

# **Mining Knowledge from Data: An Information Network Analysis Approach**

**Jiawei Han<sup>†</sup>   Yizhou Sun<sup>†</sup>   Xifeng Yan<sup>§</sup>   Philip S. Yu<sup>‡</sup>**

**<sup>†</sup>University of Illinois at Urbana-Champaign**

**<sup>§</sup> University of California at Santa Barbara**

**<sup>‡</sup>University of Illinois at Chicago**


**Acknowledgements: NSF, ARL, NASA, AFOSR (MURI), Microsoft, IBM, Yahoo!, Google, HP Lab & Boeing**

**April 3, 2012**

---

# Outline

---

- **Motivation:** Why Mining Information Networks? 
  - **Part I:** Clustering, Ranking and Classification
    - Clustering and Ranking in Information Networks
    - Classification of Information Networks
  - **Part II:** Meta-Path Based Exploration of Information Networks
    - Similarity Search in Information Networks
    - Relationship Prediction in Information Networks
  - **Part III:** Advanced Topics on Information Network Analysis
    - Role Discovery and OLAP in Information Networks
    - Relation Strength Learning in Information Networks
    - Mining Evolution and Dynamics of Information Networks
  - **Conclusions**
-

# Motivation

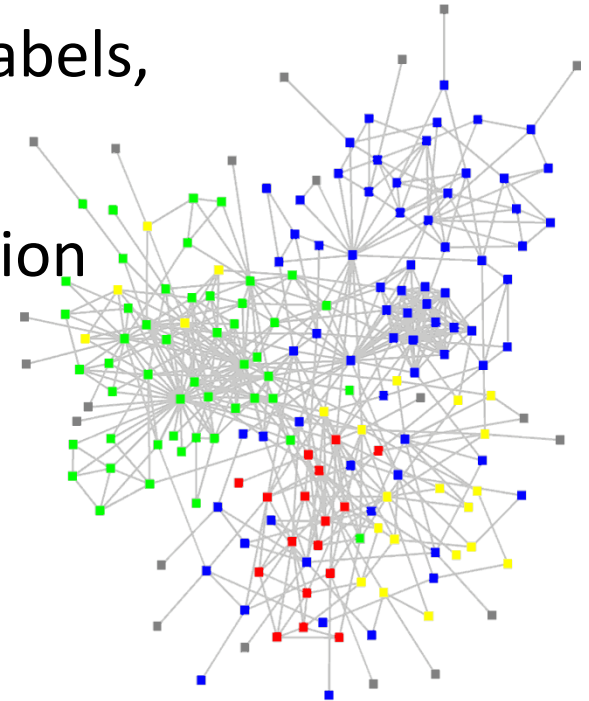
---

- Traditional view of a database
  - Database: a data repository
  - Database system: supports organized and efficient data storage, update, retrieval, management, ...
- Our view of a database: an organized info. network!
  - Information-rich, inter-related data relations and records form one or a set of gigantic, interconnected, multi-typed heterogeneous information networks
  - Surprisingly rich knowledge can be derived from such database-information network (DB-InfoNet)
- How to uncover knowledge “buried” in databases?
  - Exploring the power of multi-typed, heterogeneous links
  - Mining “semi-structured” heterogeneous information networks!

# What Are Information Networks?

---

- Information network: A network where each node represents an entity (e.g., actor in a social network) and each link (e.g., tie) a relationship between entities
  - Each node/link may have attributes, labels, and weights
  - Link may carry rich semantic information





# Information Networks Are Everywhere

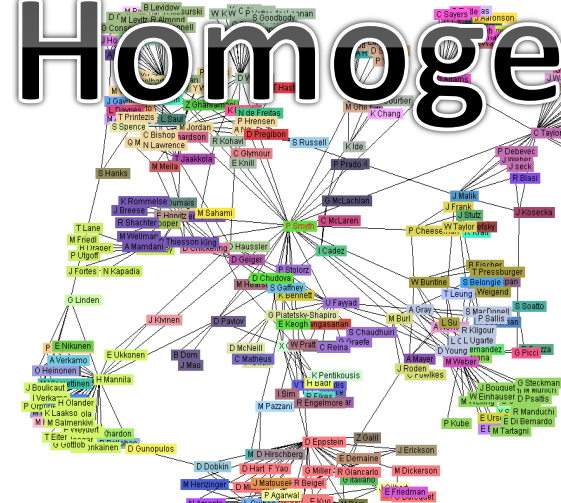


They are all treated as

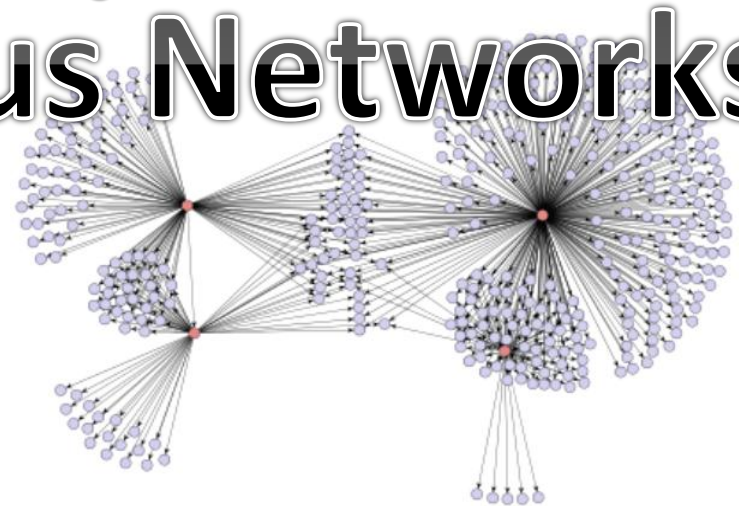
Social Networking Websites

Biological Network: Protein Interaction

# Homogeneous Networks!



Research Collaboration Network



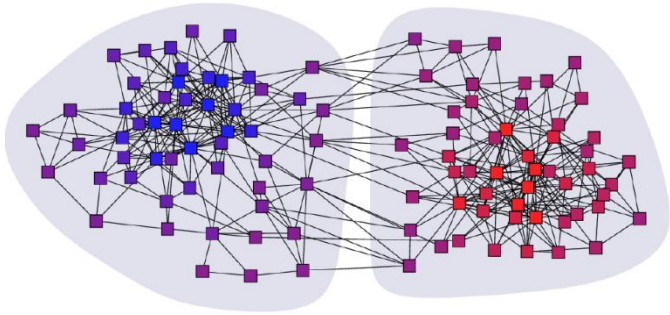
Product Recommendation Network via Emails

# Homogeneous Information Networks

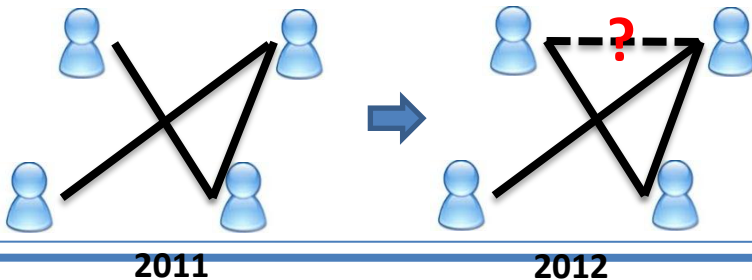
- Single object type and single link type
  - Link analysis based applications



**Ranking** web pages [Brin and Page, 1998]



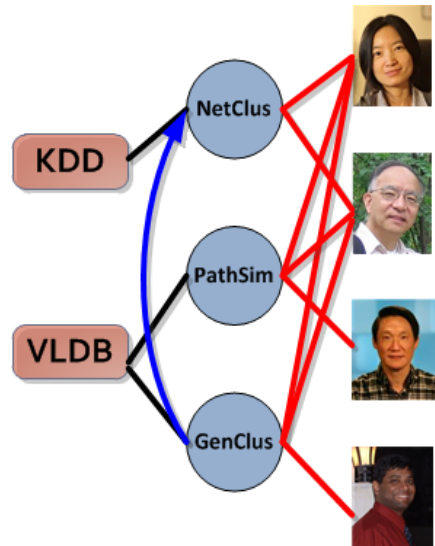
**Clustering** books about politics [Newman, 2006]



**Link Prediction** [Kleinberg, 2003]

# Heterogeneous Information Networks

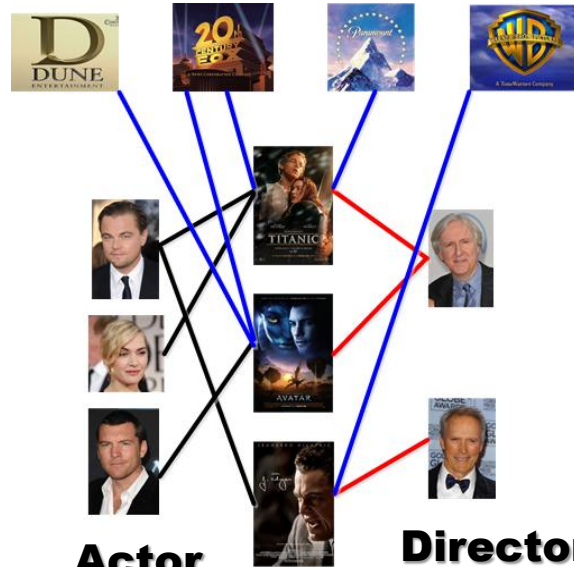
- Multiple object types and/or multiple link types



**Venue Paper Author**

DBLP Bibliographic Network

**Movie Studio**



**Actor**

**Movie**

**Director**

The IMDB Movie Network



The Facebook Network

- Homogeneous networks are **Information loss** projection of heterogeneous networks!
- New problems** are emerging in heterogeneous networks!

# Heterogeneous Networks Are Ubiquitous

---

- Healthcare
  - Doctor, patient, disease, treatment
- Content sharing websites
  - Video, image, user, comment
- E-Commerce
  - Seller, buyer, product, review
- News
  - Person, organization, location, text





# What Can be Mined from Heterogeneous Networks?

- DBLP: A Computer Science bibliographic database




Yizhou Sun, Jiawei Han, Charu C. Aggarwal, Nitesh V. Chawla: When will it happen?: relationship prediction in heterogeneous information networks. WSDM 2012: 663-672

A sample publication record in DBLP (>1.8 M papers, >0.7 M authors, >10 K venues)

Knowledge hidden in DBLP Network	Mining Functions	Publications
How are CS research areas <b>structured</b> ?	Clustering	EDBT'09, KDD'09, ICDM'09
Who are the <b>leading</b> researchers on Web search?	Ranking	EDBT'09, KDD'09,
Who are the <b>peer</b> researchers of Jure Leskovec?	Similarity Search	VLDB'11
Whom <b>will</b> Christos Faloutsos <b>collaborate with</b> in the future?	Relationship Prediction	ASONAM'11
Whether <b>will</b> an author <b>publish</b> a paper in KDD, and <b>when</b> ?	Relationship Prediction with Time	WSDM'12
Which types of <b>relationships</b> are most <b>influential</b> for an author to decide her topics?	Relation Strength Learning	VLDB'12, KDD'12 submission

# Outline

---

- **Motivation:** Why Mining Information Networks?
  - **Part I:** Clustering, Ranking and Classification
    - Clustering and Ranking in Information Networks 
    - Classification of Information Networks
  - **Part II:** Meta-Path Based Exploration of Information Networks
    - Similarity Search in Information Networks
    - Relationship Prediction in Information Networks
  - **Part III:** Advanced Topics on Information Network Analysis
    - Role Discovery and OLAP in Information Networks
    - Relation Strength Learning in Information Networks
    - Mining Evolution and Dynamics of Information Networks
  - **Conclusions**
-

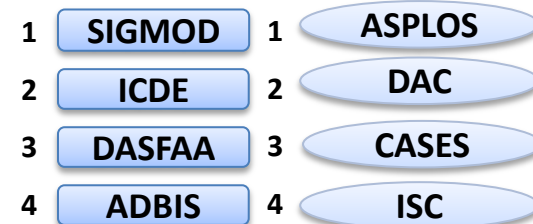
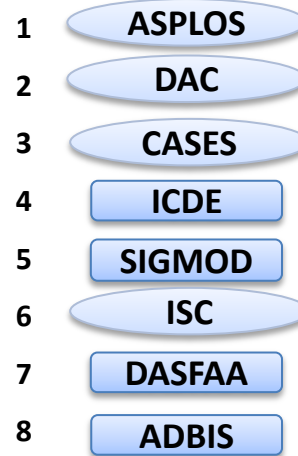
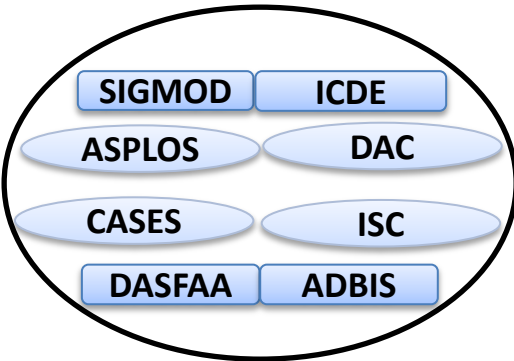
# Ranking and Clustering: Two Critical Functions

## Ranking

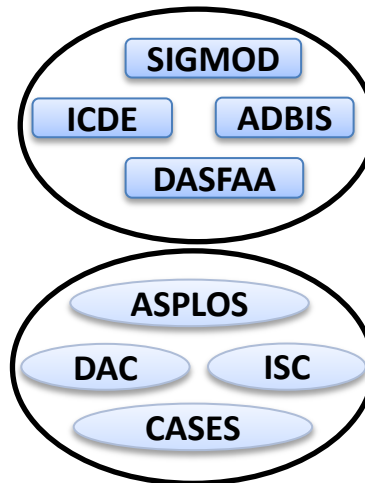
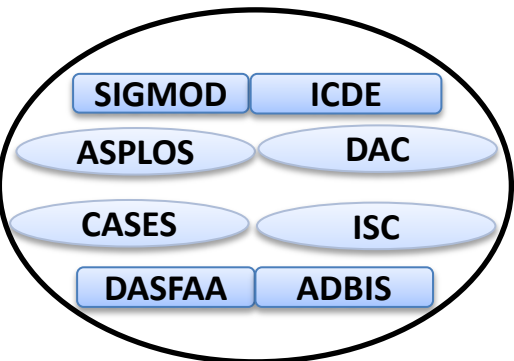
Comparing apples and oranges?

 : Database Conferences

 : Hardware and Architecture Conferences



## Clustering



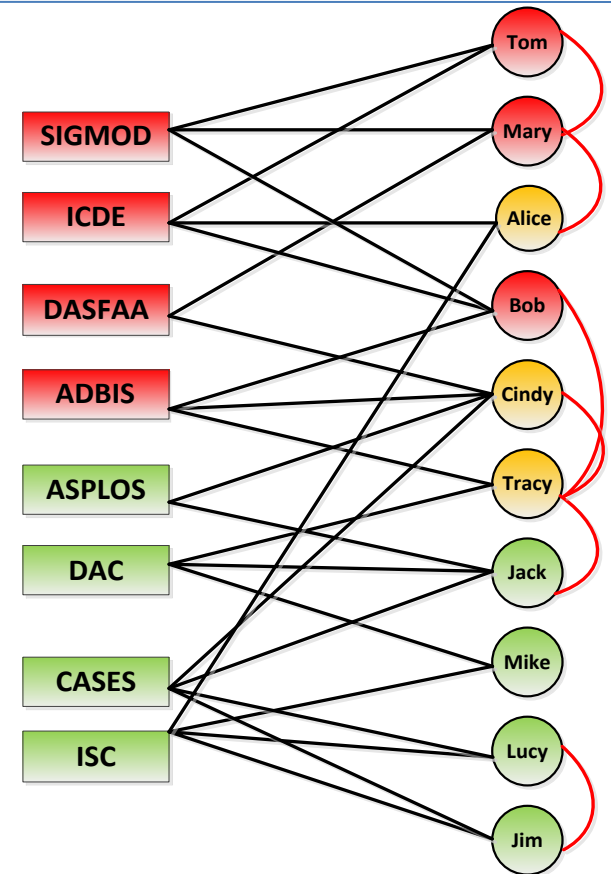
**A better solution:  
Integrating clustering  
with ranking**

Not distinguishing objects in each cluster?

# RankClus: Integrating Clustering with Ranking

## [Sun, EDBT'09]

- A case study on bi-typed DBLP network
  - Links exist between
    - Conference (X) and author (Y)
    - Author (Y) and author (Y)
  - A matrix denoting the weighted links
    - $W = \begin{bmatrix} W_{XX} & W_{XY} \\ W_{YX} & W_{YY} \end{bmatrix}$
  - Goal:
    - Clustering and ranking conferences via authors
      - **Simple solution: Project the bi-typed network into homogeneous conference network + spectral clustering [Shi & Malik, 2000]**

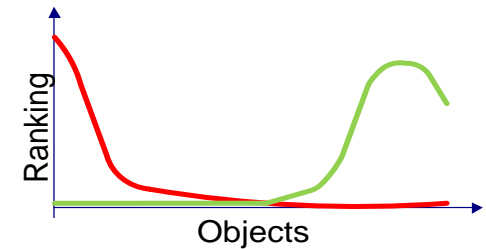




# Idea: Ranking and Clustering Mutually Enhance Each Other

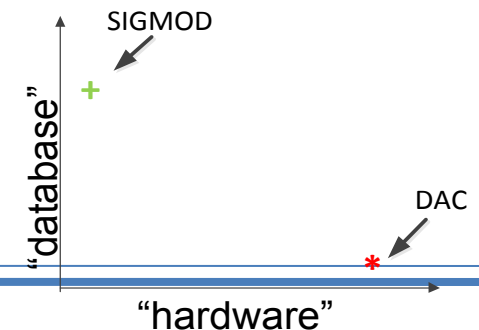
- Better clustering => Conditional ranking distributions are more distinguishing from each other
  - Conditional ranking distribution serves as the **feature** of each cluster

$P(\bullet | \text{area} = \text{"database"})$  vs.  $P(\bullet | \text{area} = \text{"hardware"})$

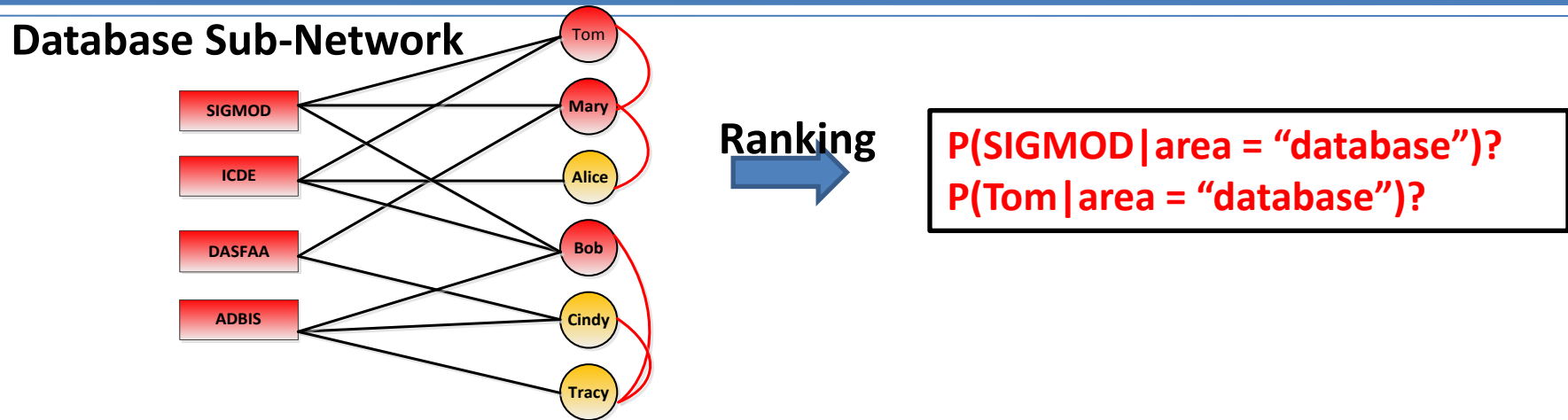


- Better ranking => Better metric for objects can be learned from the ranking for better clustering
  - Posterior probabilities for each object in each cluster serves as the **new metric** for each object

$( P(\text{area} = \text{"database"} | \text{SIGMOD}), P(\text{area} = \text{"hardware"} | \text{SIGMOD}) )$



# Simple Ranking vs. Authority Ranking



- Simple Ranking

- Proportional to # of publications of an author / a conference
- Considers only **immediate neighborhood** in the network

What about an author publishing 100 papers in low reputation conferences?

- Authority Ranking:

- More sophisticated “rank rules” are needed
- **Propagate** the ranking scores in the network over different types

# Rules for Authority Ranking

---

- Rule 1: Highly ranked authors publish *many* papers in highly ranked conferences

$$\vec{r}_Y(j) = \sum_{i=1}^m W_{YX}(j, i) \vec{r}_X(i)$$

- Rule 2: Highly ranked conferences attract *many* papers from *many* highly ranked authors

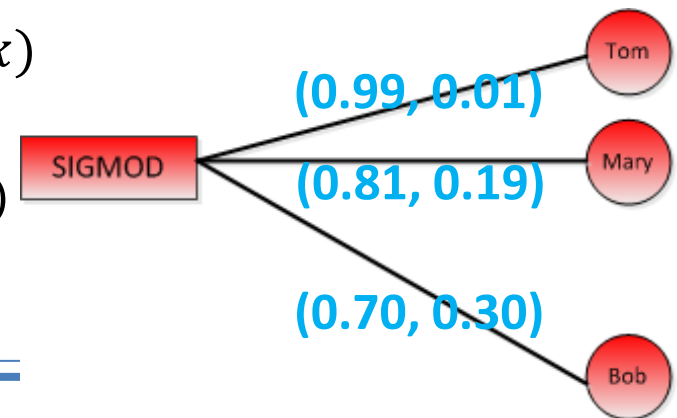
$$\vec{r}_X(i) = \sum_{j=1}^n W_{XY}(i, j) \vec{r}_Y(j)$$

- Rule 3: The rank of an author is enhanced if he or she co-authors with *many* highly ranked authors

$$\vec{r}_Y(i) = \alpha \sum_{j=1}^m W_{YX}(i, j) \vec{r}_X(j) + (1 - \alpha) \sum_{j=1}^n W_{YY}(i, j) \vec{r}_Y(j)$$

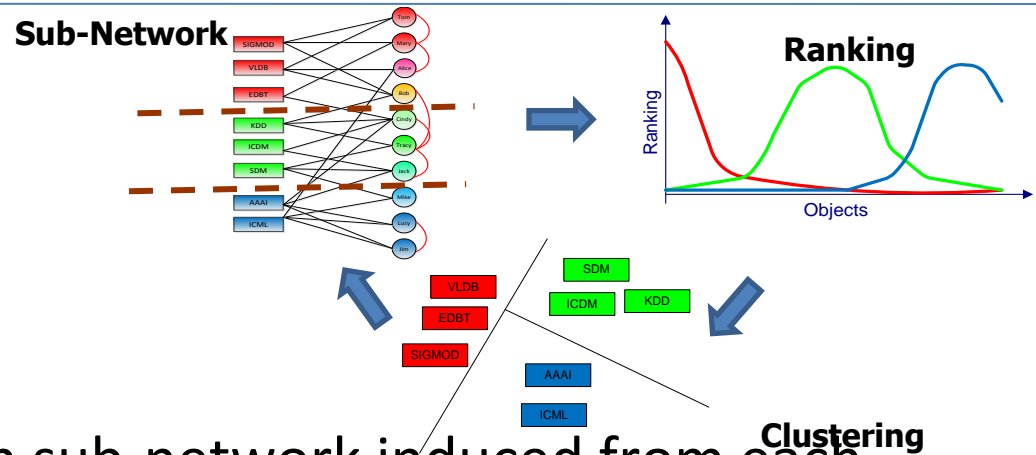
# Generating New Measure Space

- Input: Conditional ranking distributions for each cluster
  - $P_X(i|k)$ : e.g.,  $P_X(\text{SIGMOD}|\text{area} = \text{"database"})$
- Output: Each conference  $i$  is mapped into a new measure space
  - $i: (\pi_{i,1}, \dots, \pi_{i,K})$ , where  $\pi_{i,k} = P_X(k|i)$ 
    - E.g., SIGMOD:  $(P(\text{"database"}|\text{SIGMOD}), P(\text{"hardware"}|\text{SIGMOD}))$
- Solution
  - $P_X(k|i) \propto P(k) \times P_X(i|k)$
  - Calculate cluster size  $P(k)$ 
    - Maximize the log-likelihood of generating all the links
      - $P(i,j) = \sum_k P(k) \times P_X(i|k) \times P_Y(j|k)$
    - EM algorithm
      - $P(k|i,j) \propto P(k) \times P_X(i|k) \times P_Y(j|k)$
      - $P(k) \propto \sum_{ij} W_{XY}(i,j)P(k|i,j)$



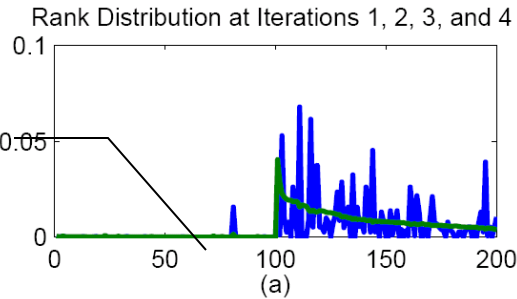
# The Algorithm Framework

- Step 0: Initialization
  - Randomly partition
- Step 1: Ranking
  - Ranking objects in each sub-network induced from each cluster
- Step 2: Generating new measure space
  - Estimate **mixture model coefficients** for each target object
- Step 3: Adjusting cluster
- Step 4: Repeating Steps 1-3 until stable

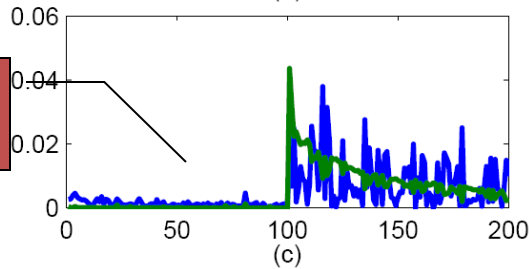


# Step-by-Step Running Case Illustration

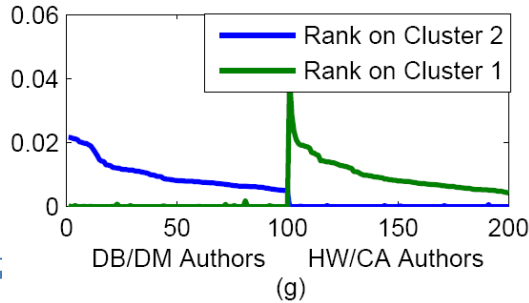
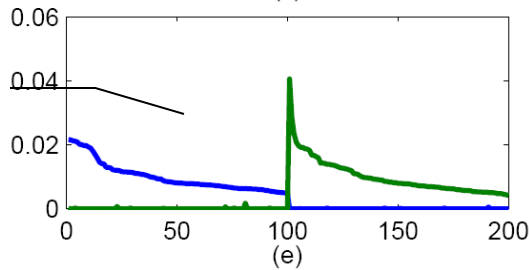
Initially, ranking distributions are mixed together



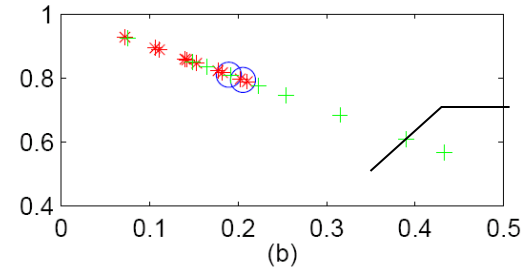
Improved a little



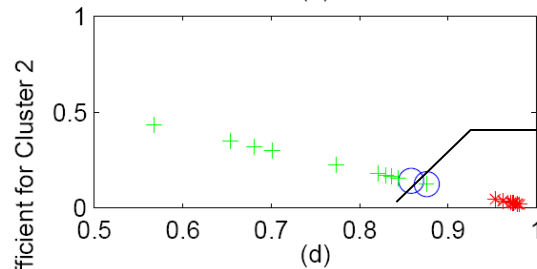
Improved significantly



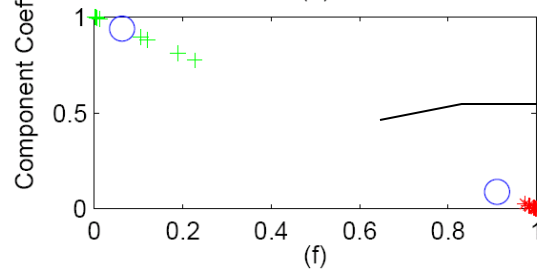
Scatter Plot for Conf. at Iterations 1, 2, 3, and 4



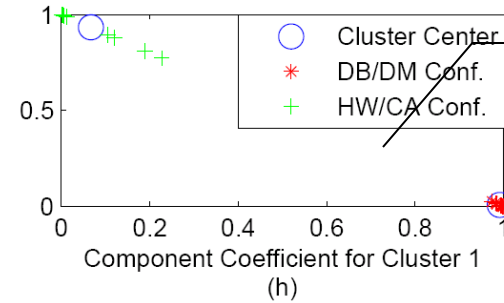
Two clusters of objects mixed together, but preserve similarity somehow



Two clusters are almost well separated



Well separated

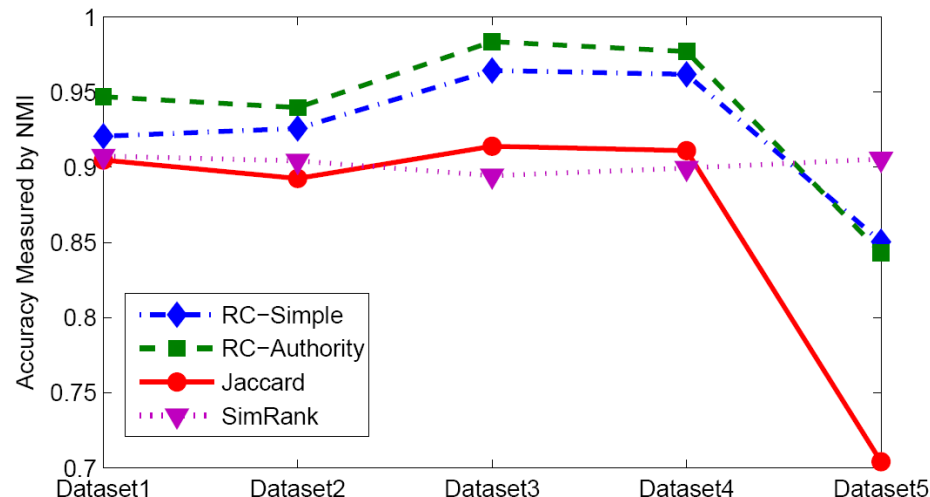


Stable

# Clustering and Ranking CS Conferences by RankClus

	DB	Network	AI	Theory	IR
1	VLDB	INFOCOM	AAMAS	SODA	SIGIR
2	ICDE	SIGMETRICS	IJCAI	STOC	ACM Multimedia
3	SIGMOD	ICNP	AAAI	FOCS	CIKM
4	KDD	SIGCOMM	Agents	ICALP	TREC
5	ICDM	MOBICOM	AAAI/IAAI	CCC	JCDL
6	EDBT	ICDCS	ECAI	SPAA	CLEF
7	DASFAA	NETWORKING	RoboCup	PODC	WWW
8	PODS	MobiHoc	IAT	CRYPTO	ECDL
9	SSDBM	ISCC	ICMAS	APPROX-RANDOM	ECIR
10	SDM	SenSys	CP	EUROCRYPT	CIVR

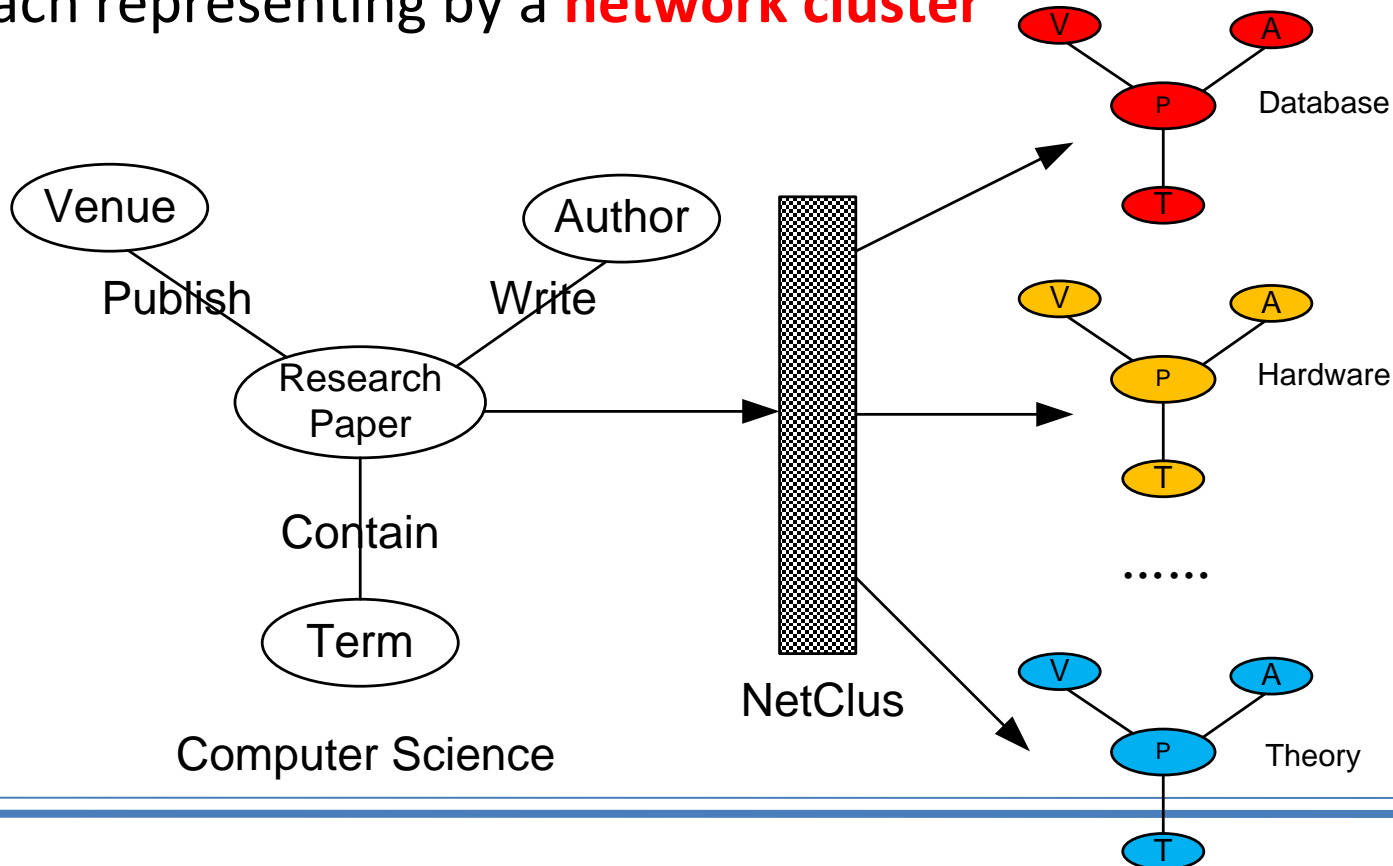
**Top-10 conferences in 5 clusters using RankClus in DBLP**



**RankClus outperforms *spectral clustering* [Shi and Malik, 2000] algorithms on projected homogeneous networks**

# NetClus [Sun, KDD'09]: Beyond Bi-Typed Networks

- Beyond bi-typed information network
  - A Star Network Schema [**richer information**]
- Split a network into different layers
  - Each representing by a **network cluster**





# Multi-Typed Networks Lead to Better Results

- The network cluster for **database area**: Conferences, Authors, and Terms
  - Better clustering and ranking than RankClus

Conference	Rank Score	Author	Rank Score	Term	Rank Score
SIGMOD	0.315	Michael Stonebraker	0.0063	database	0.0529
VLDB	0.306	Surajit Chaudhuri	0.0057	system	0.0322
ICDE	0.194	C. Mohan	0.0053	query	0.0313
PODS	0.109	Michael J. Carey	0.0052	data	0.0251
EDBT	0.046	David J. DeWitt	0.0051	object	0.0138
CIKM	0.019	H. V. Jagadish	0.0043	management	0.0113
...	...	...	...	...	...

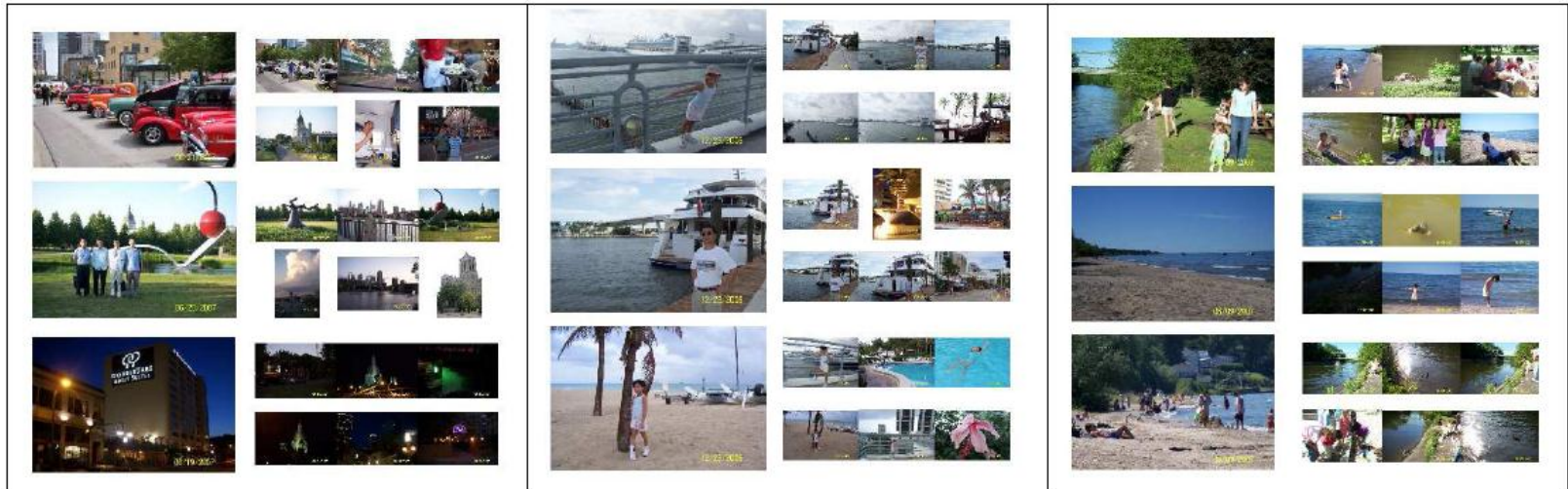
- NetClus vs. RankClus: **16%** higher accuracy on conference clustering in terms of Normalized Mutual Information

# Impact of RankClus Methodology

---

- RankCompete [Cao et al., WWW'10]
  - Extend to the domain of web images
- RankClus in Medical Literature [Li et al., Working paper]
  - Ranking treatments for diseases
- RankClass [Ji et al., KDD'11]
  - Integrate classification with ranking
- Trustworthy Analysis [Gupta et al., WWW'11] [Khac Le et al., IPSN'11]
  - Integrate clustering with trustworthiness score
- Topic Modeling in Heterogeneous Networks [Deng et al., KDD'11]
  - Propagate topic information among different types of objects
- ...

# Interesting Results from Other Domains




## RankCompete: Organize images automatically!

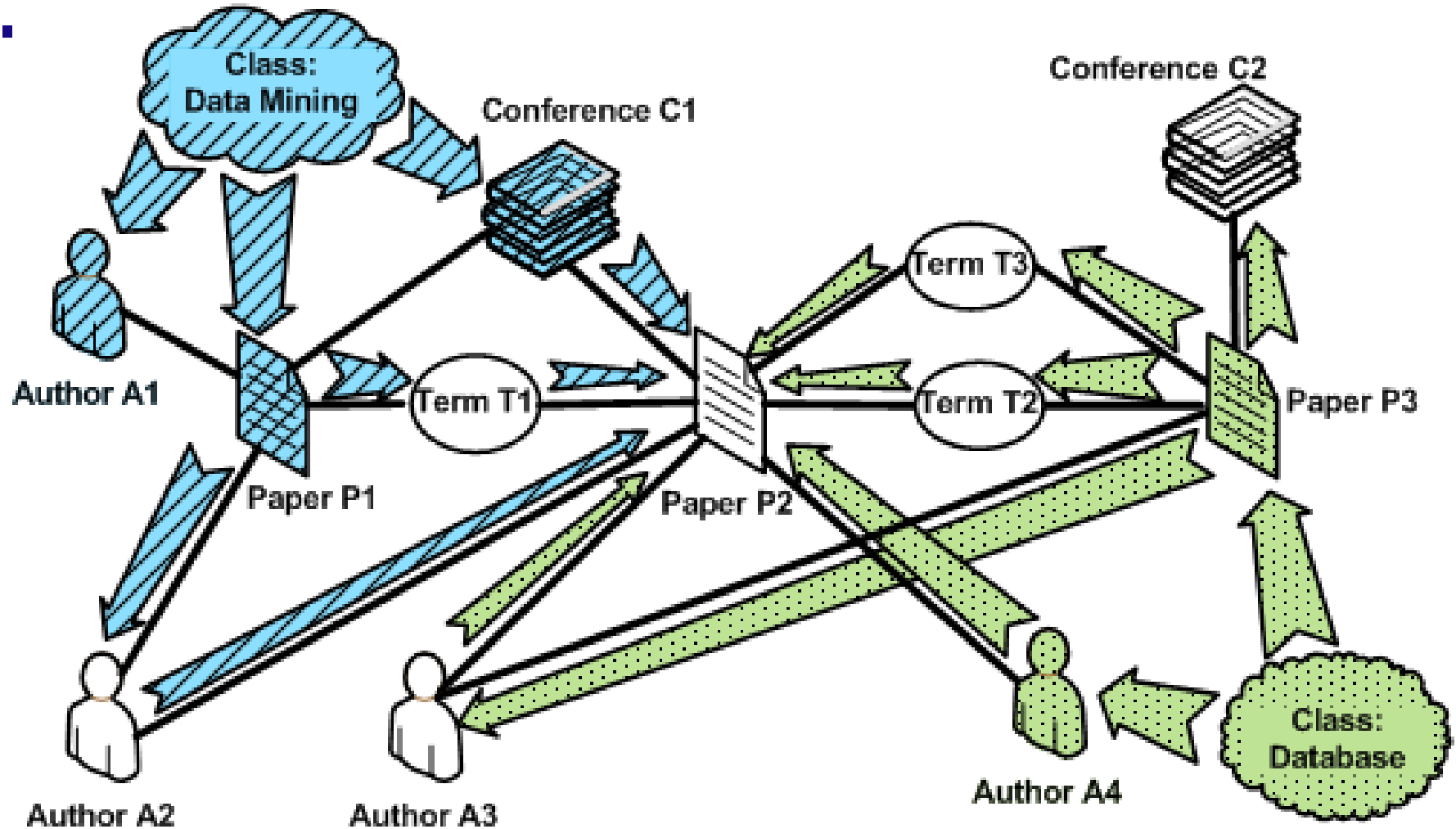
	Top 10 Treatments	Ranking
1	Zidovudine/therapeutic use	0.1679
2	Anti-HIV Agents/therapeutic use	0.1340
3	Antiretroviral Therapy, Highly Active	0.0977
4	Antiviral Agents/therapeutic use	0.0718
5	Anti-Retroviral Agents/therapeutic use	0.0236
6	Interferon Type I/therapeutic use	0.0147
7	Didanosine/therapeutic use	0.0132
8	Ganciclovir/therapeutic use	0.0114
9	HIV Protease Inhibitors/therapeutic use	0.0105
10	Antineoplastic Combined Chemotherapy	0.0103

# Outline

---

- **Motivation:** Why Mining Information Networks?
  - **Part I:** Clustering, Ranking and Classification
    - Clustering and Ranking in Information Networks
    - Classification of Information Networks 
  - **Part II:** Meta-Path Based Exploration of Information Networks
    - Similarity Search in Information Networks
    - Relationship Prediction in Information Networks
  - **Part III:** Advanced Topics on Information Network Analysis
    - Role Discovery and OLAP in Information Networks
    - Relation Strength Learning in Information Networks
    - Mining Evolution and Dynamics of Information Networks
  - **Conclusions**
-

# Classification: Knowledge Propagation



M. Ji, M. Danilevski, et al., "Graph Regularized Transductive Classification on Heterogeneous Information Networks", ECMLPKDD'10

# GNetMine: Graph-Based Regularization [Ji, PKDD'10]

- Minimize the objective function

$$J(\mathbf{f}_1^{(k)}, \dots, \mathbf{f}_m^{(k)})$$

User preference: how much do you value this relationship / ground truth?

$$= \sum_{i,j=1}^m \lambda_{ij} \sum_{p=1}^{n_i} \sum_{q=1}^{n_j} R_{ij,pq} \left( \frac{1}{\sqrt{D_{ij,pp}}} f_{ip}^{(k)} - \frac{1}{\sqrt{D_{ji,qq}}} f_{jq}^{(k)} \right)^2$$

$$+ \sum_{i=1}^m \alpha_i (\mathbf{f}_i^{(k)} - \mathbf{y}_i^{(k)})^T (\mathbf{f}_i^{(k)} - \mathbf{y}_i^{(k)})$$

*Smoothness constraints:* objects linked together should share similar estimations of confidence belonging to class  $k$

Normalization term applied to each type of link separately:  
reduce the impact of popularity of nodes

Confidence estimation on labeled data and their pre-given labels should be similar

# From RankClus to GNetMine & RankClass

---

- ❑ **RankClus [EDBT'09]: Clustering and ranking working together**
  - ❑ No training, no available class labels, no expert knowledge
- ❑ **GNetMine [PKDD'10]: Incorp. prior knowledge in networks**
  - ❑ Classification in heterog. networks, but objects treated equally
- ❑ **RankClass [KDD'11]: Integration of ranking and classification in heterogeneous network analysis**
  - ❑ Ranking: informative understanding & summary of each class
  - ❑ Class membership is critical information when ranking objects
  - ❑ Let ranking and classification mutually enhance each other!
  - ❑ Output: Classification results + ranking list of objects within each class

# Experiments on DBLP

---

- ❑ Class: Four research areas (communities)
  - Database, data mining, AI, information retrieval
- ❑ Four types of objects
  - Paper (14376), Conf. (20), Author (14475), Term (8920)
- ❑ Three types of relations
  - Paper-conf., paper-author, paper-term
- ❑ Algorithms for comparison
  - Learning with Local and Global Consistency (LLGC) [Zhou et al. NIPS 2003] – also the homogeneous version of our method
  - Weighted-vote Relational Neighbor classifier (wvRN) [Macskassy et al. JMLR 2007]
  - Network-only Link-based Classification (nLB) [Lu et al. ICML 2003, Macskassy et al. JMLR 2007]



# Performance Study on the DBLP Data Set

Table 3: Comparison of classification accuracy on authors (%)

( $a\%$ , $p\%$ ) of authors and papers labeled	nLB (A-A)	nLB (A-C-P-T)	wvRN (A-A)	wvRN (A-C-P-T)	LLGC (A-A)	LLGC (A-C-P-T)	GNetMine (A-C-P-T)	RankClass (A-C-P-T)
(0.1%, 0.1%)	25.4	26.0	40.8	34.1	41.4	61.3	82.9	<b>83.9</b>
(0.2%, 0.2%)	28.3	26.0	46.0	41.2	44.7	62.2	83.4	<b>85.6</b>
(0.3%, 0.3%)	28.4	27.4	48.6	42.5	48.8	65.7	86.7	<b>88.3</b>
(0.4%, 0.4%)	30.7	26.7	46.3	45.6	48.7	66.0	87.2	<b>88.8</b>
(0.5%, 0.5%)	29.8	27.3	49.0	51.4	50.6	68.9	87.5	<b>89.2</b>
average	28.5	26.7	46.3	43.0	46.8	64.8	85.5	<b>87.2</b>

Table 4: Comparison of classification accuracy on papers (%)

( $a\%$ , $p\%$ ) of authors and papers labeled	nLB (P-P)	nLB (A-C-P-T)	wvRN (P-P)	wvRN (A-C-P-T)	LLGC (P-P)	LLGC (A-C-P-T)	GNetMine (A-C-P-T)	RankClass (A-C-P-T)
(0.1%, 0.1%)	49.8	31.5	62.0	42.0	67.2	62.7	<b>79.2</b>	77.7
(0.2%, 0.2%)	73.1	40.3	71.7	49.7	72.8	65.5	<b>83.5</b>	83.0
(0.3%, 0.3%)	77.9	35.4	77.9	54.3	76.8	66.6	83.2	<b>83.6</b>
(0.4%, 0.4%)	79.1	38.6	78.1	54.4	77.9	70.5	83.7	<b>84.7</b>
(0.5%, 0.5%)	80.7	39.3	77.9	53.5	79.0	73.5	84.1	<b>84.8</b>
average	72.1	37.0	73.5	50.8	74.7	67.8	82.7	<b>82.8</b>

Table 5: Comparison of classification accuracy on conferences (%)

( $a\%$ , $p\%$ ) of authors and papers labeled	nLB (A-C-P-T)	wvRN (A-C-P-T)	LLGC (A-C-P-T)	GNetMine (A-C-P-T)	RankClass (A-C-P-T)
(0.1%, 0.1%)	25.5	43.5	79.0	81.0	<b>84.5</b>
(0.2%, 0.2%)	22.5	56.0	83.5	85.0	<b>85.5</b>
(0.3%, 0.3%)	25.0	59.0	<b>87.0</b>	<b>87.0</b>	<b>87.0</b>
(0.4%, 0.4%)	25.0	57.0	86.5	89.5	<b>90.5</b>
(0.5%, 0.5%)	25.0	68.0	90.0	94.0	<b>95.0</b>
average	24.6	56.7	85.2	87.3	<b>88.5</b>


# Experiments with Very Small Training Set

- ❑ DBLP: 4-fields data set (DB, DM, AI, IR) forming a heterog. info. network
- ❑ Rank objects within each class (with extremely limited label information)
- ❑ Obtain High classification accuracy and excellent rankings within each class

	Database	Data Mining	AI	IR
Top-5 ranked conferences	VLDB	KDD	IJCAI	SIGIR
	SIGMOD	SDM	AAAI	ECIR
	ICDE	ICDM	ICML	CIKM
	PODS	PKDD	CVPR	WWW
	EDBT	PAKDD	ECML	WSDM
Top-5 ranked terms	data	mining	learning	retrieval
	database	data	knowledge	information
	query	clustering	reasoning	web
	system	classification	logic	search
	xml	frequent	cognition	text

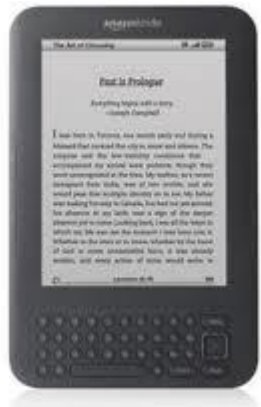
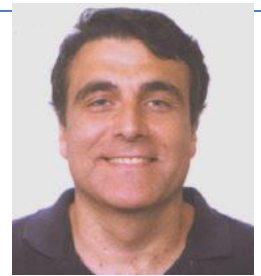
# Outline

---

- **Motivation:** Why Mining Information Networks?
  - **Part I:** Clustering, Ranking and Classification
    - Clustering and Ranking in Information Networks
    - Classification of Information Networks
  - **Part II:** Meta-Path Based Exploration of Information Networks
    - Similarity Search in Information Networks 
    - Relationship Prediction in Information Networks
  - **Part III:** Advanced Topics on Information Network Analysis
    - Role Discovery and OLAP in Information Networks
    - Relation Strength Learning in Information Networks
    - Mining Evolution and Dynamics of Information Networks
  - **Conclusions**
-

# Similarity Search: Find Similar Objects in Networks [Sun, VLDB'11]

- DBLP
  - Who are the most similar to “Christos Faloutsos”?
- IMDB
  - Which movies are the most similar to “Little Miss Sunshine”?
- E-Commerce
  - Which products are the most similar to “Kindle”?



**How to systematically answer these questions in heterogeneous information networks?**

# Existing Link-based Similarity Functions

---

- Existing similarity functions in networks
  - Personalized PageRank (P-PageRank) [Jeh and Widom, 2003]
  - SimRank [Jeh and Widom, 2002]
- Drawbacks
  - Do not distinguish object type and link type
  - Limitations on the similarity measures
    - To return highly visible objects or pure objects in the network

# Network Schema and Meta-Path

Objects are connected together via different types of relationships!

**“Jim-P1-Ann”**

**“Mike-P2-Ann”**

**“Mike-P3-Bob”**

*Author-Paper-Author*

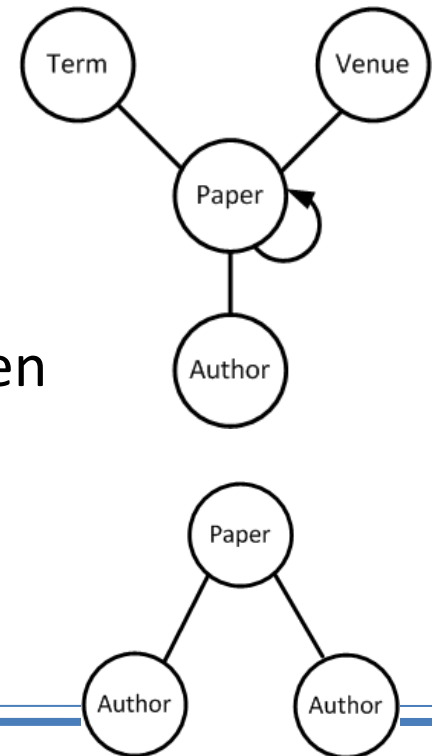
**“Jim-P1-SIGMOD-P2-Ann”**

**“Mike-P3-SIGMOD-P2-Ann”**

**“Mike-P4-KDD-P5-Bob”**

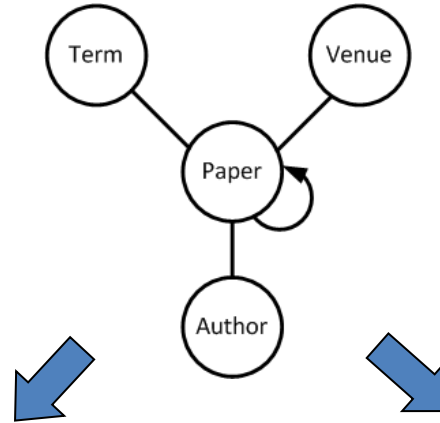
*Author-Paper-Venue-Paper-Author*

- Network schema
  - Meta-level description of a network
- Meta-Path
  - **Meta-level description** of a path between two objects
  - **A path** on network schema
  - Denote an existing or concatenated **relation** between two object types



# Different Meta-Paths Tell Different Semantics

- Who are most similar to Christos Faloutsos?



**Meta-Path: Author-Paper-Author**

Rank	Author	Score
1	Christos Faloutsos	1
2	Spiros Papadimitriou	0.127
3	Jimeng Sun	0.12
4	Jia-Yu Pan	0.114
5	Agma J. M. Traina	0.110
6	Jure Leskovec	0.096
7	Caetano Traina Jr.	0.096
8	Hanghang Tong	0.091
9	Deepayan Chakrabarti	0.083
10	Flip Korn	0.053

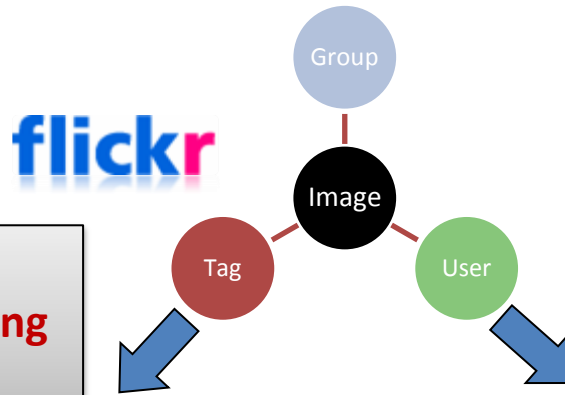
**Meta-Path: Author-Paper-Venue-Paper-Author**

Rank	Author	Score
1	Christos Faloutsos	1
2	Jiawei Han	0.842
3	Rakesh Agrawal	0.838
4	Jian Pei	0.8
5	Charu C. Aggarwal	0.739
6	H. V. Jagadish	0.705
7	Raghu Ramakrishnan	0.697
8	Nick Koudas	0.689
9	Surajit Chaudhuri	0.677
10	Divesh Srivastava	0.661

**Christos's students or close collaborators**      **Work on similar topics and have similar reputation**

# Some Meta-Path Is “Better” Than Others

- Which pictures are most similar to  ?



**Evaluate the similarity  
between images according  
to their linked tags**

**Meta-Path: *Image-Tag-Image***



(a) top-1

(b) top-2

(c) top-3



(d) top-4

(e) top-5

(f) top-6

**Evaluate the similarity  
between images according  
to tags and groups**

**Meta-Path: *Image-Tag-Image-Group-Image-Tag-Image***



(a) top-1

(b) top-2

(c) top-3



(d) top-4

(e) top-5

(f) top-6

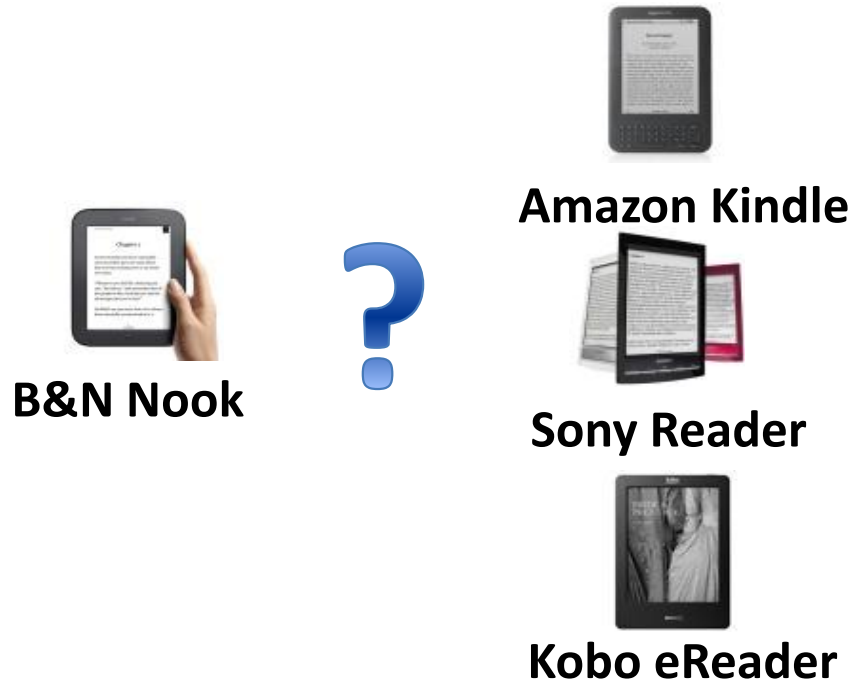




# PathSim: Similarity in Terms of “Peers”

---

- Why peers?
  - Strongly connected, while **similar visibility**

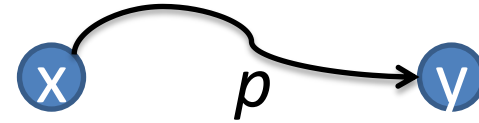


- In addition to meta-path
  - Need to consider **similarity measures**

# Limitations of Existing Similarity Measures

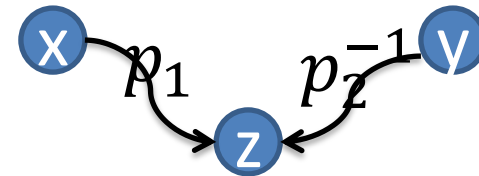
- Random walk (RW)

- $s(x, y) = \sum_{p \in \mathcal{P}} \text{Prob}(p)$
- Used in **Personalized PageRank (P-PageRank)**
- Favor **highly visible** objects
  - objects with large degrees

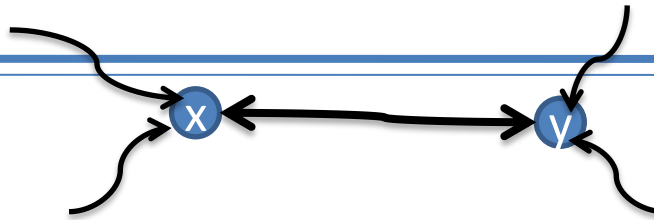


- Pairwise random walk (PRW)

- $s(x, y) = \sum_{(p_1, p_2) \in (\mathcal{P}_1, \mathcal{P}_2)} \text{Prob}(p_1) \text{Prob}(p_2^{-1})$
- Used in **SimRank**
- Favor **“pure”** objects
  - objects with highly skewed distribution in their in-links or out-links



# Only PathSim Can Find Peers



- PathSim

- Normalized path count between x and y following meta-path  $\mathcal{P}$

$$s(x, y) = \frac{2 \times |\{p_{x \rightsquigarrow y} : p_{x \rightsquigarrow y} \in \mathcal{P}\}|}{|\{p_{x \rightsquigarrow x} : p_{x \rightsquigarrow x} \in \mathcal{P}\}| + |\{p_{y \rightsquigarrow y} : p_{y \rightsquigarrow y} \in \mathcal{P}\}|}$$

Visibility of x

Visibility of y

- Favor **“peers”**:

- objects with strong connectivity and similar visibility under the given meta-path

- Calculation

- For  $\mathcal{P}: A_1 - A_2 - \dots - A_l - A_{l-1} - \dots - A_1$

- $M = W_{A_1 A_2} W_{A_2 A_3} \dots W_{A_{l-1} A_l} W_{A_l A_{l-1}} \dots W_{A_3 A_2} W_{A_2 A_1}$

- $s(x, y) = \frac{2M_{xy}}{M_{xx} + M_{yy}}$

- A co-clustering based pruning algorithm is provided
  - » 18.23% - 68.04% efficiency improvement over the baseline

# Properties of PathSim

---

- Symmetric
  - $s(x, y) = s(y, x)$
- Self-Maximum
  - $s(x, y) \in [0, 1]$ , and  $s(x, x) = 1$
- Balance of visibility
  - $$s(x, y) \leq \frac{2}{\sqrt{M_{xx}/M_{yy}} + \sqrt{M_{yy}/M_{xx}}}$$
    - $M_{xx}$  is the number of path instances starting from  $x$  and ending with  $x$  following the given meta path
- Limiting behavior
  - If repeating a pattern of meta path infinite times, PathSim degenerates to authority ranking comparison

**Long meta path without introducing new relationships is not that helpful!**

# Find Academic Peers by PathSim

- Anhai Doan

- CS, Wisconsin
- Database area
- PhD: 2002



- Jignesh Patel

- CS, Wisconsin
- Database area
- PhD: 1998

**Meta-Path: *Author-Paper-Venue-Paper-Author*** ↓

Rank	P-PageRank	SimRank	PathSim
1	AnHai Doan	AnHai Doan	AnHai Doan
2	Philip S. Yu	Douglas W. Cornell	<u>Jignesh M. Patel</u>
3	Jiawei Han	Adam Silberstein	<u>Amol Deshpande</u>
4	Hector Garcia-Molina	Samuel DeFazio	<u>Jun Yang</u>
5	Gerhard Weikum	Curt Ellmann	Renée J. Miller



- Amol Deshpande

- CS, Maryland
- Database area
- PhD: 2004



- Jun Yang

- CS, Duke
- Database area
- PhD: 2001


# Meta-Path: A Key Concept for Mining Heterogeneous Networks

---

- Meta-path based mining
  - PathPredict [Sun et al., ASONAM'11]
    - Co-authorship prediction using meta-path based similarity
  - PathPredict\_when [Sun et al., WSDM'12]
    - When a relationship will happen
  - Citation prediction [Yu et al., SDM'12]
    - Meta-path + topic
- Meta-path learning
  - User Guided Meta-Path Selection [Sun et al., KDD'12 Submission]
    - Meta-path selection + clustering

# Outline

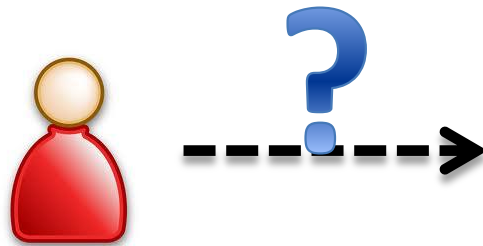
---

- **Motivation:** Why Mining Information Networks?
  - **Part I:** Clustering, Ranking and Classification
    - Clustering and Ranking in Information Networks
    - Classification of Information Networks
  - **Part II:** Meta-Path Based Exploration of Information Networks
    - Similarity Search in Information Networks
    - Relationship Prediction in Information Networks 
  - **Part III:** Advanced Topics on Information Network Analysis
    - Role Discovery and OLAP in Information Networks
    - Relation Strength Learning in Information Networks
    - Mining Evolution and Dynamics of Information Networks
  - **Conclusions**
-

# PathPredict: Meta-Path Based Relationship Prediction

---

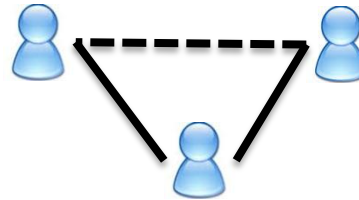
- Wide applications
  - Whom should I **collaborate** with?
  - Which paper should I **cite** for this topic?
  - Whom else should I **follow** on Twitter?
  - Whether Ann will **buy** the book “Steve Jobs”?
  - Whether Bob will **click** the ad on hotel?
  - ...





# Relationship Prediction vs. Link Prediction

- Link prediction in homogeneous networks [Liben-Nowell and Kleinberg, 2003, Hasan et al., 2006]
  - E.g., friendship prediction



- Relationship prediction in heterogeneous networks
  - **Target:** Different types of relationships need different prediction models

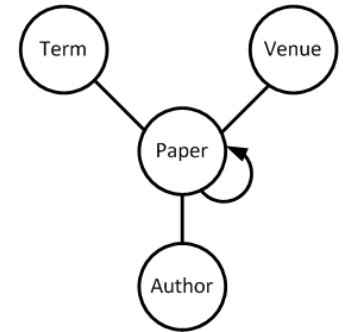


- **Features:** Different connection paths need to be treated separately!
  - **Meta-path based approach** to define topological features.



# PathPredict: Meta-Path Based Co-authorship Prediction in DBLP [Sun, ASONAM'11]

- Co-authorship prediction problem
  - Whether two authors are going to collaborate for the first time
- Co-authorship encoded in meta-path
  - Author-Paper-Author
- Topological features encoded in meta-paths



Meta-Path	Semantic Meaning
$A - P \rightarrow P - A$	$a_i$ cites $a_j$
$A - P \leftarrow P - A$	$a_i$ is cited by $a_j$
$A - P - V - P - A$	$a_i$ and $a_j$ publish in the same venues
$A - P - A - P - A$	$a_i$ and $a_j$ are co-authors of the same authors
$A - P - T - P - A$	$a_i$ and $a_j$ write the same topics
$A - P \rightarrow P \rightarrow P - A$	$a_i$ cites papers that cite $a_j$
$A - P \leftarrow P \leftarrow P - A$	$a_i$ is cited by papers that are cited by $a_j$
$A - P \rightarrow P \leftarrow P - A$	$a_i$ and $a_j$ cite the same papers
$A - P \leftarrow P \rightarrow P - A$	$a_i$ and $a_j$ are cited by the same papers

Meta-paths between authors under length 4

# The Power of PathPredict

- Explain the prediction power of each meta-path

- Wald Test for logistic regression

Social relations play very important role?

- Higher prediction accuracy than using projected homogeneous network

- **11%** higher in prediction accuracy

Meta Path	p-value	significance level <sup>1</sup>
$A - P \rightarrow P - A$	0.0378	**
$A - P \leftarrow P - A$	0.0077	***
$A - P - V - P - A$	1.2974e-174	****
$A - P - A - P - A$	1.1484e-126	****
$A - P - T - P - A$	3.4867e-51	****
$A - P \rightarrow P \rightarrow P - A$	0.7459	
$A - P \leftarrow P \leftarrow P - A$	0.0647	*
$A - P \rightarrow P \leftarrow P - A$	9.7641e-11	****
$A - P \leftarrow P \rightarrow P - A$	0.0966	*

<sup>1</sup> \*:  $p < 0.1$ ; \*\*:  $p < 0.05$ ; \*\*\*:  $p < 0.01$ , \*\*\*\*:  $p < 0.001$

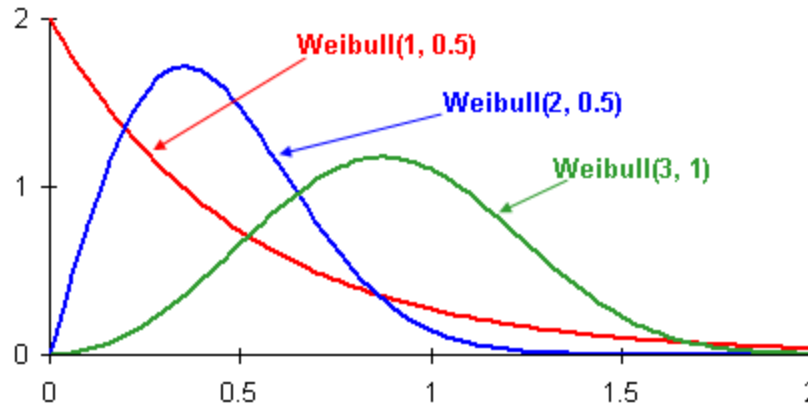
Rank	Hybrid heterogeneous features	# Shared authors
1	<b>Philip S. Yu</b>	<b>Philip S. Yu</b>
2	<b>Raymond T. Ng</b>	Ming-Syan Chen
3	Osmar R. Zaiane	Divesh Srivastava
4	<b>Ling Feng</b>	Kotagiri Ramamohanarao
5	<b>David Wai-Lok Cheung</b>	<b>Jeffrey Xu Yu</b>

**Co-author prediction for Jian Pei: Only 42 among 4809 candidates are true first-time co-authors!**  
(Feature collected in [1996, 2002]; Test period in [2003,2009])

# When Will It Happen? [Sun, WSDM'12]

- From “whether” to “when”
  - “Whether”: Will *Jim* rent the movie “*Avatar*” in Netflix?
  - “When”: When will *Jim* rent the movie “*Avatar*”?

Output:  $P(X=1)=?$



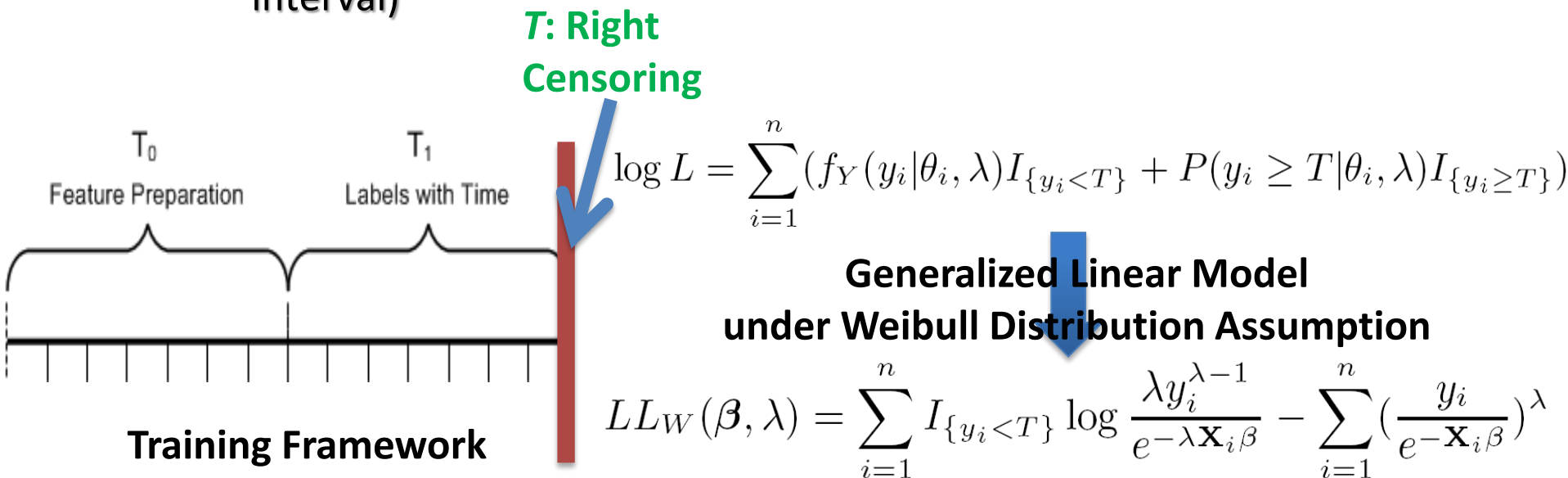
Output: A distribution of time!

- What is the probability Jim will rent “Avatar” *within 2 months*?
  - $P(Y \leq 2)$
- By when** Jim will rent “Avatar” with 90% probability?
  - $t: P(Y \leq t) = 0.9$
- What is the **expected time** it will take for Jim to rent “Avatar”?
  - $E(Y)$

May provide useful information to supply chain management

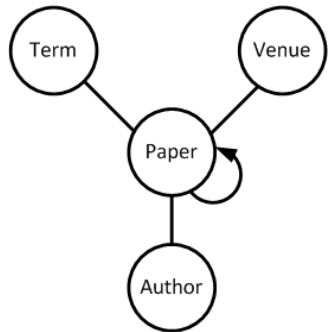
# The Relationship Building Time Prediction Model

- Solution
  - Directly **model relationship building time**:  $P(Y=t)$ 
    - Geometric distribution, Exponential distribution, Weibull distribution
  - Use **generalized linear model**
    - Deal with censoring (relationship builds beyond the observed time interval)



# Author Citation Time Prediction in DBLP

- Top-4 meta-paths for author citation time prediction



$A - P - T - P - A$

Study the same topic

$A - P \leftarrow P \rightarrow P - A$

Co-cited by the same paper

$A - P - A - P \rightarrow P - A$

Follow co-authors' citation

$A - P - T - P - A - P \rightarrow P - A$

Follow the citations of authors who study the same topic

**Social relations are less important in author citation prediction than in co-author prediction.**


- Predict when Philip S. Yu will cite a new author

$a_i$	$a_j$	Ground Truth	Median	Mean	25% quantile	75% quantile
Philip S. Yu	Ling Liu	1	2.2386	3.4511	0.8549	4.7370
Philip S. Yu	Christian S. Jensen	3	2.7840	4.2919	1.0757	5.8911
Philip S. Yu	C. Lee Giles	0	8.3985	12.9474	3.2450	17.7717
Philip S. Yu	Stefano Ceri	0	0.5729	0.8833	0.2214	1.2124
Philip S. Yu	David Maier	9+	2.5675	3.9581	0.9920	5.4329
Philip S. Yu	Tong Zhang	9+	9.5371	14.7028	3.6849	20.1811
Philip S. Yu	Rudi Studer	9+	9.7752	15.0698	3.7769	20.6849

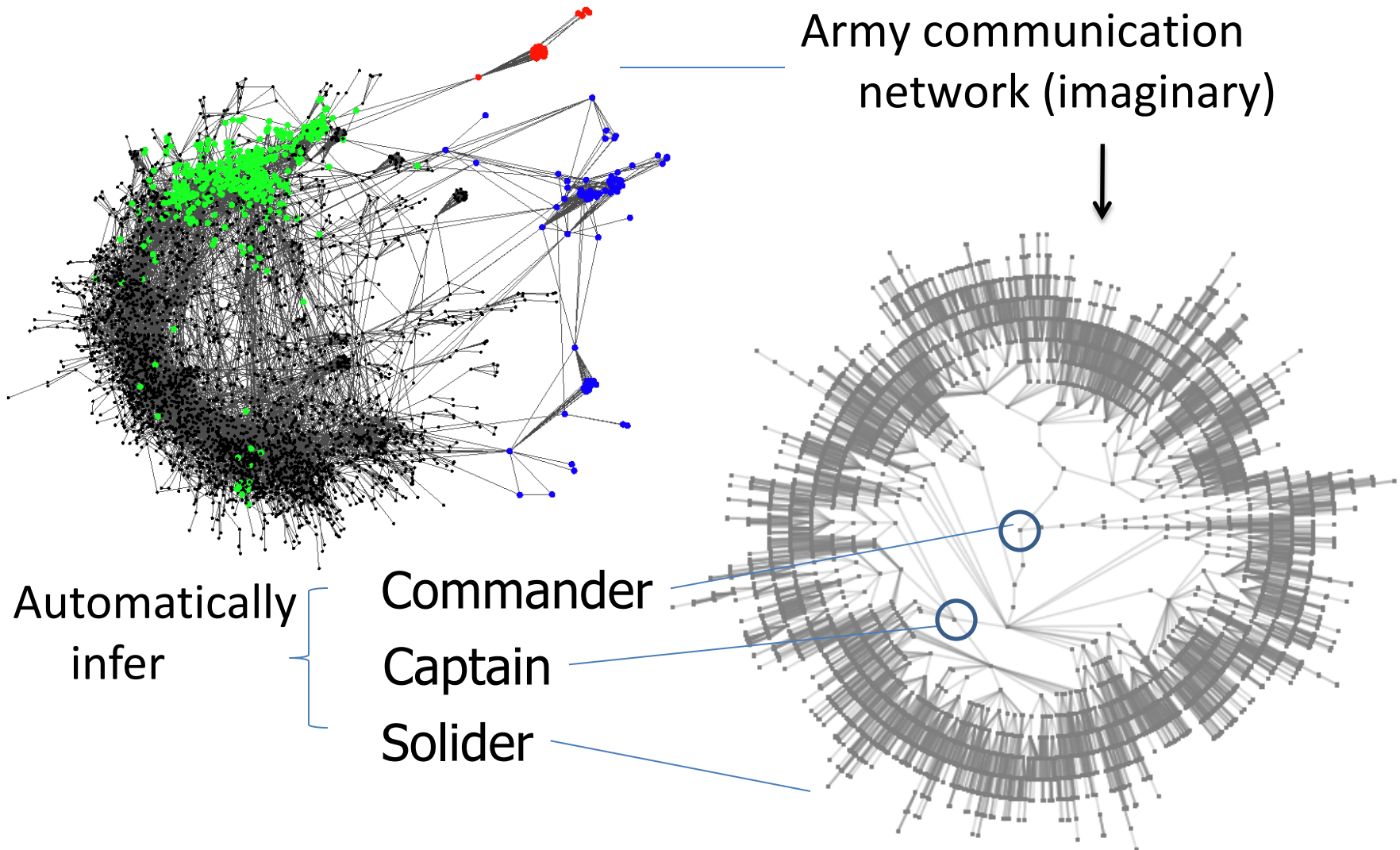
**Under Weibull distribution assumption**

# Outline

---

- **Motivation:** Why Mining Information Networks?
  - **Part I:** Clustering, Ranking and Classification
    - Clustering and Ranking in Information Networks
    - Classification of Information Networks
  - **Part II:** Meta-Path Based Exploration of Information Networks
    - Similarity Search in Information Networks
    - Relationship Prediction in Information Networks
  - **Part III:** Advanced Topics on Information Network Analysis
    - Role Discovery and OLAP in Information Networks 
    - Relation Strength Learning in Information Networks
    - Mining Evolution and Dynamics of Information Networks
  - **Conclusions**
-

# Role Discovery in Network: Why It Matters?





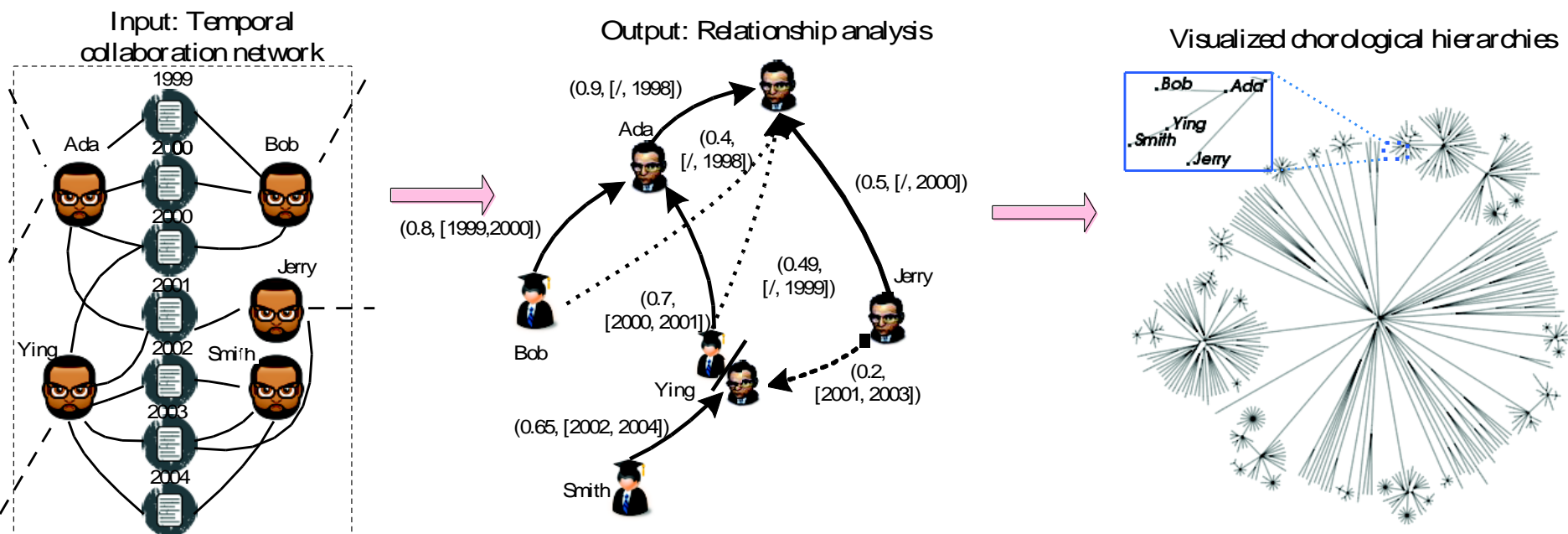
# Role Discovery: Extraction Semantic Information from Links

---

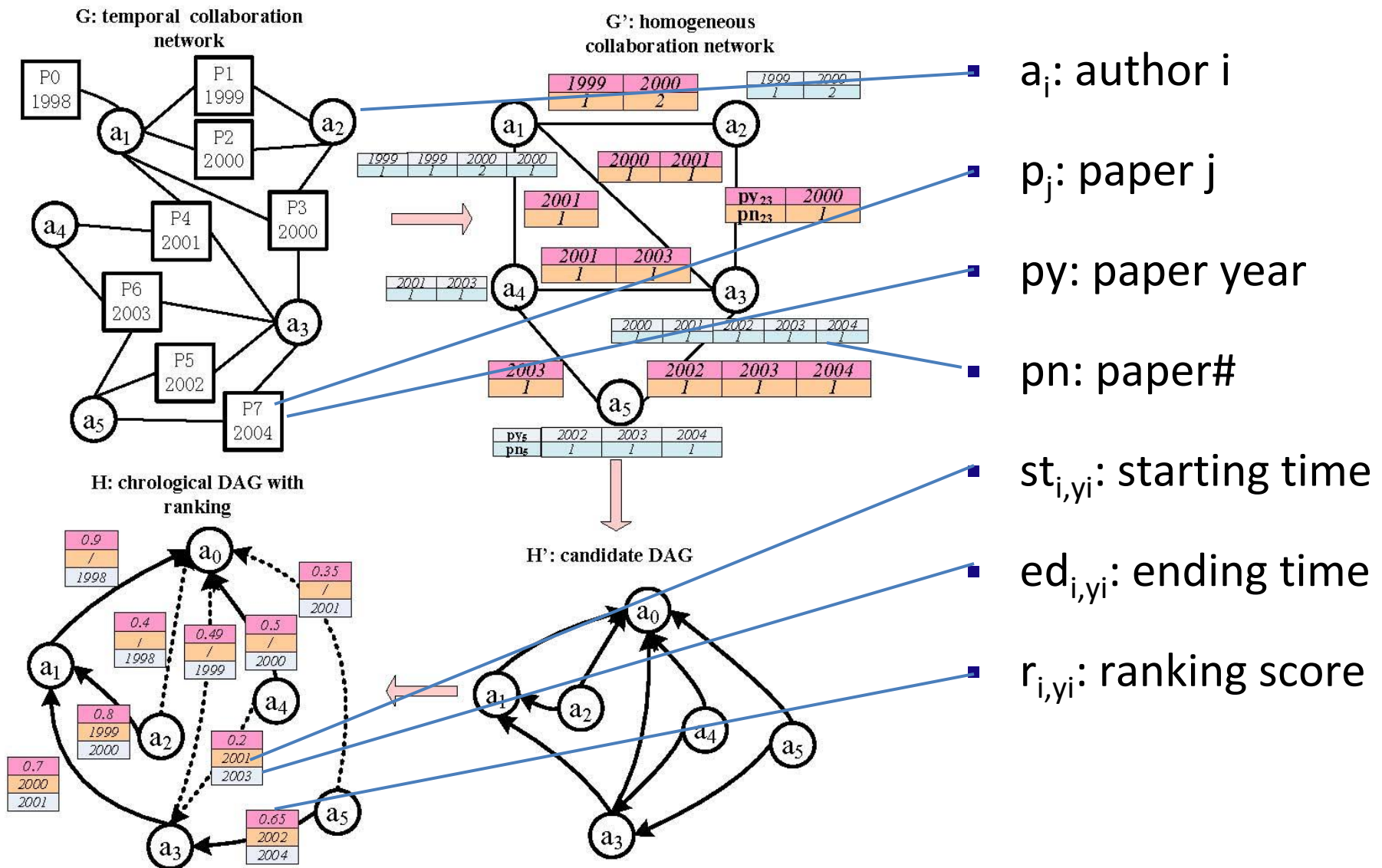
- Objective: Extract semantic meaning from plain links to finely model and better organize information networks
- Challenges
  - Latent semantic knowledge
  - Interdependency
  - Scalability
- Opportunity
  - Human intuition
  - Realistic constraint
  - Crosscheck with collective intelligence
- Methodology: propagate simple intuitive rules and constraints over the whole network

# Discovery of Advisor-Advisee Relationships in DBLP Network [Wang, KDD'10]

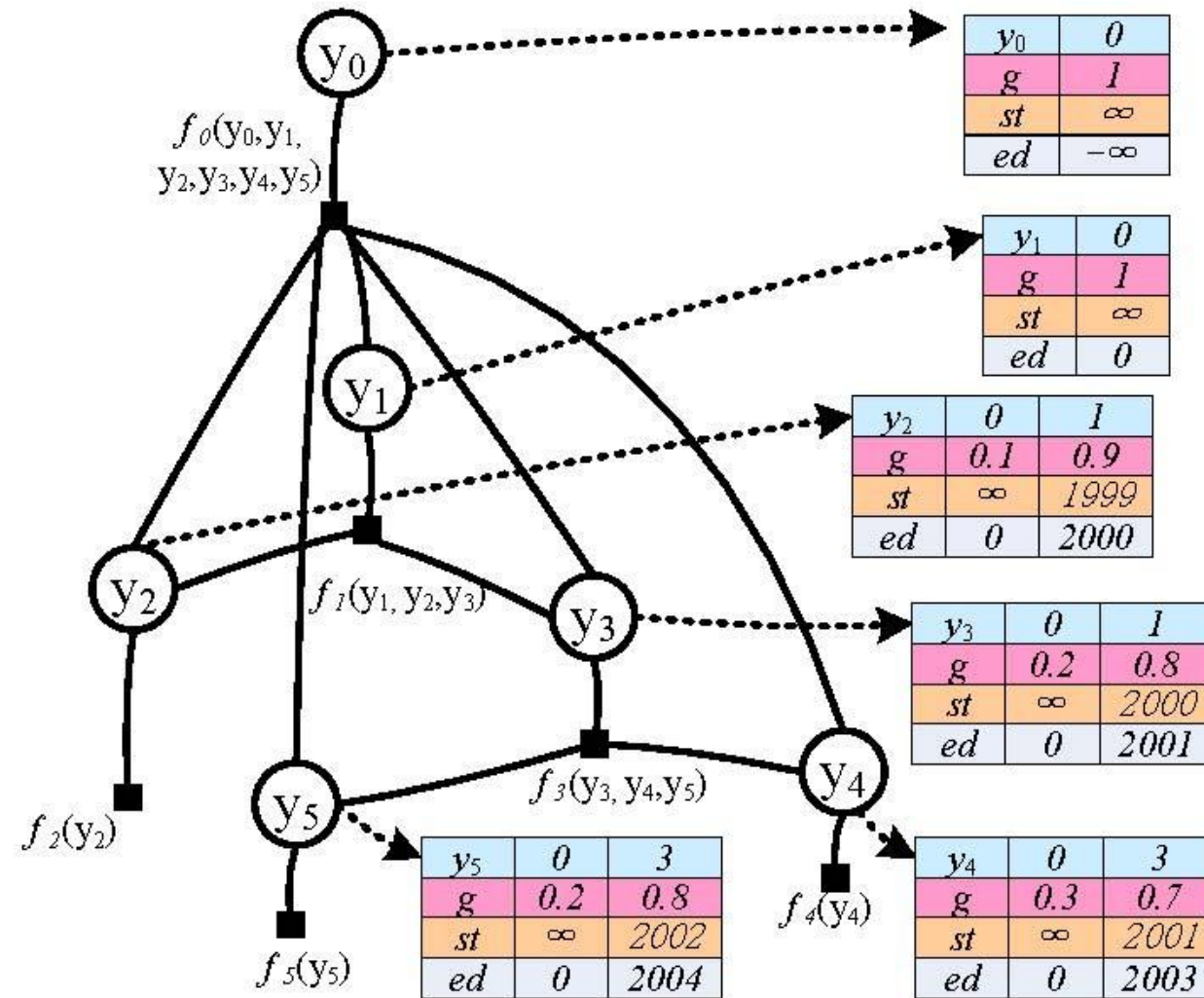
- Input: DBLP research publication network
- Output: Potential advising relationship and its ranking  $(r, [st, ed])$
- Ref. C. Wang, J. Han, et al., “Mining Advisor-Advisee Relationships from Research Publication Networks”, SIGKDD 2010



# Overall Framework



# Time-Constrained Probabilistic Factor Graph (TPFG)

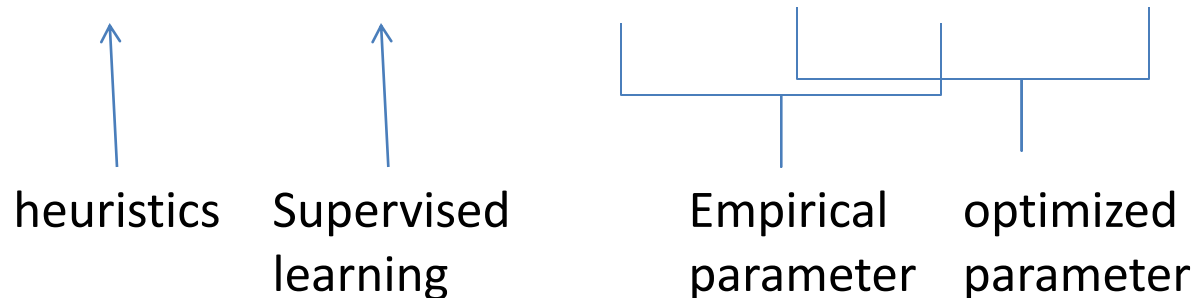


- $y_x$ :  $a_x$ 's advisor
- $st_{x,yx}$ : starting time  
 $ed_{x,yx}$ : ending time
- $g(y_x, st_x, ed_x)$  is predefined local feature
- $f_x(y_x, Z_x) = \max g(y_x, st_x, ed_x)$  under time constraint
- Objective function  $P(\{y_x\}) = \prod_x f_x(y_x, Z)$
- $Z = \{z \mid x \in Y_z\}$
- $Y_x$ : set of potential advisors of  $a_x$

# Experiment Results

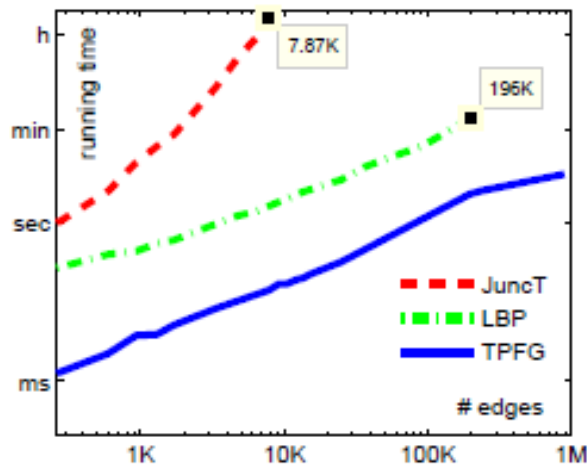
- DBLP data: 654, 628 authors, 1076,946 publications, years provided
- Labeled data: Math Genealogy Project; AI Genealogy Project; Homepage

Datasets	RULE	SVM	IndMAX		TPFG	
TEST1	69.9%	73.4%	75.2%	78.9%	80.2%	<b>84.4%</b>
TEST2	69.8%	74.6%	74.6%	79.0%	81.5%	<b>84.3%</b>
TEST3	80.6%	86.7%	83.1%	90.9%	88.8%	<b>91.3%</b>

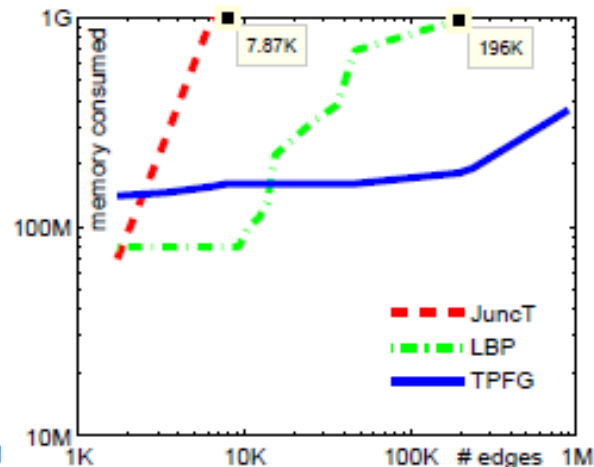


# Case Study & Scalability

Advisee	Top Ranked Advisor	Time	Note
David M. Blei	1. Michael I. Jordan	01-03	PhD advisor, 2004 grad
	2. John D. Lafferty	05-06	Postdoc, 2006
Hong Cheng	1. Qiang Yang	02-03	MS advisor, 2003
	2. Jiawei Han	04-08	PhD advisor, 2008
Sergey Brin	1. Rajeev Motawani	97-98	“Unofficial advisor”




(a) Time



(b) Space

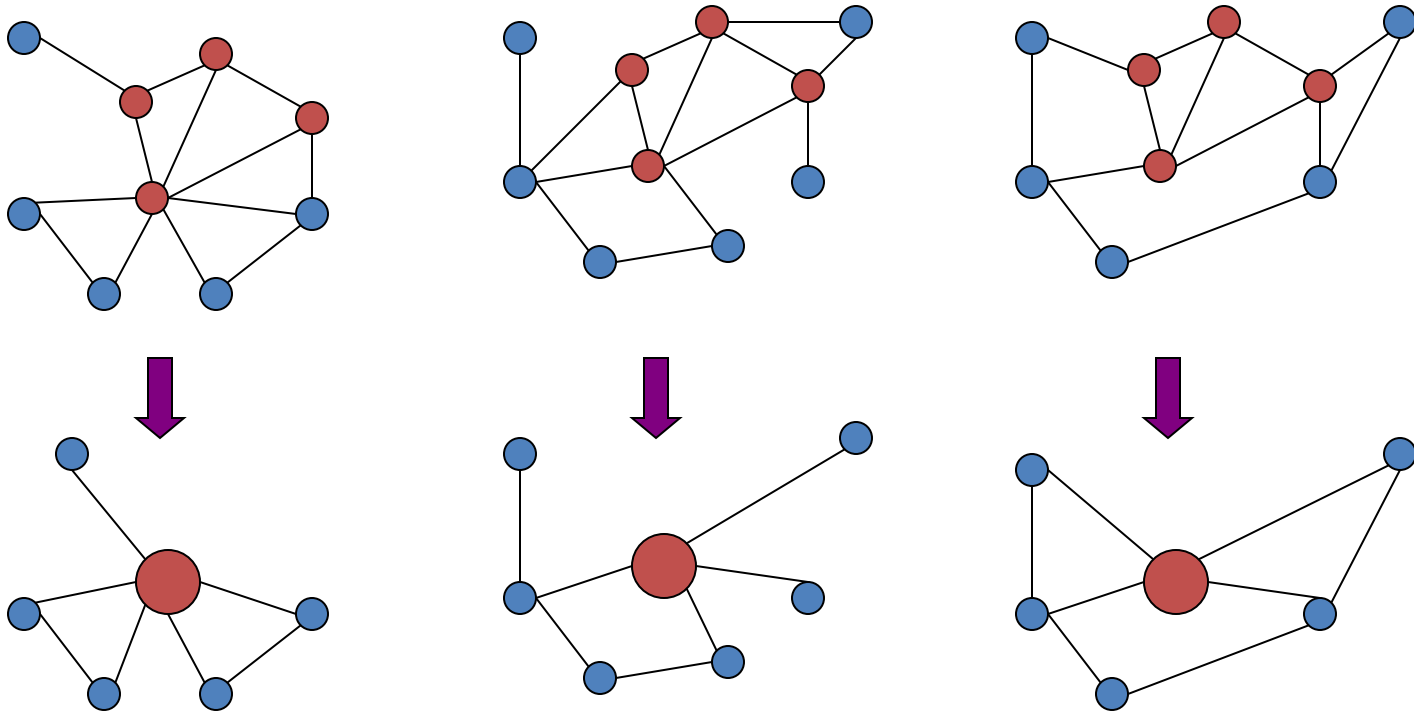
# Outline

---

- **Motivation:** Why Mining Information Networks?
  - **Part I:** Clustering, Ranking and Classification
    - Clustering and Ranking in Information Networks
    - Classification of Information Networks
  - **Part II:** Meta-Path Based Exploration of Information Networks
    - Similarity Search in Information Networks
    - Relationship Prediction in Information Networks
  - **Part III:** Advanced Topics on Information Network Analysis
    - Role Discovery and OLAP in Information Networks 
    - Relation Strength Learning in Information Networks
    - Mining Evolution and Dynamics of Information Networks
  - **Conclusions**
-

# Graph/Network Summarization: Graph Compression

- Extract common subgraphs and simplify graphs by condensing these subgraphs into nodes





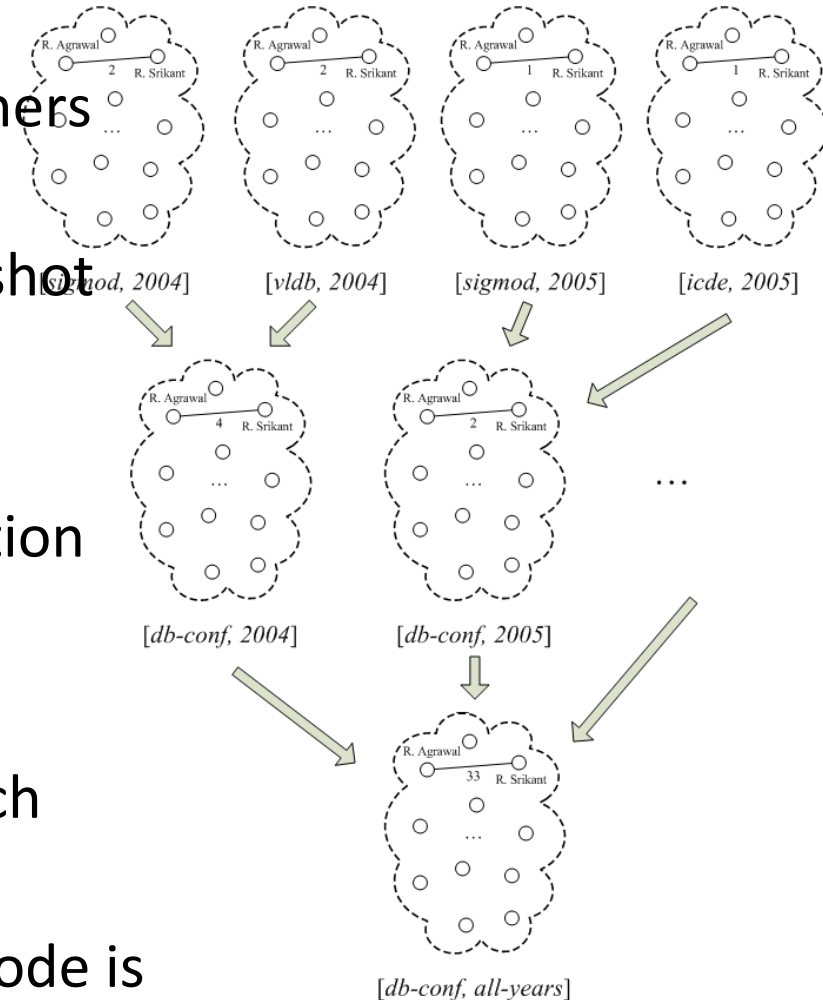
# OLAP on Information Networks [Chen, ICDM'08]

---

- Why OLAP information networks?
- Advantages of OLAP: Interactive exploration of multi-dimensional and multi-level space in a data cube Infonet
  - Multi-dimensional: Different perspectives
  - Multi-level: Different granularities
- InfoNet OLAP: Roll-up/drill-down and slice/dice on information network data
  - Traditional OLAP cannot handle this, because they ignore links among data objects
- Handling two kinds of InfoNet OLAP
  - Informational OLAP
  - Topological OLAP

# Informational OLAP

- In the DBLP network, study the collaboration patterns among researchers
- Dimensions come from informational attributes attached at the whole snapshot level, so-called *Info-Dims*
- I-OLAP Characteristics:
  - Overlay multiple pieces of information
  - No change on the objects whose interactions are being examined
    - In the underlying snapshots, each node is a researcher
    - In the summarized view, each node is still a researcher

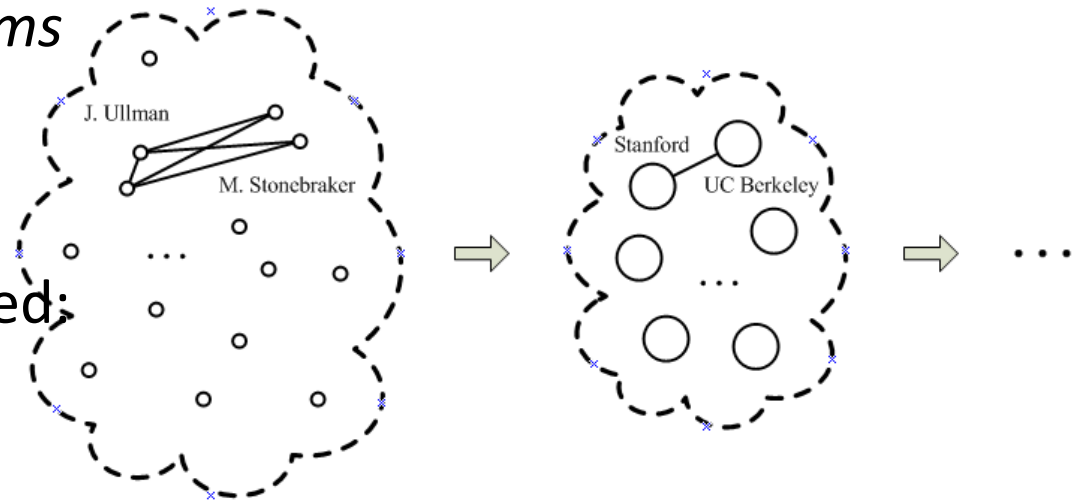


# Topological OLAP

- Dimensions come from the node/edge attributes inside individual networks, so-called *Topo-Dims*

- T-OLAP Characteristics

- Zoom in/Zoom out
- Network topology changed; “generalized” nodes and “generalized” edges
  - In the underlying network, each node is a researcher
  - In the summarized view, each node becomes an institute that comprises multiple researchers



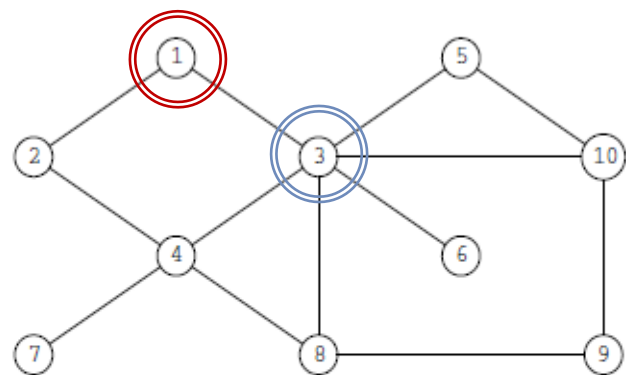
# InfoNet OLAP: Operations & Framework

	InfoNet I-OLAP	InfoNet T-OLAP
<b>Roll-up</b>	Overlay multiple snapshots to form a higher-level summary via I-aggregated network	Shrink the topology & obtain a T-aggregated network that represents a compressed view, with topological elements (i.e., nodes and/or edges) merged and replaced by corresp. higher-level ones
<b>Drill-down</b>	Return to the set of lower-level snapshots from the higher-level overlaid (aggregated) network	A reverse operation of roll-up
<b>Slice/dice</b>	Select a subset of qualifying snapshots based on Info-Dims	Select a subnetwork based on Topo-Dims

- Measure is an aggregated graph & other measures like node count, average degree, etc. can be treated as derived
- Graph plays a dual role: (1) data source, and (2) aggregate measure
- Measures could be complex, e.g., maximum flow, shortest path, centrality
- It is possible to combine I-OLAP and T-OLAP into a hybrid case

# Graph Cube: Online analytical processing in multidimensional information networks (Zhao, SIGMOD'11)

## A Multidimensional Information Network



(a) Graph

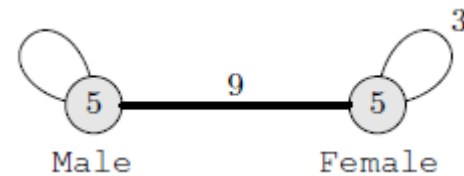
ID	Gender	Location	Profession	Income
1	Male	CA	Teacher	\$70,000
2	Female	WA	Teacher	\$65,000
3	Female	CA	Engineer	\$80,000
4	Female	NY	Teacher	\$90,000
5	Male	IL	Lawyer	\$80,000
6	Female	WA	Teacher	\$90,000
7	Male	NY	Lawyer	\$100,000
8	Male	IL	Engineer	\$75,000
9	Female	CA	Lawyer	\$120,000
10	Male	IL	Engineer	\$95,000

(b) Vertex Attribute Table

**Figure:** A Multidimensional Network Comprising a Graph Structure and a Multidimensional Vertex Attribute Table

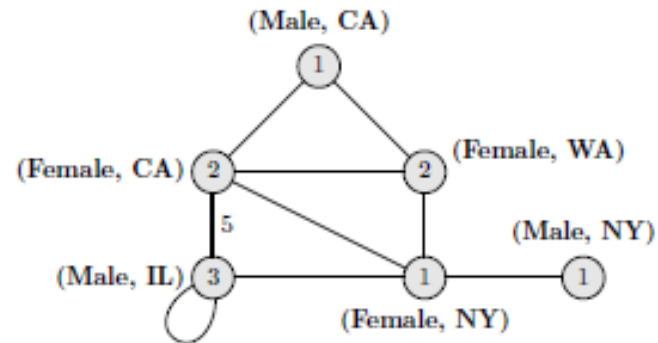
# Conventional Group-by v.s. Network Summarization

Gender	COUNT(*)
Male	5
Female	5



Group by “Gender”

Gender	Location	COUNT(*)
Male	CA	1
Female	CA	2
Female	WA	2
Male	IL	3
Male	NY	1
Female	NY	1



Group by “Gender” and “Location”

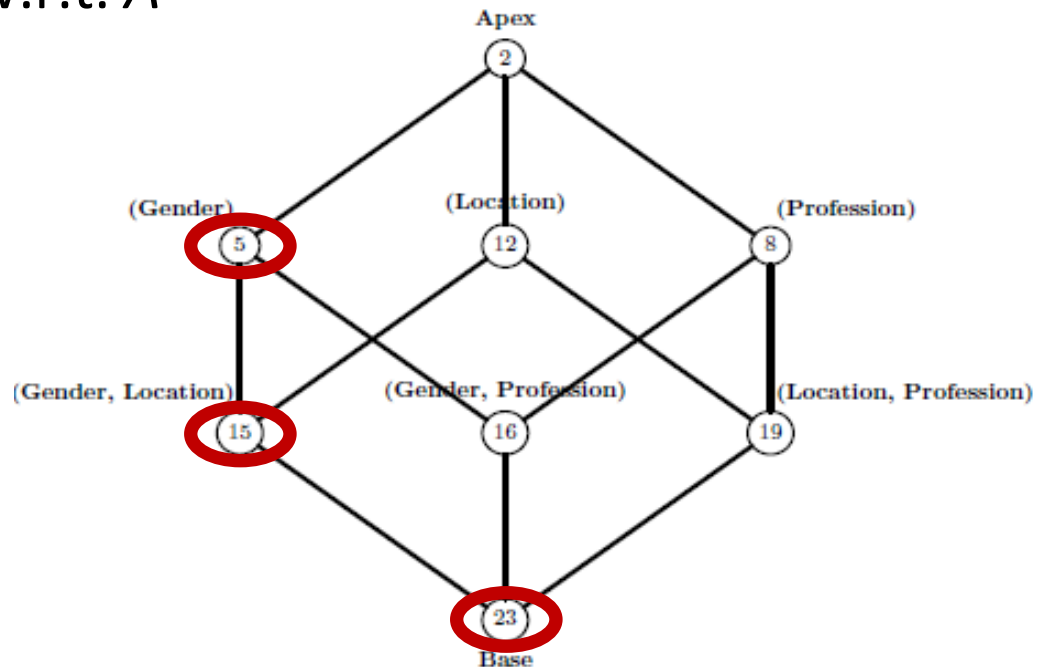
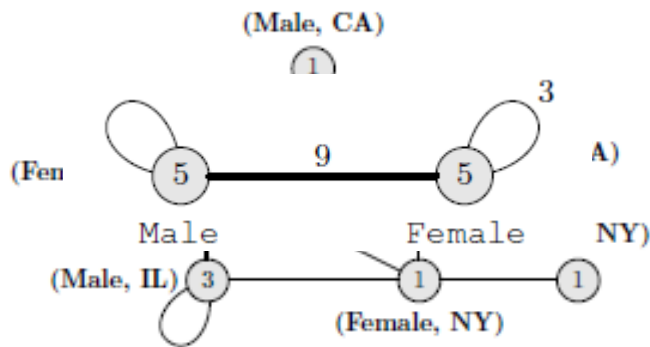
# The Graph Cube Model

---

- Multidimensional network  $N = (V, E; A)$ 
  - $A = \{A_1, A_2, \dots, A_n\}$ , the **dimension** of the network  $N$ , is a set of  $n$  vertex-specific attributes
    - Some (or all) dimension  $A_i$  could be  $*(ALL)$ , representing a super-aggregation along  $A_i$
    - there exist  **$2^n$**  multidimensional spaces (aggregations)
  - The **measure** within each possible space is no longer a simple numeric value, but an **aggregate network!**

# The Graph Cube Model

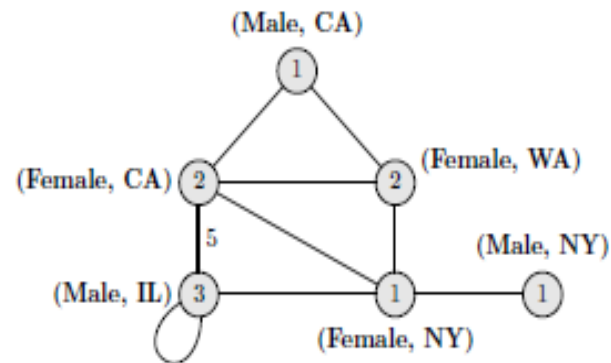
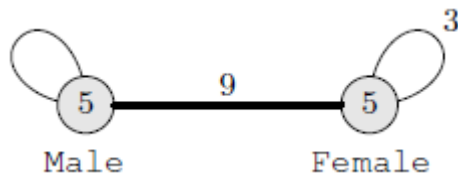
- Graph Cube
  - Restructure the network in **all possible multidimensional spaces (cuboids)** defined on  $A$
  - For each multidimensional space  $A'$ , the measure is an **aggregate network** w.r.t.  $A'$





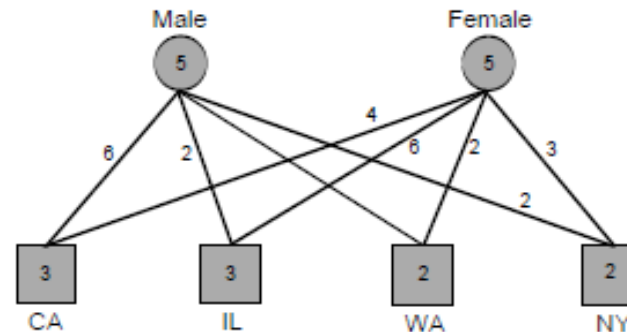
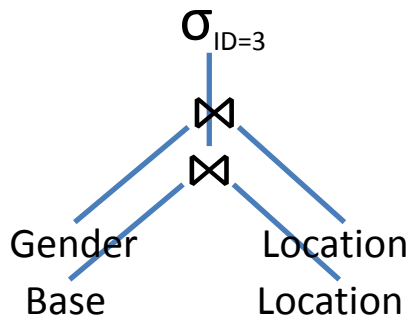
# OLAP on Graph Cube

- Cuboid query
  - Return as output the aggregate network corresponding to a specific multidimensional space (**cuboid**)
    - *What is the aggregate network between various **genders**?*
    - *What is the aggregate network between various **gender and location** combinations?*




# OLAP on Graph Cube

- Cuboid query
  - Within a single multidimensional space
- Crossboid query ( $\bowtie$ )
  - Crosses **multiple** multidimensional spaces of the network
    - *What is the network structure between user 3 and various locations?*
    - *What is the network structure between users grouped by gender v.s. users grouped by location?*



# Outline

---

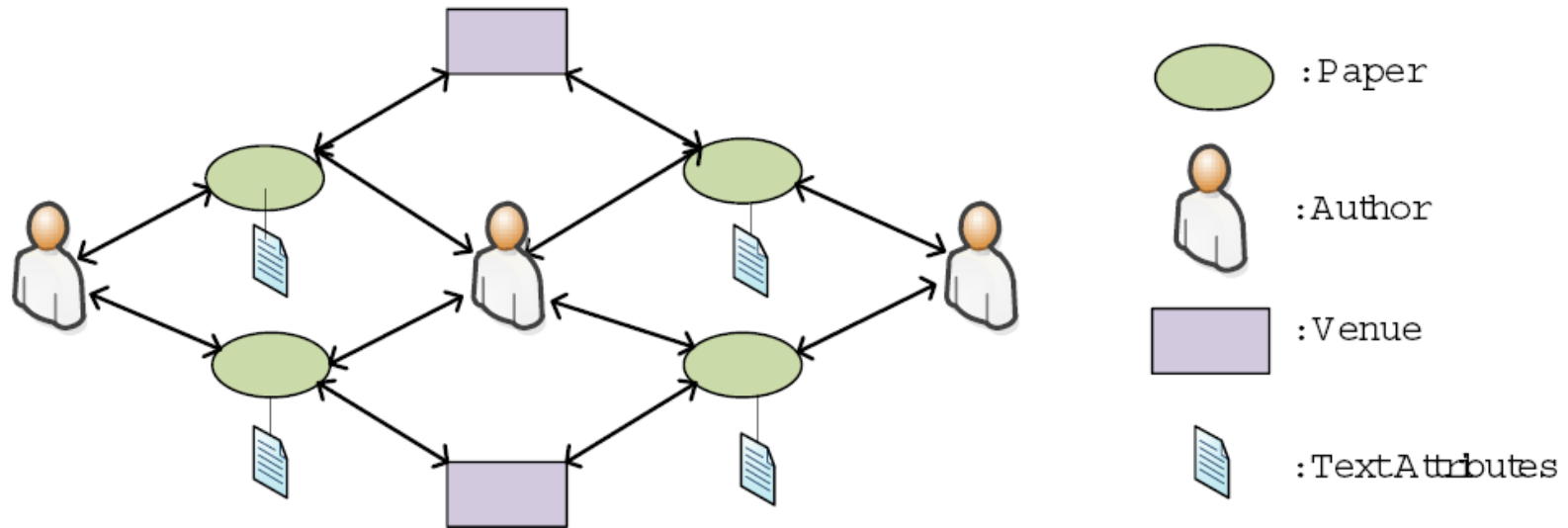
- **Motivation:** Why Mining Information Networks?
  - **Part I:** Clustering, Ranking and Classification
    - Clustering and Ranking in Information Networks
    - Classification of Information Networks
  - **Part II:** Meta-Path Based Exploration of Information Networks
    - Similarity Search in Information Networks
    - Relationship Prediction in Information Networks
  - **Part III:** Advanced Topics on Information Network Analysis
    - Role Discovery and OLAP in Information Networks
    - Relation Strength Learning in Information Networks 
    - Mining Evolution and Dynamics of Information Networks
  - **Conclusions**
-

# Relation Strength-Aware Clustering of Heterogeneous InfoNet with Incomplete Attributes [Sun, VLDB'12]

---

- **Content-Rich** Heterogeneous information networks become increasingly popular
  - Heterogeneous links + (incomplete) attributes
  - Examples
    - Social media
    - E-Commerce
    - Cyber-physical system
- **Soft clustering** objects using both link information and attribute information
  - E-Commerce: customers, products, comments, ...
  - Social websites: people, groups, books, posts, ...
- Understanding the **strengths for different relations** in determining object's cluster

# Example 1: Bibliographic Information Network



**Link type:**

- Paper-Author, Paper-Venue, (Paper->Paper)

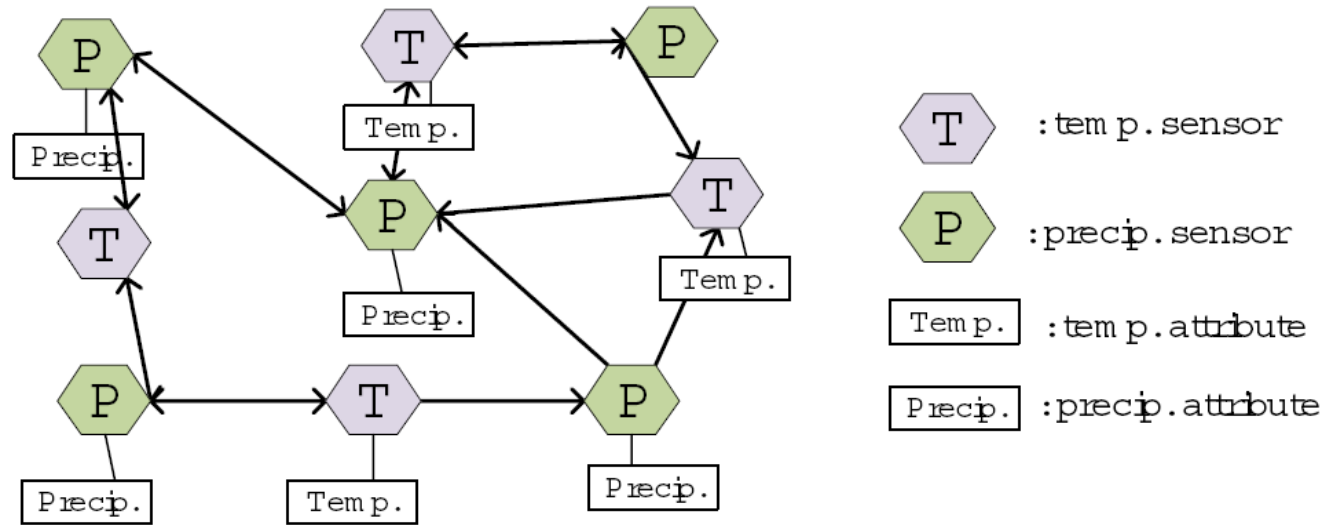
**Attribute type:**

- Text attribute for Paper type

**Goal:**

- Clustering authors, venues, papers into different research areas

# Example 2: Weather Sensor Information Network



**Link type:**

- T->P, T->T, P->P, P->T (According to KNN relationships)

**Attribute type:**

- Temperature attribute for T-typed sensors, Precipitation attribute for P-typed sensors

**Goal:**

- Clustering both types of sensors into different regional weather patterns

# Challenges

---

- Attributes are **incomplete** for objects
  - Not every type of objects contained the user specified attributes
    - E.g., Temperature typed sensors are only associated with temperature attributes
  - Missing value
    - E.g., some sensor may contain no observations due to malfunctioning
- Links are **heterogeneous**
  - Different types of links carry different importance in enhancing the quality of attribute-based clustering results
    - E.g., which type of links are more trustable to determine a person's political interest: friendship or person-like-book relationship?

# Solution Overview

- Modeling **attribute generation** and **structural consistency** in a unified framework

$$p(\{\{v[X]\}_{v \in V_X}\}_{X \in \mathcal{X}}, \Theta | G, \gamma, \beta) = \prod_{X \in \mathcal{X}} p(\{v[X]\}_{v \in V_X} | \Theta, \beta) p(\Theta | G, \gamma)$$

- Attribute generation as a mixture model

- $$p(\{v[X]\}_{v \in V_X} | \Theta, \beta) = \prod_{v \in V_X} \prod_{x \in v[X]} \sum_{k=1}^K \theta_{v,k} p(x | \beta_k)$$

- $v[X]$ : *observed values for Attribute X on Object v*
- $\Theta$ : *soft clustering membership matrix*
- $\beta$ : *parameters associated with each mixture model component*

- Structural consistency as a log-linear model

- $$p(\Theta | G, \gamma) = \frac{1}{Z(\gamma)} \exp\left\{ \sum_{e=\langle v_i, v_j \rangle \in E} f(\theta_i, \theta_j, e, \gamma) \right\}$$

- $\gamma$ : *relation strength vector*

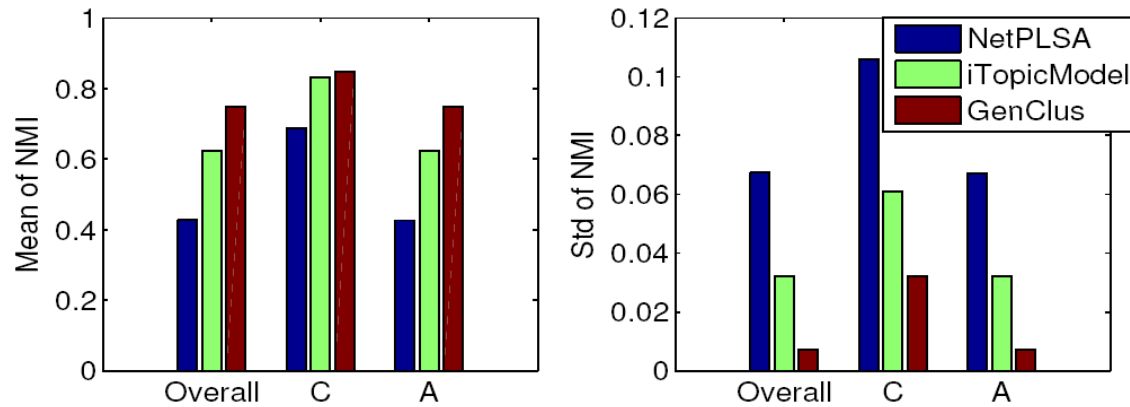


# The Objective Function and the Algorithm Overview

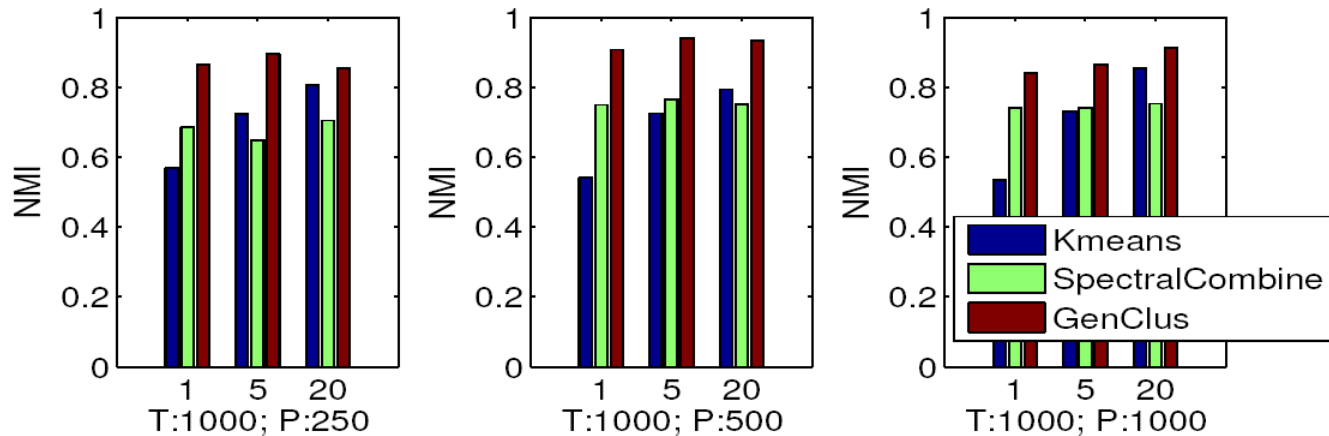
$$g(\Theta, \beta, \gamma) = \underbrace{\log \sum_{X \in \mathcal{X}} p(\{v[X]\}_{v \in V_X} | \Theta, \beta)}_{\text{Attribute Generation}} + \underbrace{\log p(\Theta | G, \gamma)}_{\text{Structural Consistency}} - \underbrace{\frac{||\gamma||^2}{2\sigma^2}}_{\text{Regularization Term}}$$

- The clustering algorithm
  - Iterative algorithm
    - Step 1: Fix the relation strength and optimize the clustering result
      - Cluster optimization
    - Step 2: Fix the clustering result and optimize the relation strength
      - Relation strength learning

# Higher Accuracy and More Stable Clustering Results

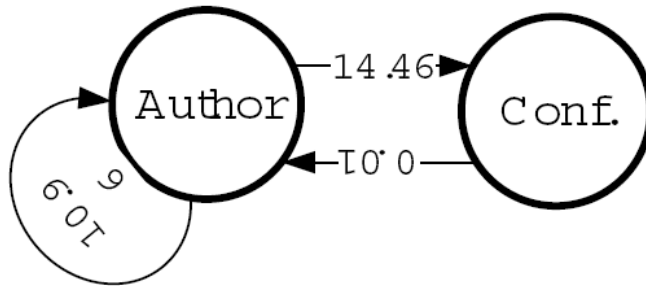


## Clustering Accuracy Comparisons for AC



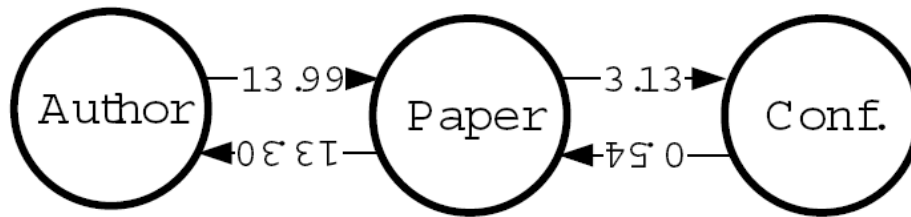
## Clustering Accuracy Comparisons for Weather Sensor Network

# Intuitive relation strength weights



(a) AC Network

An author's research area is more determined by their attended venues than their co-authors (14.46 vs. 10.69)




(b) ACP Network

A paper's research area is more determined by its authors than its venue (13.30 vs. 3.13)

# Outline

---

- **Motivation:** Why Mining Information Networks?
  - **Part I:** Clustering, Ranking and Classification
    - Clustering and Ranking in Information Networks
    - Classification of Information Networks
  - **Part II:** Meta-Path Based Exploration of Information Networks
    - Similarity Search in Information Networks
    - Relationship Prediction in Information Networks
  - **Part III:** Advanced Topics on Information Network Analysis
    - Role Discovery and OLAP in Information Networks
    - Relation Strength Learning in Information Networks
    - Mining Evolution and Dynamics of Information Networks 
  - **Conclusions**
-

# Mining Evolution and Dynamics of InfoNet

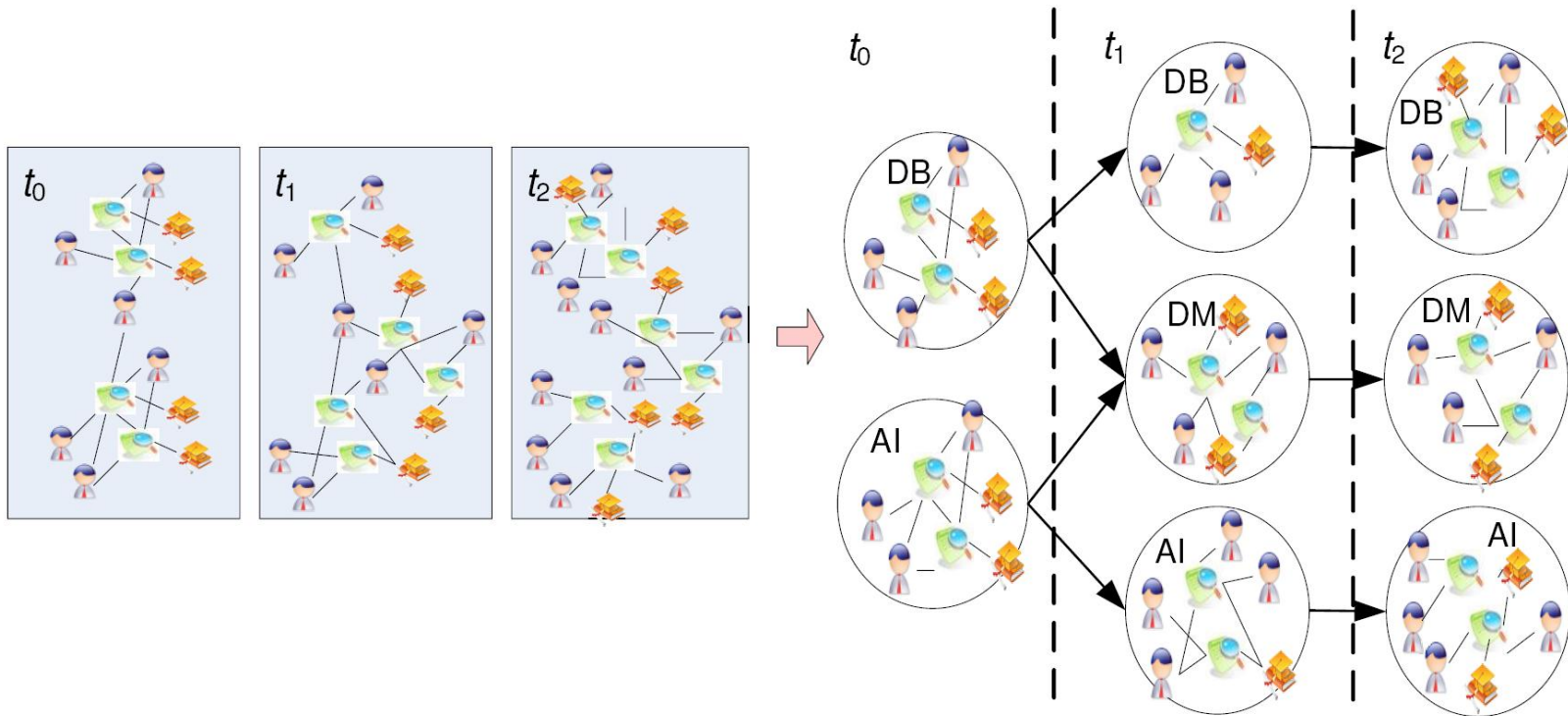
## [Sun, MLG'10]

---

- Many networks are with time information
  - E.g., according to paper publication year, DBLP networks can form network sequences
- Motivation: Model evolution of communities in heterogeneous network
  - Automatically detect the best number of communities in each timestamp
  - Model the smoothness between communities of adjacent timestamps
  - Model the evolution structure explicitly
    - Birth, death, split

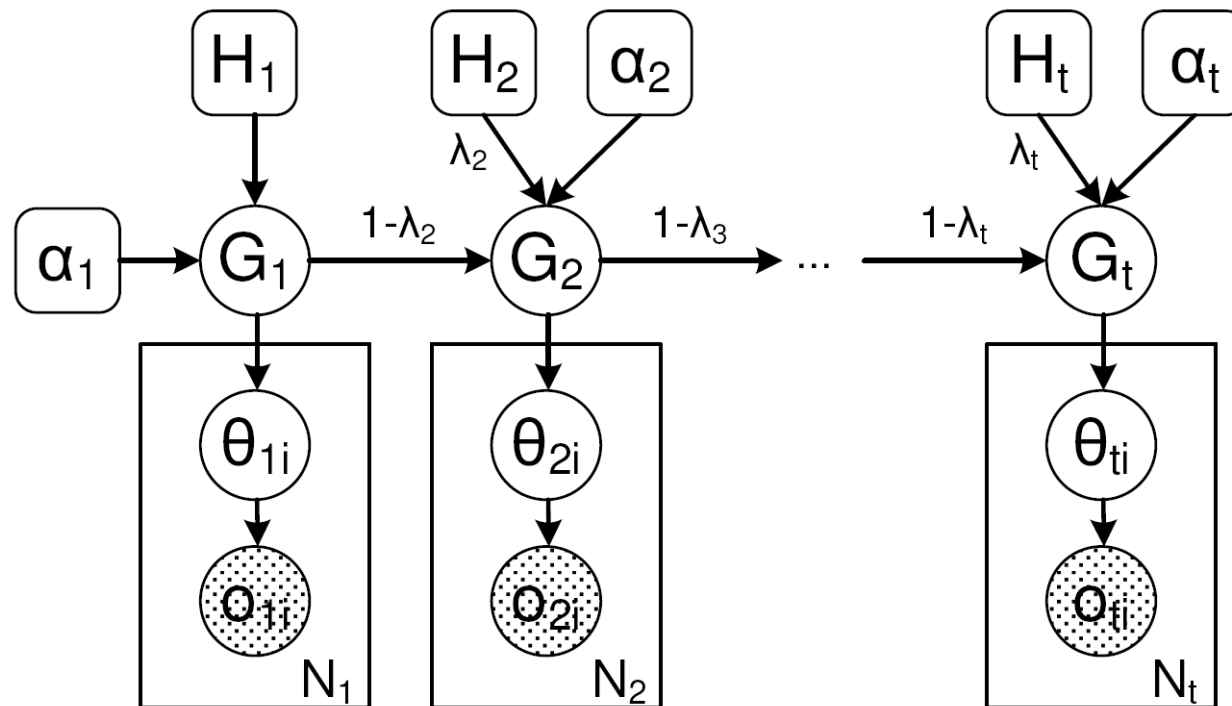
# Evolution: Idea Illustration

- From network sequences to evolutionary communities



# Graphical Model: A Generative Model

- **Dirichlet Process Mixture Model-based generative model**
  - At each timestamp, a community is dependent on historical communities and background community distribution



# Generative Model & Model Inference

- To generate a new paper  $o_i$ 
  - Decide whether to join an existing community or a new one
    - Join an existing community  $k$  with prob.  $n_k/(i-1+\alpha)$
    - Join a new community  $k$  with prob.  $\alpha/(i-1+\alpha)$ : Decide its prior, either from a background distribution ( $\lambda$ ) or historical communities  $((1-\lambda)\pi_k)$ , with different probabilities, draw the attribute distribution from the prior
  - Generate  $o_i$  according to the attribute distribution

$$\begin{aligned} p(o_{i,t} | z_{i,t} = k, \Theta_t) &= p(o_{i,t} | \theta_{k,t}) \\ &= p(\mathbf{a}_{i,t} | \theta_{k,t}^A) p(\mathbf{c}_{i,t} | \theta_{k,t}^C) p(\mathbf{d}_{i,t} | \theta_{k,t}^D) \\ &= \prod_{j=1}^{|A|} \theta_{k,t}^A(j)^{a_{ij,t}} \prod_{j=1}^{|C|} \theta_{k,t}^C(j)^{c_{ij,t}} \prod_{j=1}^{|D|} \theta_{k,t}^D(j)^{d_{ij,t}} \end{aligned}$$

- Greedy inference for each timestamp: Collapse Gibbs sampling, which is trying to sample cluster label for each target object (e.g., paper)



# Accuracy Study

- The more types of objects used, the better accuracy
- Historical prior results in better accuracy

Year	Training Type	Testing Type	Test Size 10% (cluster number $K$ )	Test Size 20% (cluster number $K$ )
1992	Term	Term	1.600 (4)	1.390 (4)
1992	Term+Author	Term+Author	2.205 (8)	1.697 (6)
1992	Term+Author+Conf.	Term+Author	2.434 (8)	2.095 (8)
1992 1991	Term+Author+Conf.	Term+Author	<b>2.8365</b> (8)	<b>2.671</b> (8)

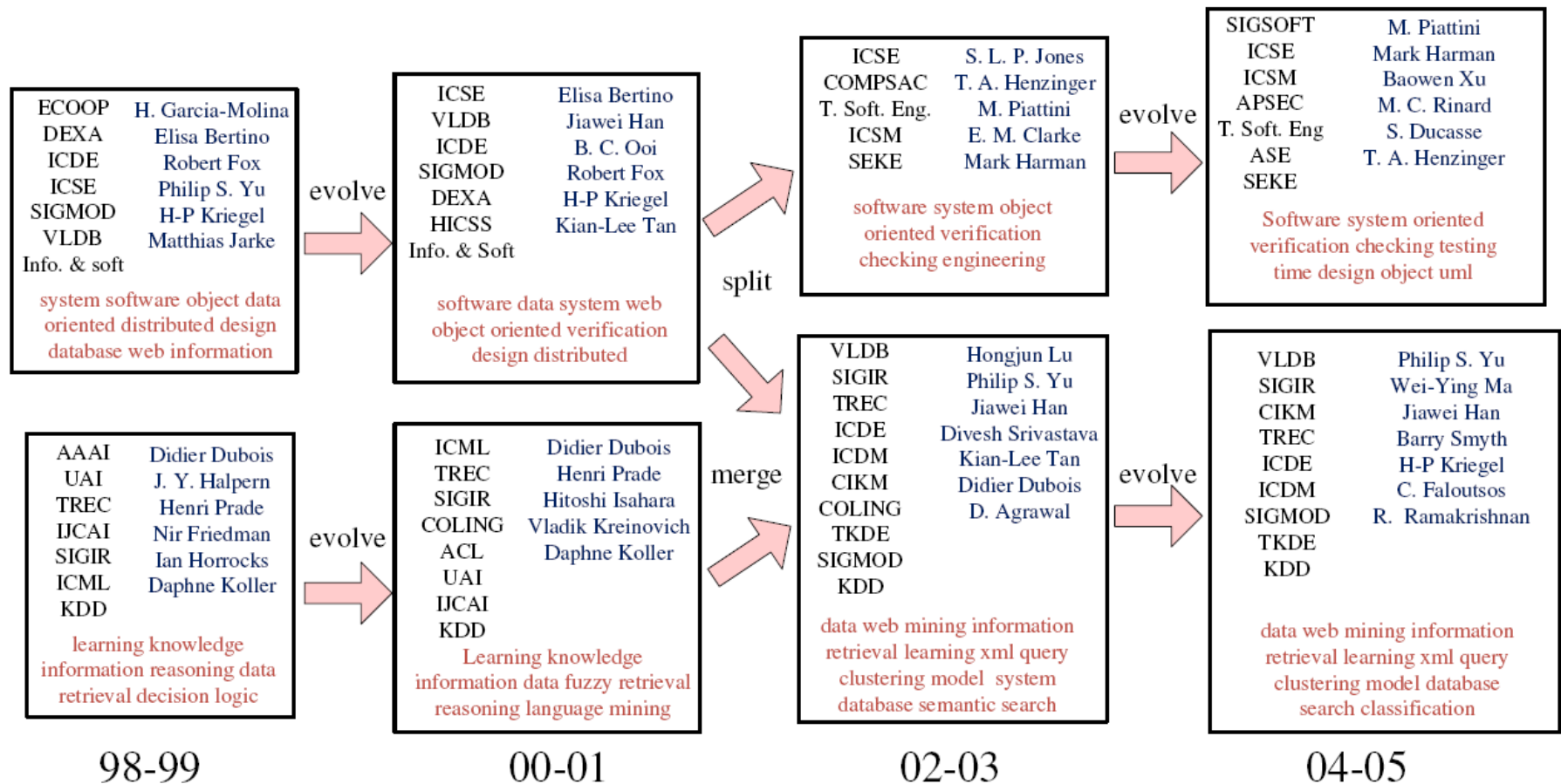
**Table 1: Conference Compactness of Different Models on Test Dataset**

Year	Training Type	Testing Type	Test Size 10%	Test Size 20%
1992	Term+Author+Conf.	Term+Author+Conf.	$3.493 \times 10^{18}$	$4.673 \times 10^{18}$
1992 1991	Term+Author+Conf.	Term+Author+Conf.	<b><math>6.384 \times 10^{17}</math></b>	<b><math>7.106 \times 10^{17}</math></b>

**Table 2: Perplexity Comparison between Models with/without Historical Prior**

# Case Study on DBLP

- Tracking database and information system community evolution



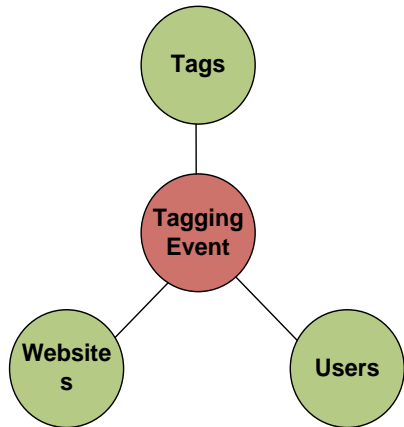
# Case Study on Delicious.com

Jan. 1 - Jan. 7

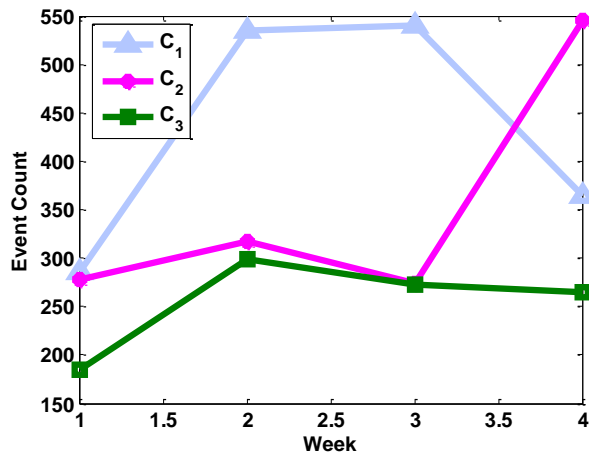
Jan. 8 - Jan. 14

Jan. 15 - Jan. 21

Jan. 22 - Jan. 28



Delicious Schema



C<sub>1</sub>:

Security
Terrorism
Politics
Travel
Usa
Airport
Israel
Obama
CIA
Afghanistan



Google
China
Security
Internet
Privacy
Politics
Censorship
Facebook
Business
Terrorism



Security
Google
China
Internet
Microsoft
Privacy
Censorship
Politics
Browser
USA



Google
Security
China
Internet
Privacy
Digg
Politics
Datenschutz
Facebook
USA

C<sub>2</sub>:

Mac
Apple
Iphone
Windows
Tablet
Ipod
Tips
Macbook
Tutorial
Drm



Iphone
Apple
Twitter
Mac
Mobile
Apps
Ratio
Blog
Newspapers
Technology



Iphone
Apple
Mac
Mobile
Twitter
Software
Apps
Business
Osx
Radio



Ipad
Apple
Iphone
Technology
Tablet
Mac
Mobile
Newspapers
Kindle
Media

C<sub>3</sub>:

Health
Depression
Sleep
Teenagers
Dubai
Tallest
BBC
Building
Architecture
Mentalhealth



Weather
UK
Photography
Photo
Haiti
Photos
2010
BBC
Snow
Earthquake




Haiti
Photography
BBC
Earthquake
Photos
UK
2010
Disaster
Travel
Wildlife



Haiti
BBC
Photography
Animals
Earthquake
2010
Photos
Nature
Funny
Theonion

# Outline

---

- **Motivation:** Why Mining Information Networks?
- **Part I:** Clustering, Ranking and Classification
  - Clustering and Ranking in Information Networks
  - Classification of Information Networks
- **Part II:** Meta-Path Based Exploration of Information Networks
  - Similarity Search in Information Networks
  - Relationship Prediction in Information Networks
- **Part III:** Advanced Topics on Information Network Analysis
  - Role Discovery and OLAP in Information Networks
  - Mining Evolution and Dynamics of Information Networks
- **Conclusions** 

# Conclusions

---

- Rich knowledge can be mined from information networks
- What is the magic?
  - ***Heterogeneous, semi-structured information networks!***
- Clustering, ranking and classification: Integrated clustering, ranking and classification: RankClus, NetClus, GNetMine, ...
- Meta-Path based similarity search and relationship prediction
- Role discovery, OLAP, relation strength learning, and evolutionary analysis
- Knowledge is power, but knowledge is hidden in massive links!
- ***Mining heterogeneous information networks:*** Much more to be explored!!

# Future Research

---

- Discovering **ontology** and structure in information networks
- Discovering and mining **hidden** information networks
- Mining information networks formed **by structured data linking with unstructured data** (text, multimedia and Web)
- Mining **cyber-physical** networks (networks formed by dynamic sensors, image/video cameras, with information networks)
- Enhancing the power of knowledge discovery by transforming massive **unstructured data**: Incremental information extraction, role discovery, ...  $\Rightarrow$  multi-dimensional structured info-net
- Mining **noisy, uncertain, un-trustable** massive datasets by information network analysis approach
- Turning **Wikipedia and/or Web** into structured or semi-structured databases by heterogeneous information network analysis

# References: Books on Network Analysis

---

- A.-L. Barabasi. Linked: How Everything Is Connected to Everything Else and What It Means. Plume, 2003.
- M. Buchanan. Nexus: Small Worlds and the Groundbreaking Theory of Networks. W. W. Norton & Company, 2003.
- P. J. Carrington, J. Scott, and S. Wasserman. Models and Methods in Social Network Analysis. Cambridge University Press, 2005.
- S. Chakrabarti. Mining the Web: Discovering Knowledge from Hypertext Data. Morgan Kaufmann, 2003.
- D. J. Cook and L. B. Holder. Mining Graph Data. John Wiley & Sons, 2007.
- J. Davies, D. Fensel, and F. van Harmelen. Towards the Semantic Web: Ontology-Driven Knowledge Management. John Wiley & Sons, 2003.
- A. Degenne and M. Forse. Introducing Social Networks. Sage Publications, 1999.
- M. O. Jackson. Social and Economic Networks. Princeton University Press, 2010.
- D. Easley and J. Kleinberg. Networks, Crowds, and Markets. Cambridge University Press, 2010.
- D. Fensel, W. Wahlster, H. Lieberman, and J. Hendler. Spinning the Semantic Web: Bringing the World Wide Web to Its Full Potential. MIT Press, 2002.
- L. Getoor and B. Taskar (eds.). Introduction to statistical learning. In MIT Press, 2007.
- B. Liu. Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data. Springer, 2006.
- M. E. J. Newman. Networks: An Introduction. Oxford University Press, 2010
- J. P. Scott. Social Network Analysis: A Handbook. Sage Publications, 2005.
- J. Watts. Six Degrees: The Science of a Connected Age. W. W. Norton & Company, 2003.
- D. J. Watts. Small Worlds: The Dynamics of Networks between Order and Randomness. Princeton University Press, 2003.
- S. Wasserman and K. Faust. Social Network Analysis: Methods and Applications. Cambridge University Press, 1994.

# References: Some Overview Papers

---

- T. Berners-Lee, J. Hendler, and O. Lassila. The semantic web. Scientific American, May 2001.
- C. Cooper and A Frieze. A general model of web graphs. Algorithms, 22, 2003.
- S. Chakrabarti and C. Faloutsos. Graph mining: Laws, generators, and algorithms. ACM Comput. Surv., 38, 2006.
- T. Dietterich, P. Domingos, L. Getoor, S. Muggleton, and P. Tadepalli. Structured machine learning: The next ten years. Machine Learning, 73, 2008
- S. Dumais and H. Chen. Hierarchical classification of web content. SIGIR'00.
- S. Dzeroski. Multirelational data mining: An introduction. ACM SIGKDD Explorations, July 2003.
- L. Getoor. Link mining: a new data mining challenge. SIGKDD Explorations, 5:84{89, 2003.
- L. Getoor, N. Friedman, D. Koller, and B. Taskar. Learning probabilistic models of relational structure. ICML'01
- D. Jensen and J. Neville. Data mining in networks. In Papers of the Symp. Dynamic Social Network Modeling and Analysis, National Academy Press, 2002.
- T. Washio and H. Motoda. State of the art of graph-based data mining. SIGKDD Explorations, 5, 2003.



# References: Some Influential Papers

---

- A. Z. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. L. Wiener. Graph structure in the web. *Computer Networks*, 33, 2000.
- S. Brin and L. Page. The anatomy of a large-scale hyper-textual web search engine. *WWW'98*.
- S. Chakrabarti, B. E. Dom, S. R. Kumar, P. Raghavan, S. Rajagopalan, A. Tomkins, D. Gibson, and J. M. Kleinberg. Mining the web's link structure. *COMPUTER*, 32, 1999.
- M. Faloutsos, P. Faloutsos, and C. Faloutsos. On power-law relationships of the internet topology. *ACM SIGCOMM'99*
- M. Girvan and M. E. J. Newman. Community structure in social and biological networks. In *Proc. Natl. Acad. Sci. USA* 99, 2002.
- B. A. Huberman and L. A. Adamic. Growth dynamics of world-wide web. *Nature*, 399:131, 1999.
- G. Jeh and J. Widom. SimRank: a measure of structural-context similarity. *KDD'02*
- D. Kempe, J. Kleinberg, and E. Tardos. Maximizing the spread of influence through a social network. *KDD'03*
- J. M. Kleinberg, R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. The web as a graph: Measurements, models, and methods. *COCOON'99*
- J. M. Kleinberg. Small world phenomena and the dynamics of information. *NIPS'01*
- R. Kumar, P. Raghavan, S. Rajagopalan, D. Sivakumar, A. Tomkins, and E. Upfal. Stochastic models for the web graph. *FOCS'00*
- M. E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45, 2003.

# References: Clustering and Ranking (1)

---

- E. Airoldi, D. Blei, S. Fienberg and E. Xing, “Mixed Membership Stochastic Blockmodels”, JMLR’08
- Liangliang Cao, Andrey Del Pozo, Xin Jin, Jiebo Luo, Jiawei Han, and Thomas S. Huang, “[RankCompete: Simultaneous Ranking and Clustering of Web Photos](#)”, WWW’10
- G. Jeh and J. Widom, “SimRank: a measure of structural-context similarity”, KDD’02
- Jing Gao, Feng Liang, Wei Fan, Chi Wang, Yizhou Sun, and Jiawei Han, “[Community Outliers and their Efficient Detection in Information Networks](#)”, KDD’10
- M. E. J. Newman and M. Girvan, “Finding and evaluating community structure in networks”, Physical Review E, 2004
- M. E. J. Newman and M. Girvan, “Fast algorithm for detecting community structure in networks”, Physical Review E, 2004
- J. Shi and J. Malik, “Normalized cuts and image Segmentation”, CVPR’97
- Yizhou Sun, Yintao Yu, and Jiawei Han, "Ranking-Based Clustering of Heterogeneous Information Networks with Star Network Schema", KDD’09
- Yizhou Sun, Jiawei Han, Peixiang Zhao, Zhijun Yin, Hong Cheng, and Tianyi Wu, "RankClus: Integrating Clustering with Ranking for Heterogeneous Information Network Analysis", EDBT’09

# References: Clustering and Ranking (2)

---

- Yizhou Sun, Jiawei Han, Jing Gao, and Yintao Yu, "iTopicModel: Information Network-Integrated Topic Modeling", ICDM'09
- Yizhou Sun, Charu C. Aggarwal, and Jiawei Han, "*Relation Strength-Aware Clustering of Heterogeneous Information Networks with Incomplete Attributes*", PVLDB 5(5), 2002
- A. Wu, M. Garland, and J. Han. Mining scale-free networks using geodesic clustering. KDD'04
- Z. Wu and R. Leahy, "An optimal graph theoretic approach to data clustering: Theory and its application to image segmentation", IEEE Trans. Pattern Anal. Mach. Intell., 1993.
- X. Xu, N. Yuruk, Z. Feng, and T. A. J. Schweiger. SCAN: A structural clustering algorithm for networks. KDD'07
- Xiaoxin Yin, Jiawei Han, Philip S. Yu. "[LinkClus: Efficient Clustering via Heterogeneous Semantic Links](#)", VLDB'06.
- Yintao Yu, Cindy X. Lin, Yizhou Sun, Chen Chen, Jiawei Han, Binbin Liao, Tianyi Wu, ChengXiang Zhai, Duo Zhang, and Bo Zhao, "iNextCube: Information Network-Enhanced Text Cube", VLDB'09 (demo)
- X. Yin, J. Han, and P. S. Yu. Cross-relational clustering with user's guidance. KDD'05

# References: Network Classification (1)

---

- A. Appice, M. Ceci, and D. Malerba. Mining model trees: A multi-relational approach. ILP'03
- Jing Gao, Feng Liang, Wei Fan, Yizhou Sun, and Jiawei Han, "Bipartite Graph-based Consensus Maximization among Supervised and Unsupervised Models ", NIPS'09
- L. Getoor, N. Friedman, D. Koller and B. Taskar, “Learning Probabilistic Models of Link Structure”, JMLR’02.
- L. Getoor, E. Segal, B. Taskar and D. Koller, “Probabilistic Models of Text and Link Structure for Hypertext Classification”, IJCAI WS ‘Text Learning: Beyond Classification’, 2001.
- L. Getoor, N. Friedman, D. Koller, and A. Pfeffer, “Learning Probabilistic Relational Models”, chapter in Relation Data Mining, eds. S. Dzeroski and N. Lavrac, 2001.
- M. Ji, Y. Sun, M. Danilevsky, J. Han, and J. Gao, “Graph-based classification on heterogeneous information networks”, ECMLPKDD’10.
- M. Ji, J. Jan, and M. Danilevsky, “Ranking-based Classification of Heterogeneous Information Networks”, KDD’11.
- Q. Lu and L. Getoor, “Link-based classification”, ICML'03
- D. Liben-Nowell and J. Kleinberg, “The link prediction problem for social networks”, CIKM'03

# References: Network Classification (2)

---

- J. Neville, B. Gallaher, and T. Eliassi-Rad. Evaluating statistical tests for within-network classifiers of relational data. ICDM'09.
- J. Neville, D. Jensen, L. Friedland, and M. Hay. Learning relational probability trees. KDD'03
- Jennifer Neville, David Jensen, “Relational Dependency Networks”, JMLR'07
- M. Szummer and T. Jaakkola, “Partially labeled classification with markov random walks”, In NIPS, volume 14, 2001.
- M. J. Rattigan, M. Maier, and D. Jensen. Graph clustering with network structure indices. ICML'07
- P. Sen, G. M. Namata, M. Galileo, M. Bilgic, L. Getoor, B. Gallagher, and T. Eliassi-Rad. Collective classification in network data. AI Magazine, 29, 2008.
- B. Taskar, E. Segal, and D. Koller. Probabilistic classification and clustering in relational data. IJCAI'01
- B. Taskar, P. Abbeel, M.F. Wong, and D. Koller, “[Relational Markov Networks](#)”, chapter in L. Getoor and B. Taskar, editors, [Introduction to Statistical Relational Learning](#), 2007
- X. Yin, J. Han, J. Yang, and P. S. Yu, “[CrossMine: Efficient Classification across Multiple Database Relations](#)”, ICDE'04.
- D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Scholkopf, “Learning with local and global consistency”, In *NIPS 16*, Vancouver, Canada, 2004.
- X. Zhu and Z. Ghahramani, “Learning from labeled and unlabeled data with label propagation”, Technical Report, 2002.

# References: Social Network Analysis

---

- B. Aleman-Meza, M. Nagarajan, C. Ramakrishnan, L. Ding, P. Kolari, A. P. Sheth, I. B. Arpinar, A. Joshi, and T. Finin. Semantic analytics on social networks: experiences in addressing the problem of conflict of interest detection. WWW'06
- R. Agrawal, S. Rajagopalan, R. Srikant, and Y. Xu. Mining newsgroups using networks arising from social behavior. WWW'03
- P. Boldi and S. Vigna. The WebGraph framework I: Compression techniques. WWW'04
- D. Cai, Z. Shao, X. He, X. Yan, and J. Han. Community mining from multi-relational networks. PKDD'05
- P. Domingos. Mining social networks for viral marketing. IEEE Intelligent Systems, 20, 2005.
- P. Domingos and M. Richardson. Mining the network value of customers. KDD'01
- P. DeRose, W. Shen, F. Chen, A. Doan, and R. Ramakrishnan. Building structured web community portals: A top-down, compositional, and incremental approach. VLDB'07
- G. Flake, S. Lawrence, C. L. Giles, and F. Coetzee. Self-organization and identification of web communities. IEEE Computer, 35, 2002.
- J. Kubica, A. Moore, and J. Schneider. Tractable group detection on large link data sets. ICDM'03

# References: Data Quality & Search in Networks

---

- I. Bhattacharya and L. Getoor, “Iterative record linkage for cleaning and integration”, Proc. SIGMOD 2004 Workshop on Research Issues on Data Mining and Knowledge Discovery (DMKD'04)
- Xin Luna Dong, Laure Berti-Equille, and Divesh Srivastava, “Integrating conflicting data: The role of source dependence”, PVLDB, 2(1):550–561, 2009.
- Xin Luna Dong, Laure Berti-Equille, and Divesh Srivastava, “Truth discovery and copying detection in a dynamic world”, PVLDB, 2(1):562–573, 2009.
- H. Han, L. Giles, H. Zha, C. Li, and K. Tsoutsoulis, “Two supervised learning approaches for name disambiguation in author citations”, ICDL'04.
- Y. Sun, J. Han, T. Wu, X. Yan, and Philip S. Yu, “PathSim: Meta Path-Based Top-K Similarity Search in Heterogeneous Information Networks”, VLDB'11.
- X. Yin, J. Han, and P. S. Yu, “Object Distinction: Distinguishing Objects with Identical Names by Link Analysis”, ICDE'07.
- X. Yin, J. Han, and P. S. Yu, “Truth Discovery with Multiple Conflicting Information Providers on the Web”, IEEE TKDE, 20(6):796-808, 2008
- P. Zhao and J. Han, “On Graph Query Optimization in Large Networks”, VLDB'10.

# References: Link and Relationship Prediction

---

- V. Leroy, B. B. Cambazoglu, and F. Bonchi, “Cold start link prediction”, *KDD '10*.
- D. Liben-Nowell and J. Kleinberg, “The link prediction problem for social networks”, *CIKM '03*,
- R. N. Lichtenwalter, J. T. Lussier, and N. V. Chawla, “New perspectives and methods in link prediction”, *KDD'10*.
- Yizhou Sun, Rick Barber, Manish Gupta, Charu C. Aggarwal and Jiawei Han, "Co-Author Relationship Prediction in Heterogeneous Bibliographic Networks", *ASONAM'11*.
- Yizhou Sun, Jiawei Han, Charu C. Aggarwal, and Nitesh V. Chawla, "When Will It Happen? --- Relationship Prediction in Heterogeneous Information Networks", *WSDM'12*.
- B. Taskar, M. fai Wong, P. Abbeel, and D. Koller, “Link prediction in relational data”, *NIPS '03*.
- Xiao Yu, Quanquan Gu, Mianwei Zhou, and Jiawei Han, "Citation Prediction in Heterogeneous Bibliographic Networks", *SDM'12*.



# References: Role Discovery, Summarization and OLAP

---

- D. Archambault, T. Munzner, and D. Auber. Topolayout: Multilevel graph layout by topological features. IEEE Trans. Vis. Comput. Graph, 2007.
- Chen Chen, Xifeng Yan, Feida Zhu, Jiawei Han, and Philip S. Yu, "Graph OLAP: Towards Online Analytical Processing on Graphs", ICDM 2008
- Chen Chen, Xifeng Yan, Feida Zhu, Jiawei Han, and Philip S. Yu, "Graph OLAP: A Multi-Dimensional Framework for Graph Data Analysis", KAIS 2009.
- Xin Jin, Jiebo Luo, Jie Yu, Gang Wang, Dhiraj Joshi, and Jiawei Han, "[\*iRIN: Image Retrieval in Image Rich Information Networks\*](#)", WWW'10 (demo paper)
- Lu Liu, Feida Zhu, Chen Chen, Xifeng Yan, Jiawei Han, Philip Yu, and Shiqiang Yang, "[\*Mining Diversity on Networks\*](#)", DASFAA'10
- Y. Tian, R. A. Hankins, and J. M. Patel. Efficient aggregation for graph summarization. SIGMOD'08
- 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
- Zhijun Yin, Manish Gupta, Tim Weninger and Jiawei Han, "[\*LINKREC: A Unified Framework for Link Recommendation with User Attributes and Graph Structure\*](#) ", WWW'10
- Peixiang Zhao, Xiaolei Li, Dong Xin, Jiawei Han. Graph Cube: On Warehousing and OLAP Multidimensional Networks, SIGMOD'11

# References: Network Evolution

---

- L. Backstrom, D. Huttenlocher, J. Kleinberg, and X. Lan. Group formation in large social networks: Membership, growth, and evolution. KDD'06
- M.-S. Kim and J. Han. A particle-and-density based evolutionary clustering method for dynamic networks. VLDB'09
- J. Leskovec, J. Kleinberg, and C. Faloutsos. Graphs over time: Densification laws, shrinking diameters and possible explanations. KDD'05
- Yizhou Sun, Jie Tang, Jiawei Han, Manish Gupta, Bo Zhao, “Community Evolution Detection in Dynamic Heterogeneous Information Networks”, KDD-MLG'10