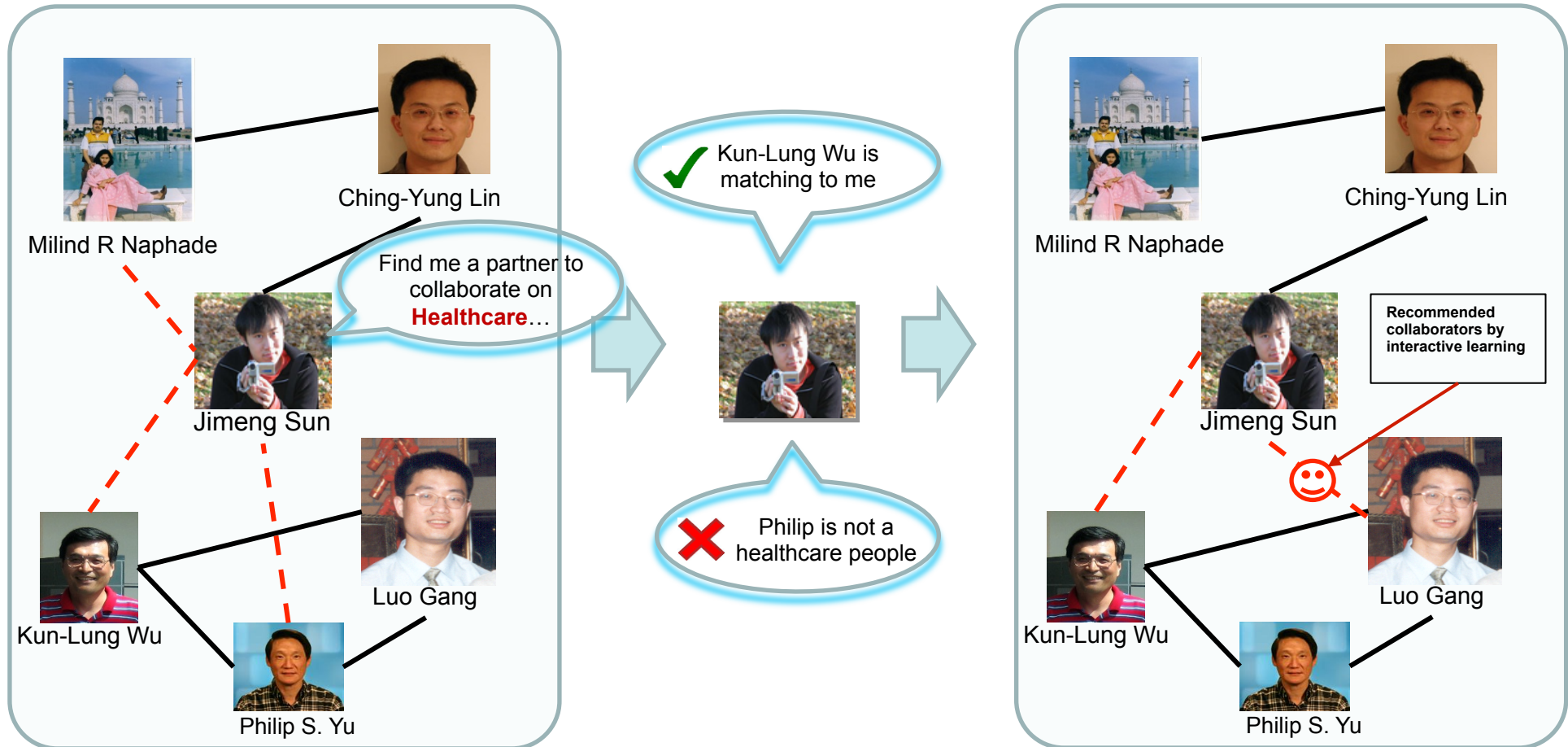
Part C:

Further incorporate user feedback
“interactive collaboration recommendation”

(ACM TKDD, TIST, WSDM 2013-14)

Example

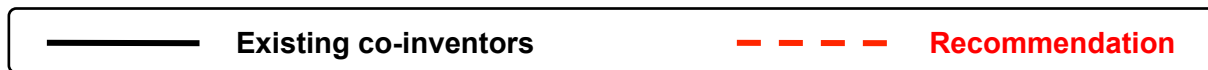
Finding co-inventors in IBM (>300,000 employees)



Recommend Candidates

Interactive feedback

Refined Recommendations





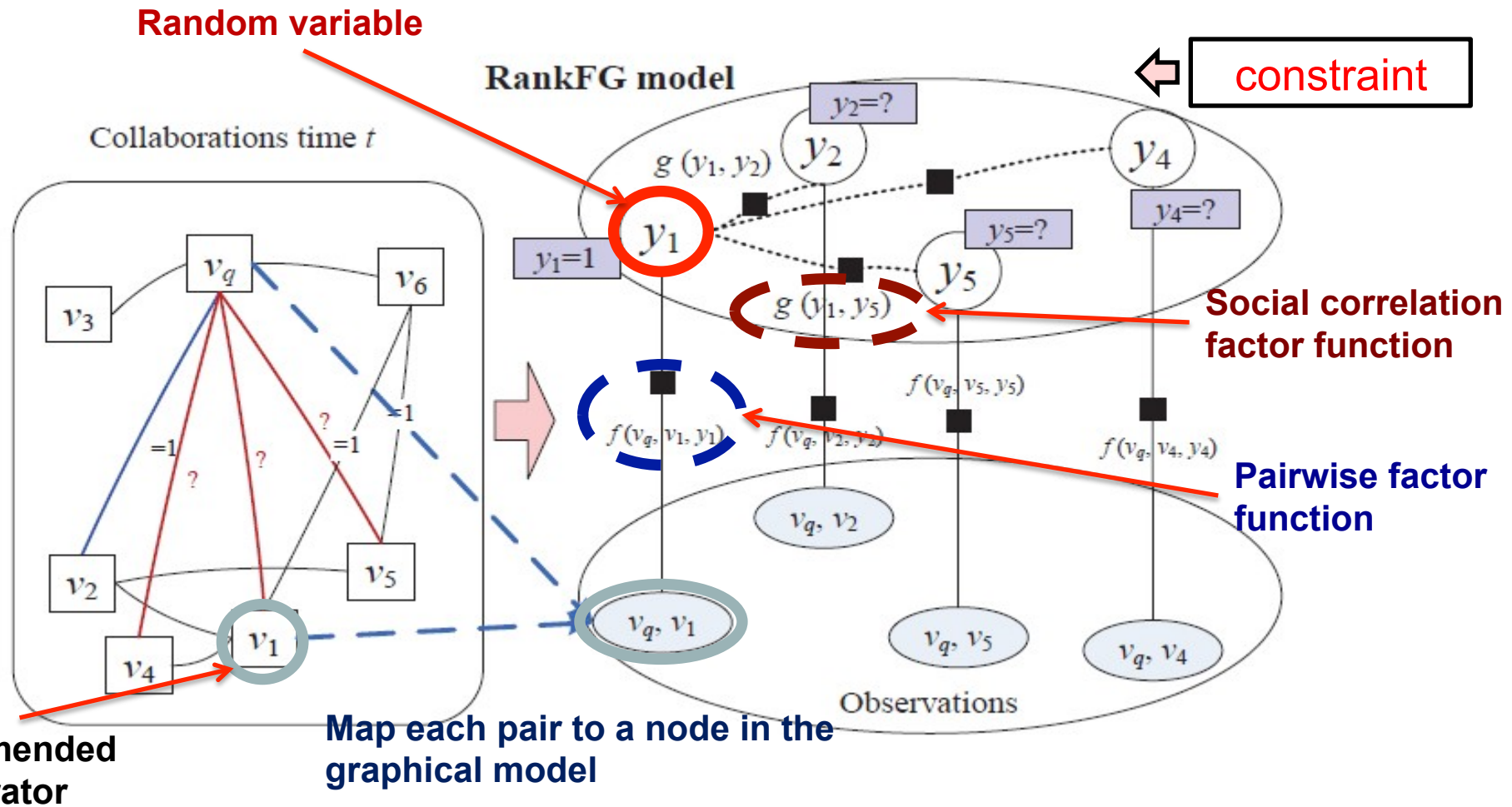
Challenges

- What are the **fundamental factors** that influence the formation of co-invention relationships?
- How to design an **interactive mechanism** so that the user can provide feedback to the system to refine the recommendations?
- How to learn the interactive recommendation framework in an **online** mode?

interactive collaboration recommendation

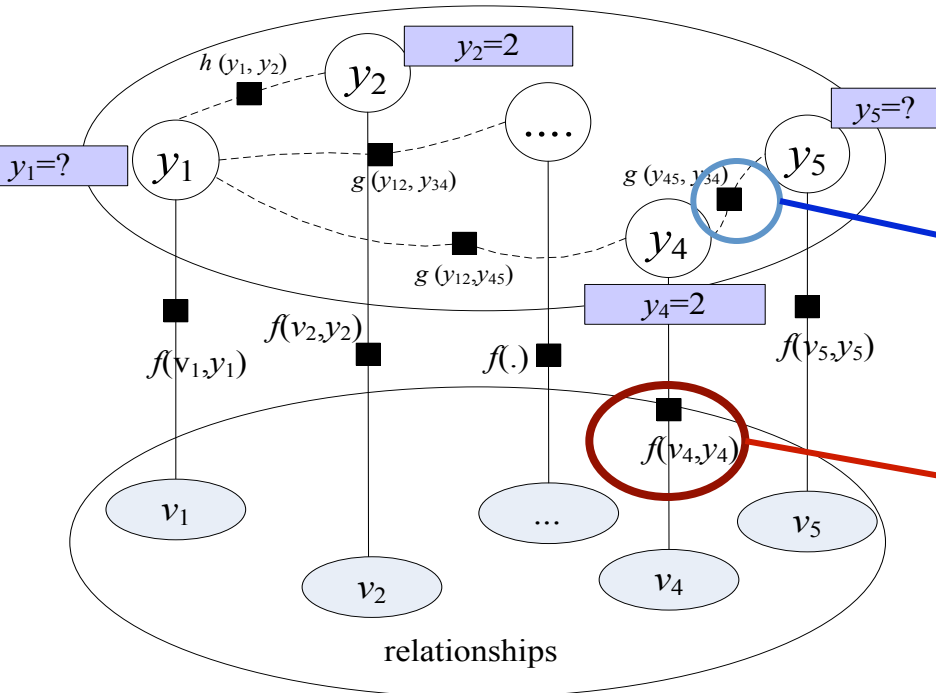
Learning framework

RankFG Model



The problem is cast as, for each relationship, identifying which type has the highest probability.

Modeling with exponential family



Partially Labeled Model

$$P(y_i | Y_{-i}) \propto \exp\left\{\sum_{c_i} \sum_k \mu_k h_k(Y_{c_i})\right\}$$

$$P(x_i | y_i) \propto \exp\left\{\sum_{j=1}^d \alpha_j g_j(x_{ij}, y_i)\right\}$$

Likelihood objective function

$$P(Y | X, G) = \frac{P(X, G | Y)P(Y)}{P(X, G)}$$

$$\propto P(X | Y) \cdot P(Y | G) = P(Y | G) \prod_i P(x_i | y_i)$$

Ranking Factor Graphs

- Pairwise factor function:

$$f(v_q, v_i, y_i) = \frac{1}{Z_a} \exp\left\{\sum_k \alpha_k \psi_k(\mathbf{x}_q, \mathbf{x}_i, y_i)\right\}$$

- Correlation factor function:

$$g(y_i, y_j) = \frac{1}{Z_b} \exp\left\{\sum_l \beta_l \phi_l(y_i, y_j)\right\}$$

- Log-likelihood objective function:

$$\begin{aligned} \log P(Y|X, \theta) &= \sum_{y_i \in Y} \sum_k \alpha_k \psi_k(\mathbf{x}_q, \mathbf{x}_i, y_i) \\ &\quad + \sum_{v_i \sim v_j} \sum_l \beta_l \phi_l(y_i, y_j) - \log Z \end{aligned}$$

- Model learning

$$\theta^* = \arg \max_{\theta} \log P(Y|X, \theta)$$

Learning Algorithm

```

Input: Query inventors  $Q = \{v_q\}$  with corresponding topics
           $\{q\}$ ,  $G = (V, E, X)$ , and the learning rate  $\eta$ ;
Output: learned parameters  $\theta$ ;
 $\theta \leftarrow \mathbf{0}$ ;
repeat
  foreach  $v_q \in Q$  and  $q$  do
    //Initialization;
     $L \leftarrow$  initialization list;
    Factor graph  $FG \leftarrow BuildFactorGraph(L)$ ;
    // Learn the parameter  $\theta$  for factor graph model;
    repeat
      foreach  $v_i \in order$  do
        | Update the messages of  $v_i$  by Eqs. 8 and 9;
      end
    until (all messages  $\mu$  do not change);
    foreach  $\theta_i \in \theta$  do
      | Calculate gradient  $\nabla_i$  according to Eq. 7;
      | Update  $\theta^{new} = \theta^{old} + \eta \cdot \nabla_i$ ;
    end
  end
until converge;

```

Expectation Computing
Loopy Belief Propagation

Algorithm 1: Learning algorithm for RankFG.



Still Challenge

How to incrementally incorporate
users' feedback?

Learning Algorithm

```

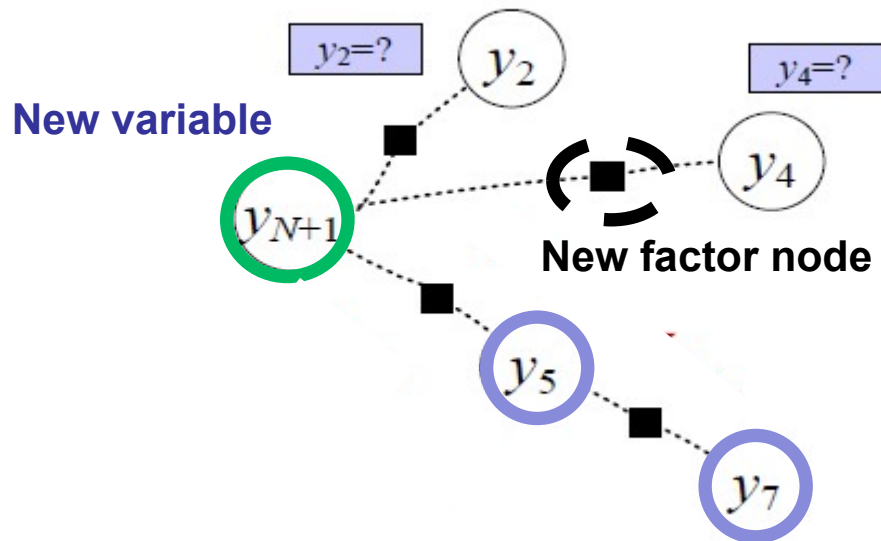
Input: Query inventors  $Q = \{v_q\}$  with corresponding topics
           $\{q\}$ ,  $G = (V, E, X)$ , and the learning rate  $\eta$ ;
Output: learned parameters  $\theta$ ;
 $\theta \leftarrow \mathbf{0}$ ;
repeat
  foreach  $v_q \in Q$  and  $q$  do
    //Initialization;
     $L \leftarrow$  initialization list;
    Factor graph  $FG \leftarrow BuildFactorGraph(L)$ ;
    // Learn the parameter  $\theta$  for factor graph model;
    repeat
      foreach  $v_i \in order$  do
        | Update the messages of  $v_i$  by Eqs. 8 and 9;
      end
    until (all messages  $\mu$  do not change);
    foreach  $\theta_i \in \theta$  do
      | Calculate gradient  $\nabla_i$  according to Eq. 7;
      | Update  $\theta^{new} = \theta^{old} + \eta \cdot \nabla_i$ ;
    end
  end
until converge;

```

Incremental estimation

Algorithm 1: Learning algorithm for RankFG.

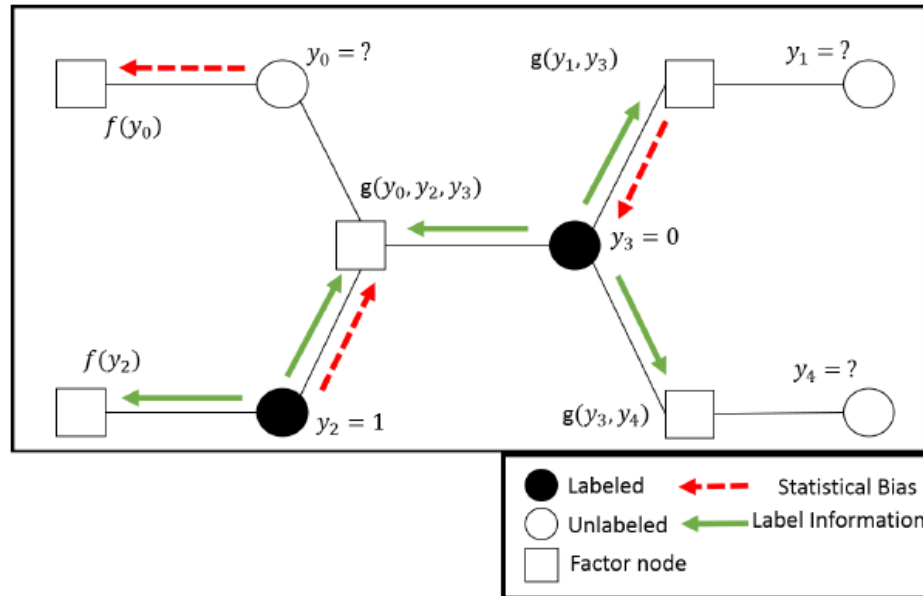
Interactive Learning



- 1) add new factor nodes to the factor graph built in the model learning process.
- 2) l -step message passing:
 - Start from the new variable node (y_{N+1} to node).
 - Send messages to all of its neighborhood factors.
 - Propagate the messages up to l -step
 - Perform a backward messages passing.
- 3) Calculate an approximate value of the marginal probabilities of the newly factors.

$$\mathbb{E}^{new}[\cdot] = \frac{N}{N+1} \mathbb{E}^{old}[\cdot] + \frac{1}{N+1} \sum_k \theta_k \phi_k(\mathbf{x}_{N+1}, \mathbf{y}_{N+1})$$

From passive interactive to active



- Entropy

$$\mu(v) = \sum_{y \in \mathcal{Y}} B_v(y) \log \frac{1}{B_v(y)}$$

- Threshold

$$t(v) = \min\{\lceil \eta(\mu_{\max} - \mu(v))d(v) \rceil, d(v)\}$$

- Influence model

$$f_{\tau}(v) = \begin{cases} 1 & \text{if } \sum_{u \in \text{NB}(v)} f_{\tau-1}(u) \geq t(v) \\ 0 & \text{if } \sum_{u \in \text{NB}(v)} f_{\tau-1}(u) < t(v) \end{cases}$$

[1] Z. Yang, J. Tang, and B. Xu. Active Learning for Networked Data Based on Non-progressive Diffusion Model. WSDM'14.

[2] L. Shi, Y. Zhao, and J. Tang. Batch Mode Active Learning for Networked Data. ACM Transactions on Intelligent Systems and Technology (TIST), Volume 3, Issue 2 (2012), Pages 33:1--33:25.

Active learning via Non-progressive diffusion model



- Maximizing the diffusion

$$\max_{V_S \subseteq V_U} \left\{ \max_{V_T \subseteq V_U} |V_T| \right\}, \quad |V_S| \leq k$$

with the constraints:

$$f_0(v) = 1 \iff v \in V_S \quad (2)$$

$$\exists \tau_M \text{ s.t. } \forall v \in V_T \forall \tau > \tau_M f_\tau(v) = 1 \quad (3)$$

$$\forall v \in V_U \setminus V_T, \forall u \in V_T, \mu(v) \leq \mu(u) \quad (4)$$

$$f_\tau(v) = 1 \iff \sum_{u \in \text{NB}(v)} f_{\tau-1}(u) \geq t(v) \quad (5)$$

NP-hard!

MinSS

- Greedily expand V_p

```
5   $V_p \leftarrow V_U \setminus V_\tau$ 
6  sort nodes in  $V_p$  in descending order of  $t(v)$  as
    $v_1, v_2, \dots, v_p$ 
7  foreach  $v \in V_\tau$  do
8  |  $w(v) \leftarrow 0$ 
9  for  $i \leftarrow 1$  to  $p$  do
10 | if  $w(u) < d(u) - t(u) \forall u \in NB(v_i) \cap V_\tau$  then
11 | | foreach  $u \in NB(v_i) \cap V_\tau$  do
12 | | |  $w(u) \leftarrow w(u) + 1$ 
13 | |  $V_p \leftarrow V_p \setminus \{v_i\}$ 
```

MinSS(cont.)

```
14  if  $V_P = \emptyset$  then
15      sort nodes in  $V_\tau$  in ascending order of  $d(v)$  as
       $v_1, v_2, \dots, v_m$ 
16      foreach  $v \in V_\tau$  do
17           $w(v) \leftarrow 0$ 
18      foreach  $i \leftarrow 1$  to  $m$  do
19          if  $\exists u \in NB(v_i) \cap V_\tau$  st.  $w(u) = d(u) - t(u)$ 
          then
20               $V_S \leftarrow V_S \cup \{v_i\}$ 
21          else
22              foreach  $u \in NB(v_i) \cap V_\tau$  do
23                   $w(u) \leftarrow w(u) + 1$ 
```

Lower Bound and Upper Bound

Theorem 3. Lower Bound. Let $D(V) = \sum_{v \in V} d(v)$, $T(V) = \sum_{v \in V} t(v)$, and suppose $t(v) \leq \beta d(v)$ for all $v \in V$. If $2T(V_U) - D(V_U) > 0$ and $V_T = V_U$, we have an lower bound for optimal solution $|V_{S,\text{opt}}|$ to problem 2.

$$|V_{S,\text{opt}}| \geq \frac{2T(V_U) - D(V_U)}{\beta\Delta} \quad (9)$$

Theorem 4. Upper Bound. Suppose $t(v) \leq \beta d(v)$ for all $v \in V$, we can derive an upper bound for MinSS algorithm.

$$|V_S| \leq \frac{\beta\Delta}{1 - \beta + \beta\Delta} |V_U|$$

Approximation Ratio

$$|V_{S,opt}| \geq \frac{2T(V_U) - D(V_U)}{\beta\Delta}$$



$$|V_S| \leq \frac{\beta\Delta}{1 - \beta + \beta\Delta} |V_U|$$

COROLLARY 2. Approximation Ratio. Let $V_{S,g}$ denote the solution given by MinSS algorithm, $V_{S,opt}$ represent the optimal solution and Δ be the maximum degree in the graph. Suppose $t(v) \leq \beta d(v)$ for all $v \in V$, if $V_T = V_U$ and $2T(V_U) > D(V_U)$, we have

$$\frac{|V_{S,g}|}{|V_{S,opt}|} \leq \frac{(\beta\Delta)^2}{(1 - \beta + \beta\Delta) \cdot \mathbb{E}[2t(v) - d(v)]} \quad (10)$$

where $\mathbb{E}[\cdot]$ represents the expectation over all samples in the network.

interactive collaboration recommendation

Experiments

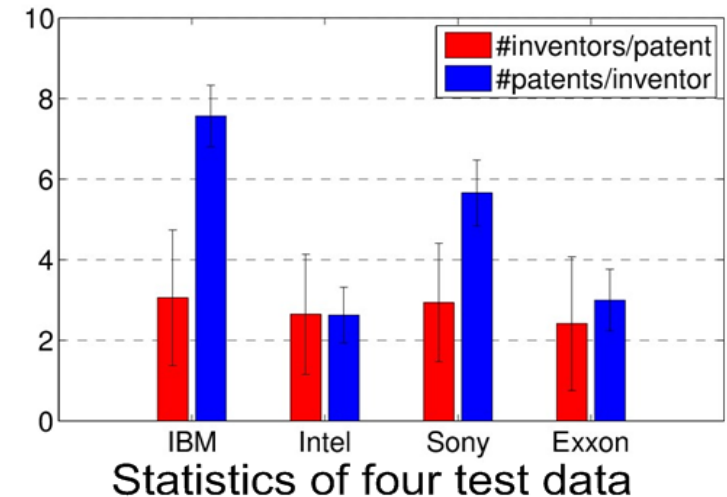
Data Set

- PatentMiner (pminer.org)

DataSet	Inventors	Patents	Average increase #patent	Average increase #co-invention
IBM	55,967	46,782	8.26%	11.9%
Intel	18,264	54,095	18.8%	35.5%
Sony	8,505	31,569	11.7%	13/0%
Exxon	19,174	53,671	10.6%	14.7%

- Baselines:

- Content Similarity (Content)
- Collaborative Filtering (CF)
- Hybrid
- SVM-Rank



Performance Analysis-IBM

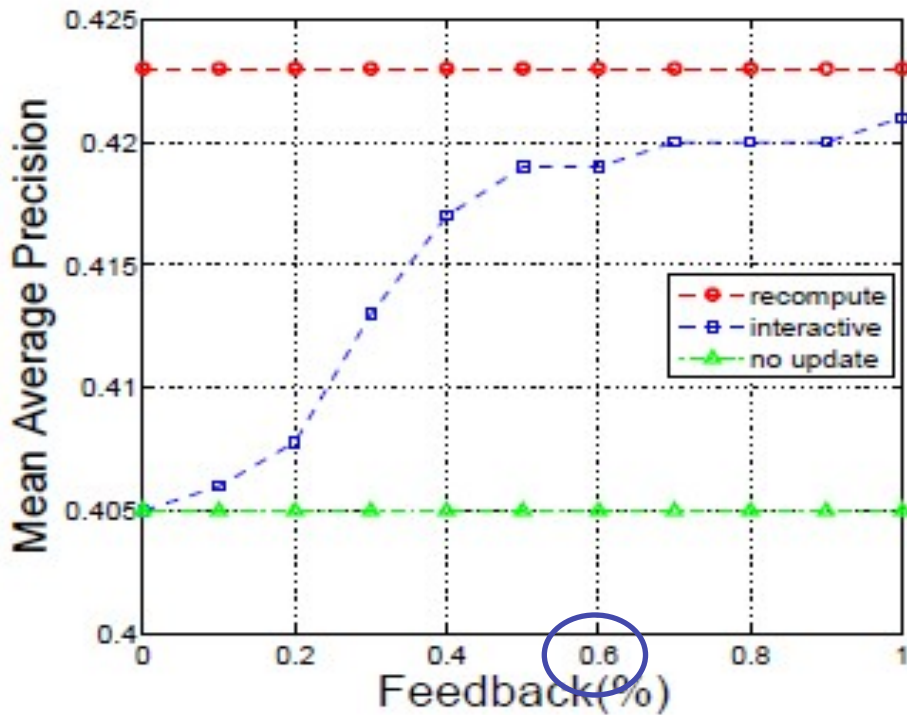
Training: collaboration before 2000 Validation: 2001-2010

Data	ALG	P@5	P@10	P@15	P@20	MAP	R@100
IBM	Content	23.0	23.3	18.8	15.6	24.0	33.7
	CF	13.8	12.8	11.3	11.5	21.7	36.4
	Hybrid	13.9	12.8	11.5	11.5	21.8	36.7
	SVMRank	13.3	11.9	9.6	9.8	22.2	43.5
	RankFG	31.1	27.5	25.6	22.4	40.5	46.8
	RankFG+	31.2	27.5	26.6	22.9	42.1	51.0

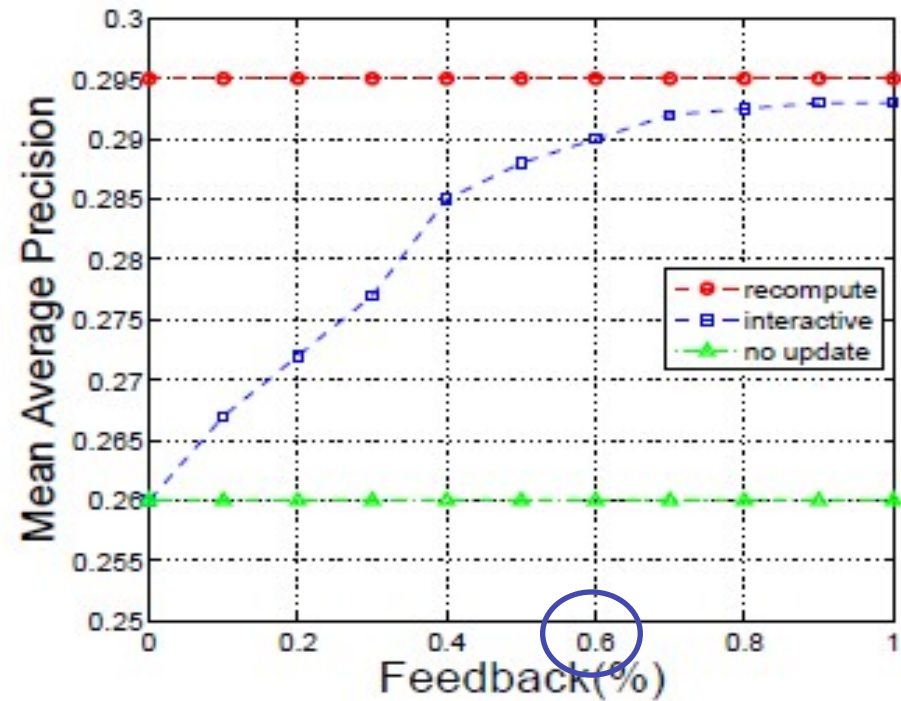


RankFG+: it uses the proposed RankFG model with 1% interactive feedback.

Interactive Learning Analysis



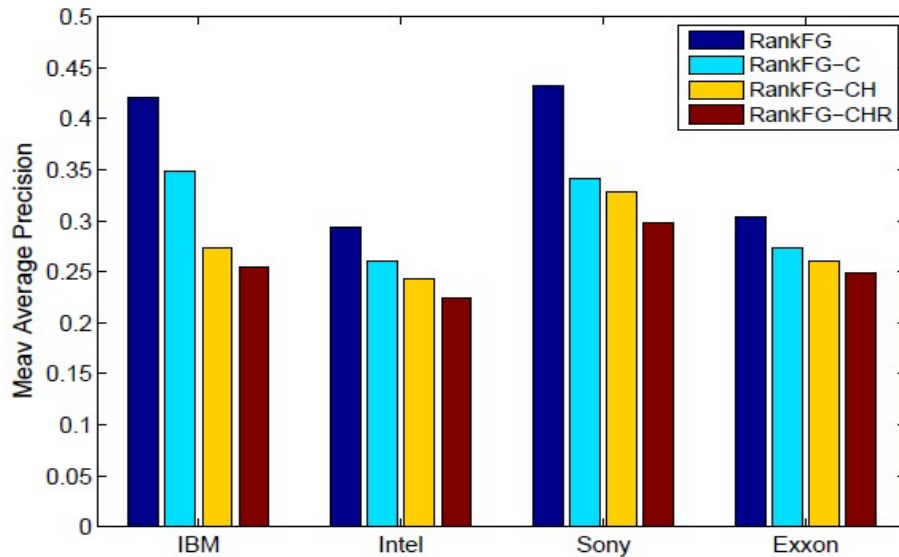
(a) IBM



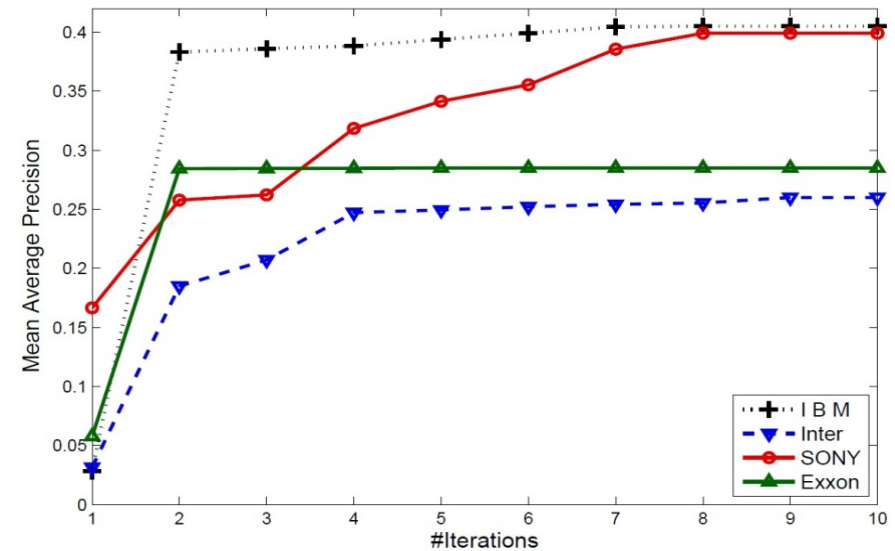
(b) Intel

Interactive learning achieves a **close performance** to the complete learning with only **1/100** of the running time used for complete training.

Parameter Analysis



Factor contribution analysis



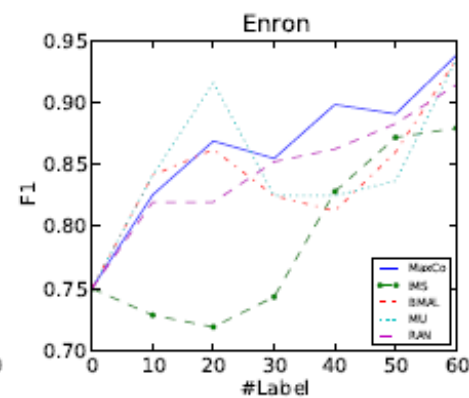
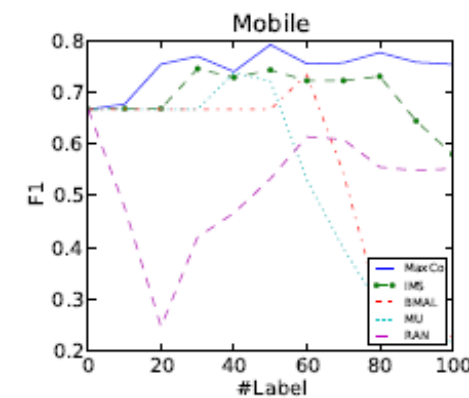
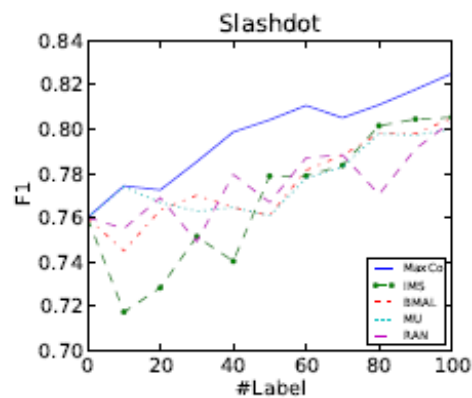
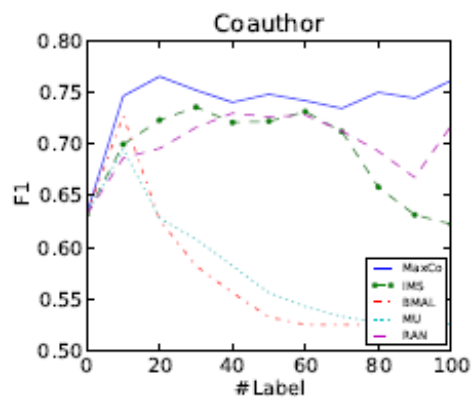
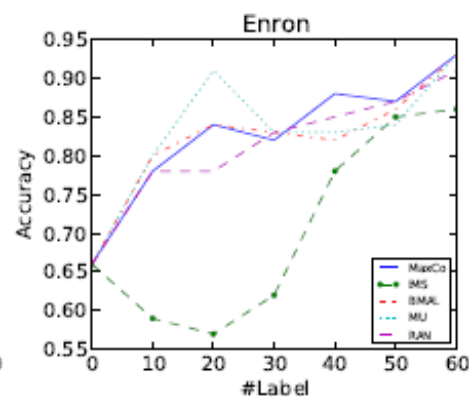
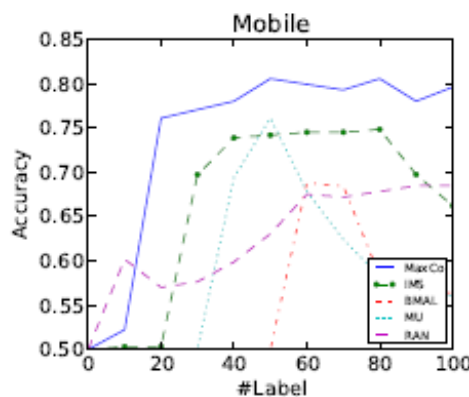
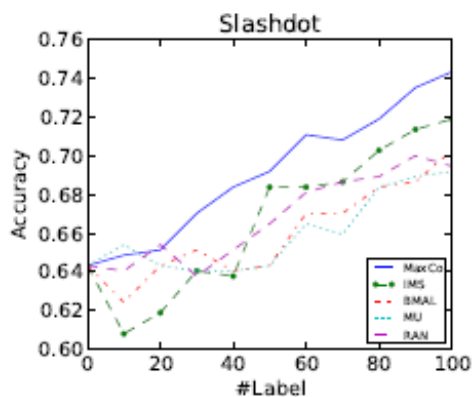
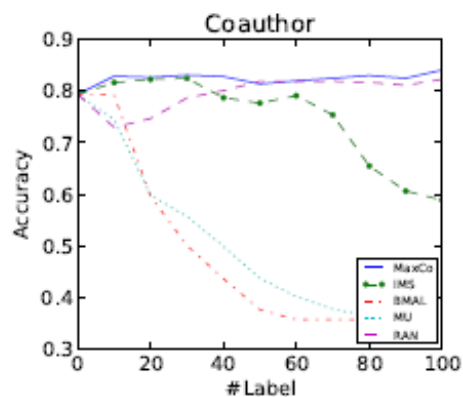
Convergence analysis

RankFG-C: stands for ignoring referral chaining factor functions.

RankFG-CH: stands for ignoring both referral chaining and homophily.

RankFG-CHR: stands for further ignoring recency.

Results of Active Learning



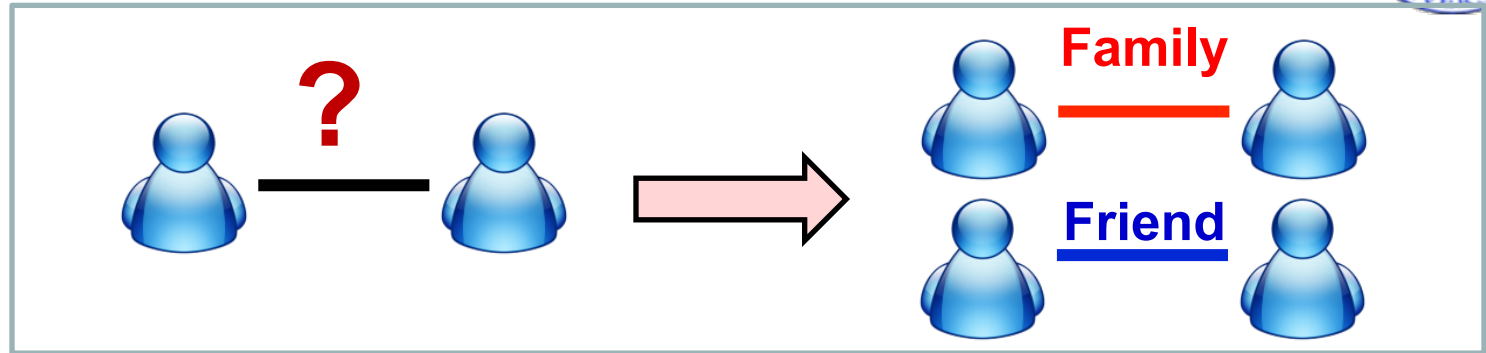


Summaries

- Inferring social ties in single network
 - Time-dependent factor graph model
- Cross-domain collaboration recommendation
 - Cross-domain topic learning
- Interactive collaboration recommendation
 - Ranking factor graph model
 - Active learning via non-progressive diffusion

Future Work

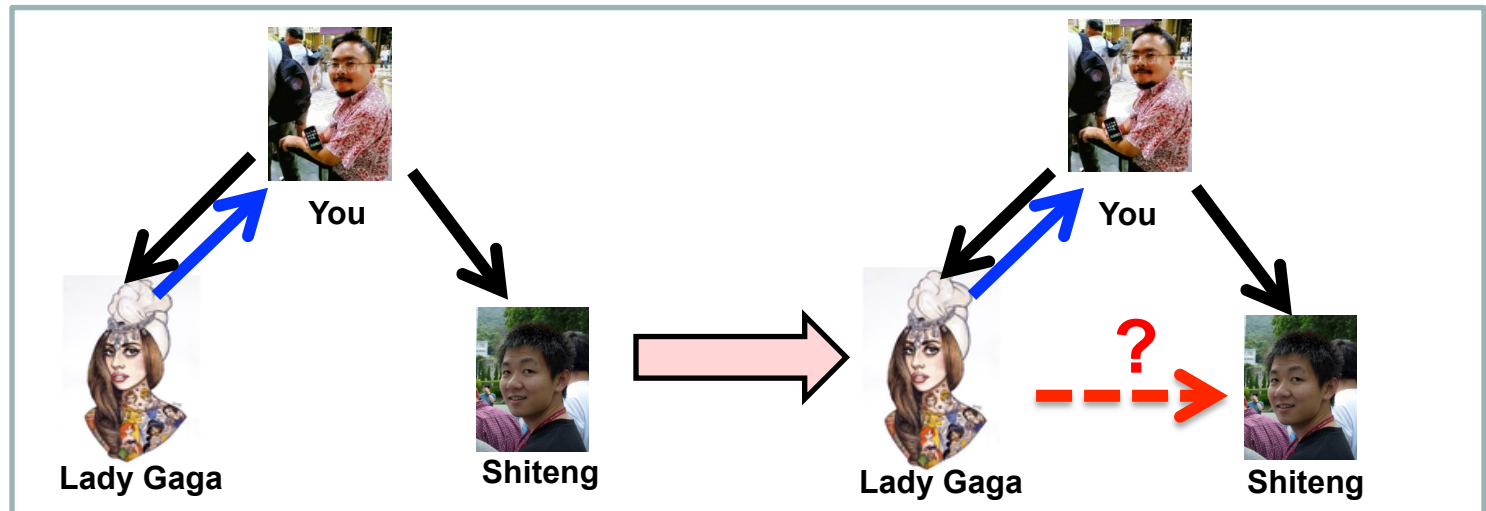
Inferring social ties



Reciprocity



Triadic Closure



References



- Tiancheng Lou, Jie Tang, John Hopcroft, Zhanpeng Fang, Xiaowen Ding. Learning to Predict Reciprocity and Triadic Closure in Social Networks. In **TKDD**, 2013.
- Yi Cai, Ho-fung Leung, Qing Li, Hao Han, Jie Tang, Juanzi Li. Typicality-based Collaborative Filtering Recommendation. *IEEE Transaction on Knowledge and Data Engineering (TKDE)*.
- Honglei Zhuang, Jie Tang, Wenbin Tang, Tiancheng Lou, Alvin Chin, and Xia Wang. Actively Learning to Infer Social Ties. **DMKD**, Vol. 25, Issue 2 (2012), pages 270-297.
- Lixin Shi, Yuhang Zhao, and Jie Tang. Batch Mode Active Learning for Networked Data. *ACM Transactions on Intelligent Systems and Technology (TIST)*, Volume 3, Issue 2 (2012), Pages 33:1--33:25.
- Jie Tang, Jing Zhang, Ruoming Jin, Zi Yang, Keke Cai, Li Zhang, and Zhong Su. Topic Level Expertise Search over Heterogeneous Networks. **Machine Learning Journal**, Vol. 82, Issue 2 (2011), pages 211-237.
- Zhilin Yang, Jie Tang, and Bin Xu. Active Learning for Networked Data Based on Non-progressive Diffusion Model. **WSDM'14**.
- Sen Wu, Jimeng Sun, and Jie Tang. Patent Partner Recommendation in Enterprise Social Networks. **WSDM'13**, pages 43-52.
- Jie Tang, Sen Wu, Jimeng Sun, and Hang Su. Cross-domain Collaboration Recommendation. **KDD'12**, pages 1285-1293. (Full Presentation & Best Poster Award)
- Jie Tang, Bo Wang, Yang Yang, Po Hu, Yanting Zhao, Xinyu Yan, Bo Gao, Minlie Huang, Peng Xu, Weichang Li, and Adam K. Usadi. PatentMiner: Topic-driven Patent Analysis and Mining. **KDD'12**, pages 1366-1374.
- Jie Tang, Tiancheng Lou, and Jon Kleinberg. Inferring Social Ties across Heterogeneous Networks. **WSDM'12**, pages 743-752.
- Chi Wang, Jiawei Han, Yuntao Jia, Jie Tang, Duo Zhang, Yintao Yu, and Jingyi Guo. Mining Advisor-Advisee Relationships from Research Publication Networks. **KDD'10**, pages 203-212.
- Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. ArnetMiner: Extraction and Mining of Academic Social Networks. **KDD'08**, pages 990-998.

Thank you !

Collaborators: John Hopcroft, Jon Kleinberg (**Cornell**)

Jiawei Han and Chi Wang (**UIUC**)

Tiancheng Lou (**Google**)

Jimeng Sun (**IBM**)

Jing Zhang, Zhanpeng Fang, Zi Yang, Sen Wu (**THU**)

Jie Tang, KEG, Tsinghua U,
Download all data & Codes,

<http://keg.cs.tsinghua.edu.cn/jietang>
<http://arnetminer.org/download>