



Data-Driven Behavioral Analytics: Observations, Representations and Models

Meng Jiang (UIUC)

Peng Cui (Tsinghua)

Jiawei Han (UIUC)

<http://www.meng-jiang.com/tutorial-cikm16.html>



What is Behavior?

- **Definition.** Interactions made by individuals in conjunction with themselves or their environment. (Wikipedia)





Behavioral Analysis

- ❑ *Significance.* What can we discover from behavioral data?
 - ❑ Ex. Given every phone call/message between military leaders, scientists, businesspersons, find...

Observations

Who, what, where, when, why, how...
(scientific view)

Representations

Graph, network, matrix, tensor...
(mathematical view)

Models

Prediction, recommendation, anomaly detection...
(application view)

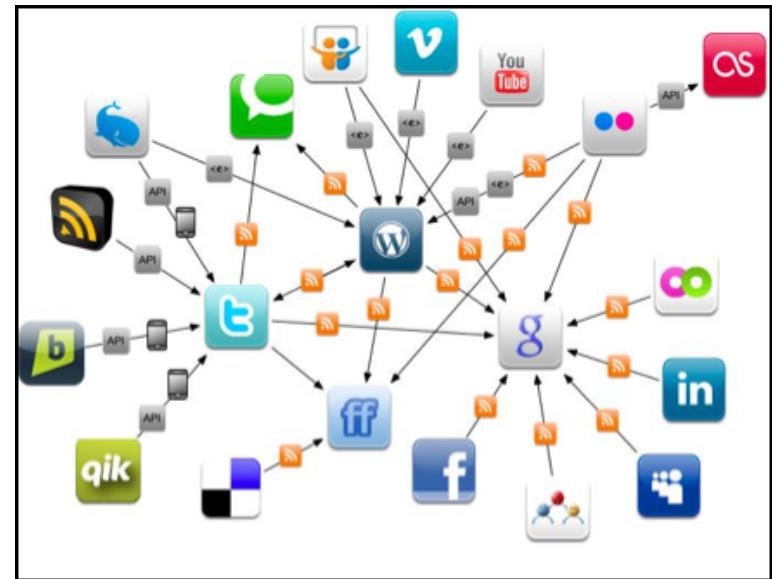
Why Behavioral Analysis Today?

- *Today.* The human behaviors are broadly recorded in an unprecedented level. Insights of sciences and society?

Physical World



Online Applications





Basic Research Areas

- Six Disruptive Basic Research Areas
 - Engineered Materials (metamaterials and plasmonics)
 - Quantum Information and Control
 - Cognitive Neuroscience
 - Nanoscience and Nanoengineering
 - Synthetic Biology
 - Computational Modeling of Human and Social Behavior



VI. Computational Models of Human Behavior



A fundamental understanding and predictive capability of human behavior dynamics from individuals to societies.

- **Enabled capabilities**

- Predictive models supporting strategic, operational, and tactical decision making and planning
- Real time cultural situational awareness
- Immersive training and mission rehearsal
- Cross cultural coalition building

- **Key research challenges:**

- Conflicting theories
- Data management and fusion
- Mathematical complexity
- Validation of models

Costly Punishment Across Human Societies

Joseph Henrich,^{1,*} Richard McElreath,² Abigail S. Alexander Bolhuis,² Juan Camilo Cardenas,³ Natalie Henrich,² Carolyn Lescroart,² Frank M.

Recent behavioral experiments aimed at understanding cooperation have suggested that a willingness to sacrifice one's own interests for the benefit of others, or "costly punishment," may be part of human psychology and evolution. However, because most experiments have been limited to small-scale laboratory populations, the generalizability of these insights to the species has been questioned. In this paper, we report results from 15 diverse populations that show that (i) the propensity to administer costly punishment is unequal between populations and (ii) the propensity to administer costly punishment varies substantially across populations, with little evidence of across-population patterns. These gene-culture correlations of human altruism and costly punishment needs to explain.

For tens of thousands of years before formal contracts, assets, and monetary human societies maintained important forms of cooperation in domains such as hunting, foraging, and food sharing. The scale of cooperation in both contemporary and past human societies remains a puzzle for the evolutionary and social sciences, because, first, neither kin selection nor reciprocity appears to readily explain altruism in very large groups of unrelated individuals and, second, conventional assumptions of self-regarding preferences in economics and related fields appear equally ill-fitted to the facts (1). Reciprocal cooperation can support altruism in large groups; however, some other mechanism is needed to explain why reciprocity should be linked to prosociality rather than selfish or neutral behavior (2). Keen theoretical work



RESEARCH ARTICLES
tions (13). Such experiments have even begun to probe the neural underpinnings of punishment (14, 15).

These results are important, because the propensity of costly punishment can explain many pieces of the puzzle of largescale cooperation. However, like previous field studies, ours was limited to small-scale populations. We conducted all our experiments among university students, not knowing whether such findings for the propensities of students and university students generalize to industrialized societies or whether they indeed capture species characteristics. Our earlier research expanded beyond a 15 diverse societies to measure costly punishment behavior (1, 16). We found that social self-interest could not explain all in any of the 15 societies studied, found much more variation in gene-culture correlations of human altruism and costly punishment than in our previous studies with university students. Similarly, until costly punishment is studied in more societies and in nonuniversity students, it is difficult to be confident in its importance for explaining human behavior.

Given our interest in understanding how widespread costly punishment is, we also examined whether costly punishment with altruistic behavior is valuable for the evolution of costly punishment. In societies in which costly punishment is common will exhibit stronger norms of altruism and prosociality, because the

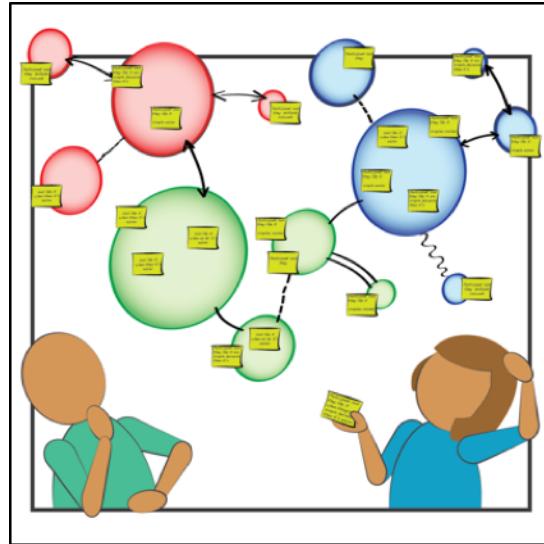


- **Measures of success**

- Early success of simple models
- Success of social network analysis
- Prediction of crowd tipping points



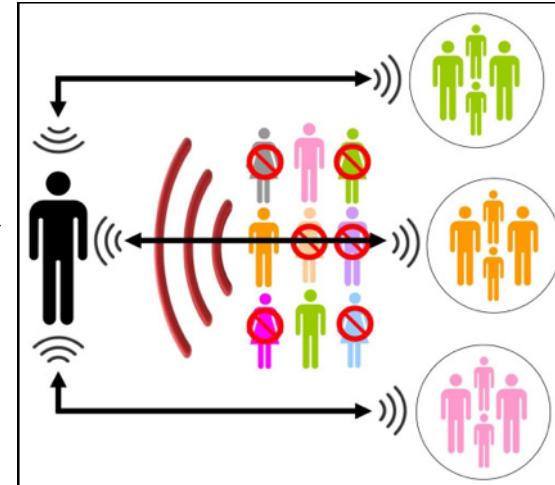
Challenges in Behavioral Analysis



Content
(preference)

Social context
(influence)

Behavioral
Analysis



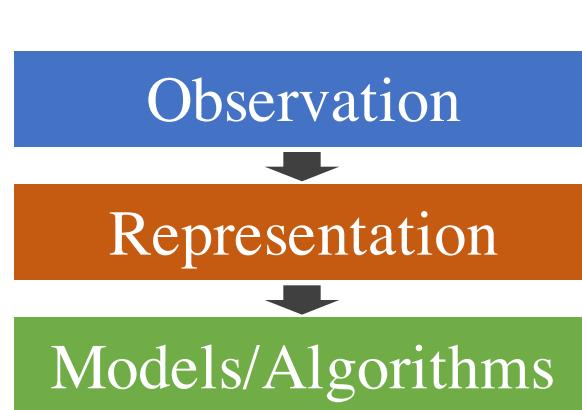
Spatiotemporal context



Intention
(suspiciousness)

REWARDS	# TICKETS GIVEN	CONSEQUENCES	# TICKETS TAKEN AWAY
Extra Math	+5	HITTING	-3
Getting along WELL with others	+3	BULLYING	-4
Good Table Manners	+4	TEASING	-1
LOVE & RESPECT	+5	LYING	-2
Obeying the FIRST TIME	+3	THROWING A FIT	-3
Calm & Quiet in STORE	+3	Ignoring Parents	-4
Extra Reading	+2	SCREAMING or YELLING	-1
CLEANING up after PLAYING	+2	BAD SPORT	-2

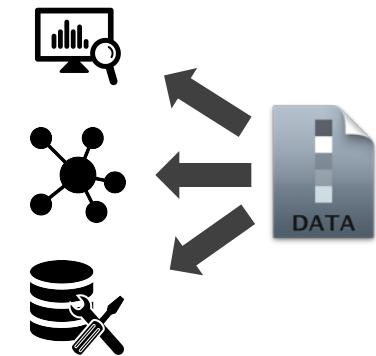
Methodology: Why Data-Driven?



Experience-Driven

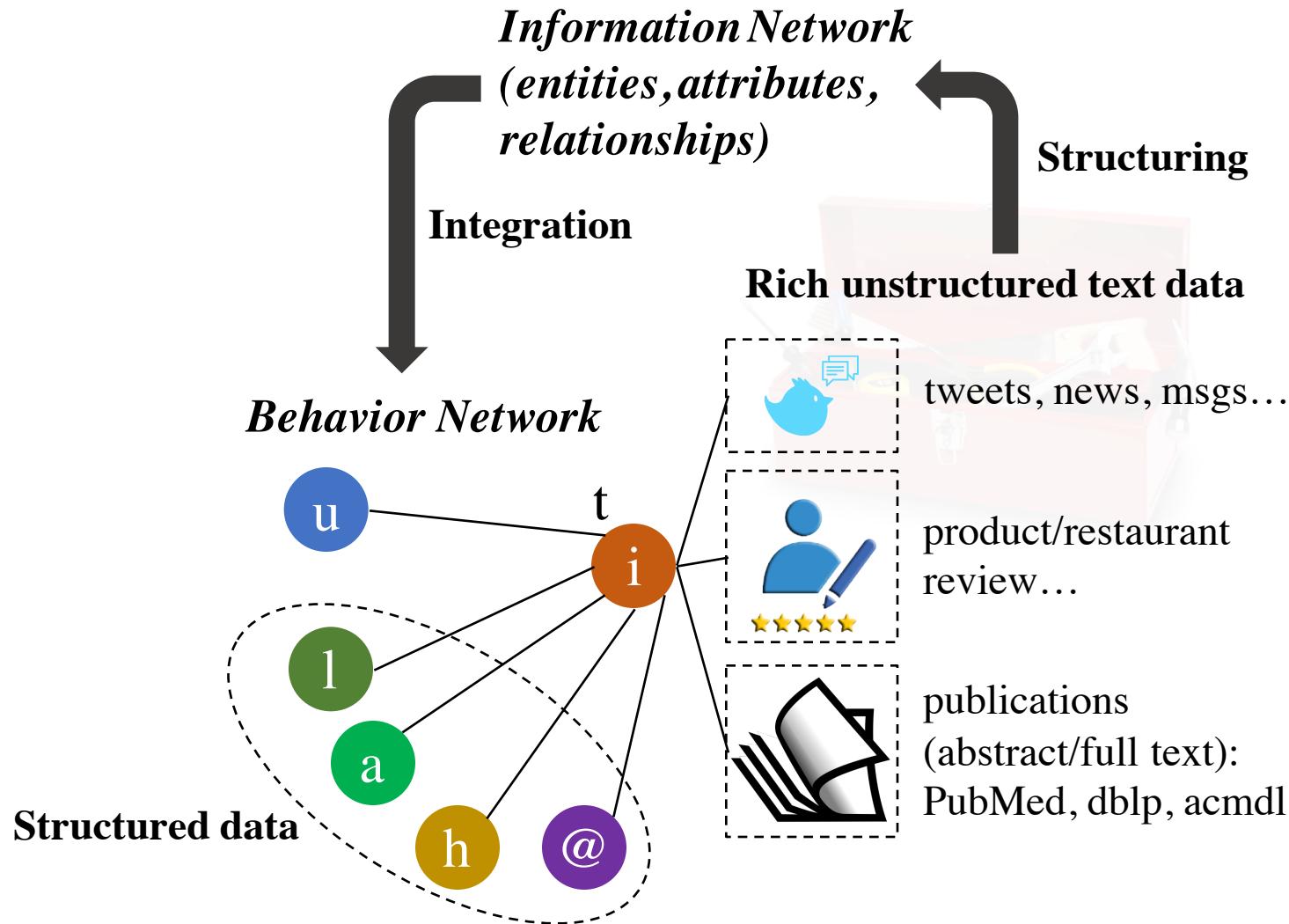


Data-Driven



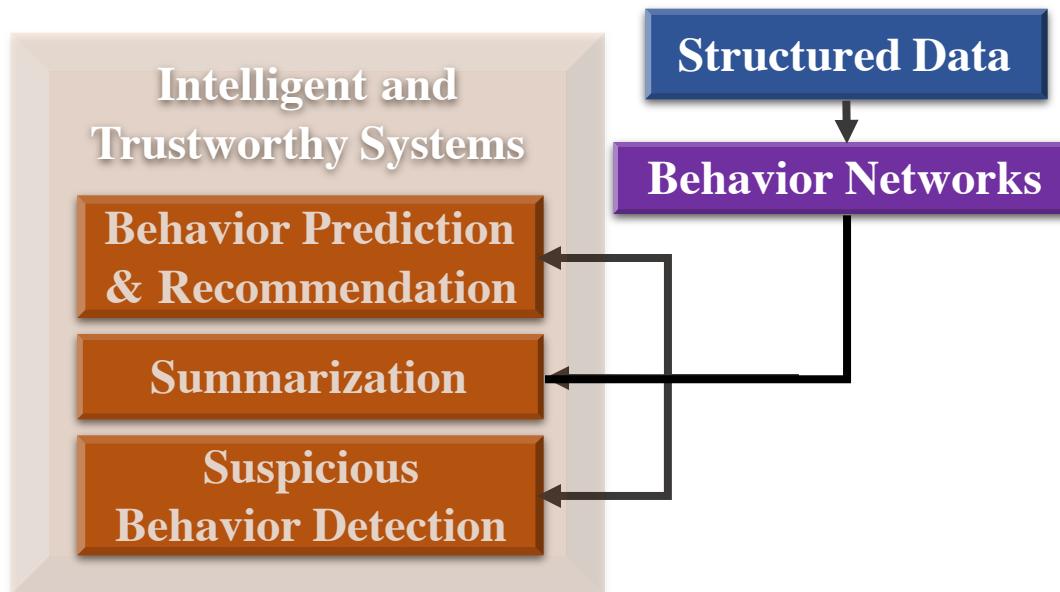
- ❑ **Applications.** Recommender systems, fraud/spam detection.
- ❑ **Representation.** Behavior Network for interaction.
 - ❑ **Nodes:** users/authors, items (*e.g.*, products, tweets, papers), *etc.*
 - ❑ **Links:** (interaction) following, purchasing, tweeting, publishing, *etc.*
 - ❑ **Node attributes:** user profiles, item properties/features, *etc.*
 - ❑ **Link attributes:** similarity, distance, weight, *etc.*

Data to Network to Knowledge



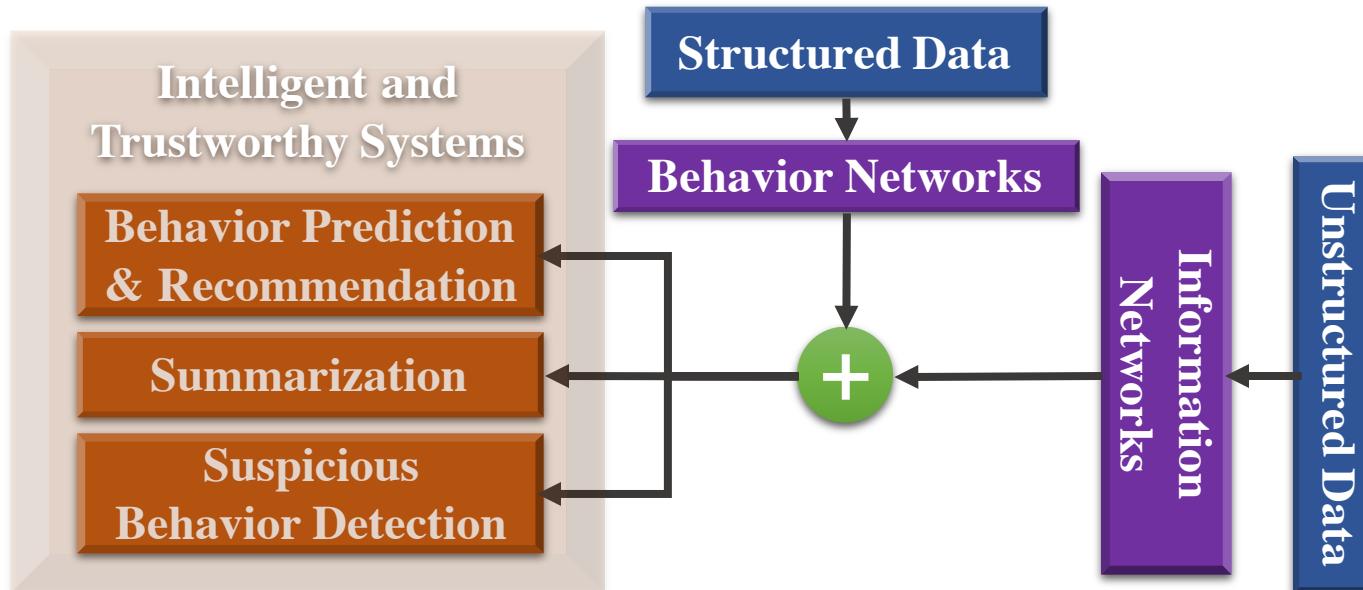
Outline: Data-Driven Behavioral Analytics

- ❑ Mining behavior networks with social and spatiotemporal contexts to support intelligent and trustworthy systems
 - ❑ Mining for behavior prediction and recommendation
 - ❑ Mining for suspicious behavior detection



Outline: Data-Driven Behavioral Analytics

- ❑ Mining behavior networks with social and spatiotemporal contexts to support intelligent and trustworthy systems
 - ❑ Mining for behavior prediction and recommendation
 - ❑ Mining for suspicious behavior detection
- ❑ Structuring behavioral content and integrating behavioral analysis with information networks





I. Mining behavior networks with social and spatiotemporal contexts

I.1. Behavior prediction and recommendation



Behavior in Social Networks

❑ Facebook: Post, Like, Comment, Share

Update Status | Add Photos/Videos | Create Photo Album

What's on your mind?

Public Post

132 Likes 20 Comments



Like



Comment



Share

❑ Twitter: Post, Reply, Retweet, Favorite

What's happening?

Media Location 140 Tweet



5



7



❑ YouTube: Upload, Subscribe, Download, Share, Comment

Upload

Notification bell icon

Top 10 NBA Plays: October 18

NBA Subscribed 6,434,753 Download 720 126,540

Add to Share More

2,468 24

Behavior in Social Networks



Like
Reply
Share
Favorite
Retweet
Comment
Subscribe
Download
Add to
Send
Pin it
Visit
.....



Social Recommender Systems

Huan Liu shared a link.
17 hrs ·

Your Child Is Not Special
We have two choices of when our children can fail: now or later. Now, they are still in a safe environment with people willing to help them succeed. Later, it will be in the context of the workplace or with their...
HUFFINGTONPOST.COM

Like Comment Share
2 people like this.
Write a comment...

Huan Liu and Jiliang Tang like Southwest Airlines.

Southwest Airlines
Sponsored ·
Since some of the other airlines charge you to print your boarding pass, "Find a guy." Or fly Southwest® where #FeesDontFly.
Low fares. Nothing to hide. That's Transparency.

Fee Hacker Tip #6
See more fee hacks
SOUTHWEST.COM
67k Views
24 Likes 2 Comments
Learn More

Like Comment Share

Microsoft Research @MSFTResearch · 3h
@MSFTResearch Labs leader Jeannette Wing on why @Microsoft cares about basic research blogs.technet.com/b/inside_microsoft...

18:06

Carnegie Mellon Retweeted
CNBC's Closing Bell @CNBCClosingBell · Sh
@Kelly_Evans goes behind the wheel of CarnegieMellon's autonomous car. #TheSpark video.cnn.com/gallery/?video...

View summary
4 5 ***

Twitter

Twitter

Pinterest

Recommended



模仿新闻联播嘲讽时弊暴红遭
封杀的相声 - 新闻晚知道 (...)
by ChinaNews360
715,377 views • 3 years ago



Yoga For Weight Loss | Strengthen and Lengthen
by Yoga With Adriene
1,522,005 views • 1 year ago



Yoga For Weight Loss - Hips & Hamstrings
by Yoga With Adriene
219,391 views • 3 months ago



Yoga For Weight Loss - Love Yoga Flow
by Yoga With Adriene
161,772 views • 3 weeks ago



Yoga for Strength and Focus
by Yoga With Adriene
406,306 views • 1 year ago



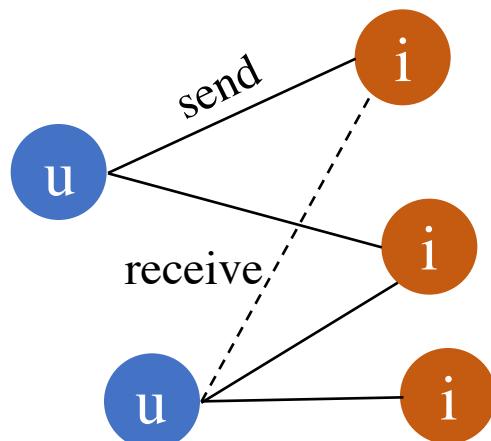
30 Days of Yoga For Your Back - Day 4
by Yoga With Adriene
921,194 views • 9 months ago

Facebook

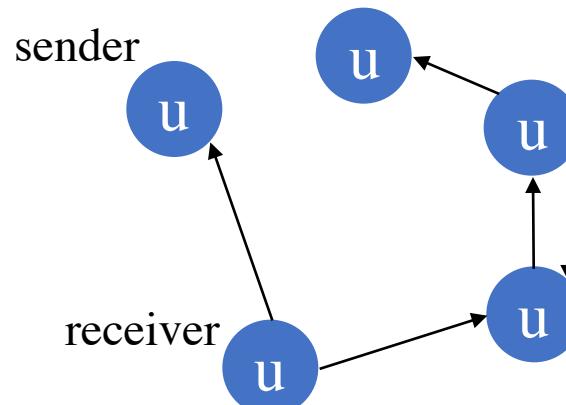
Social Recommender Systems

- ❑ April 20, 2011: Tencent Weibo visited Tsinghua University
 - ❑ Low *conversion rate* (< 6%): #retweets per feed request
 - ❑ Can we build a *social recommender system*?
 - ❑ Given

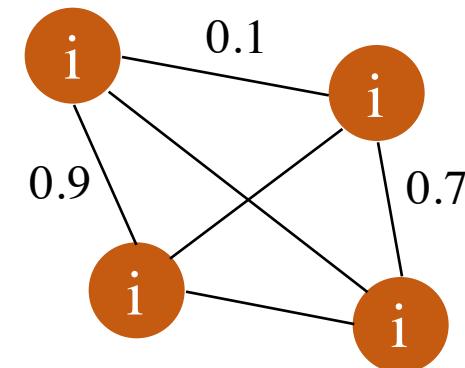
User-item behavior network



User-user social network



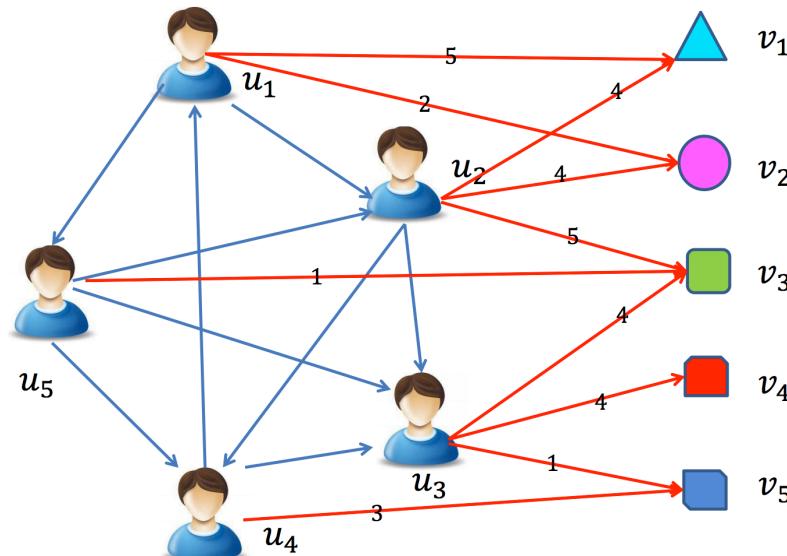
Content similarity
(topic level) [Blei *et al.*]



- ❑ Predict which tweet/item a user will retweet.

Traditional Recommender Systems

- ❑ Assumed that users are independent and identically distributed (user-movie, user-book, *etc.*)



	v_1	v_2	v_3	v_4	v_5
u_1	5	?	2	?	?
u_2	4	4	5	?	?
u_3	?	?	4	4	1
u_4	?	?	?	?	3
u_5	?	?	1	?	?

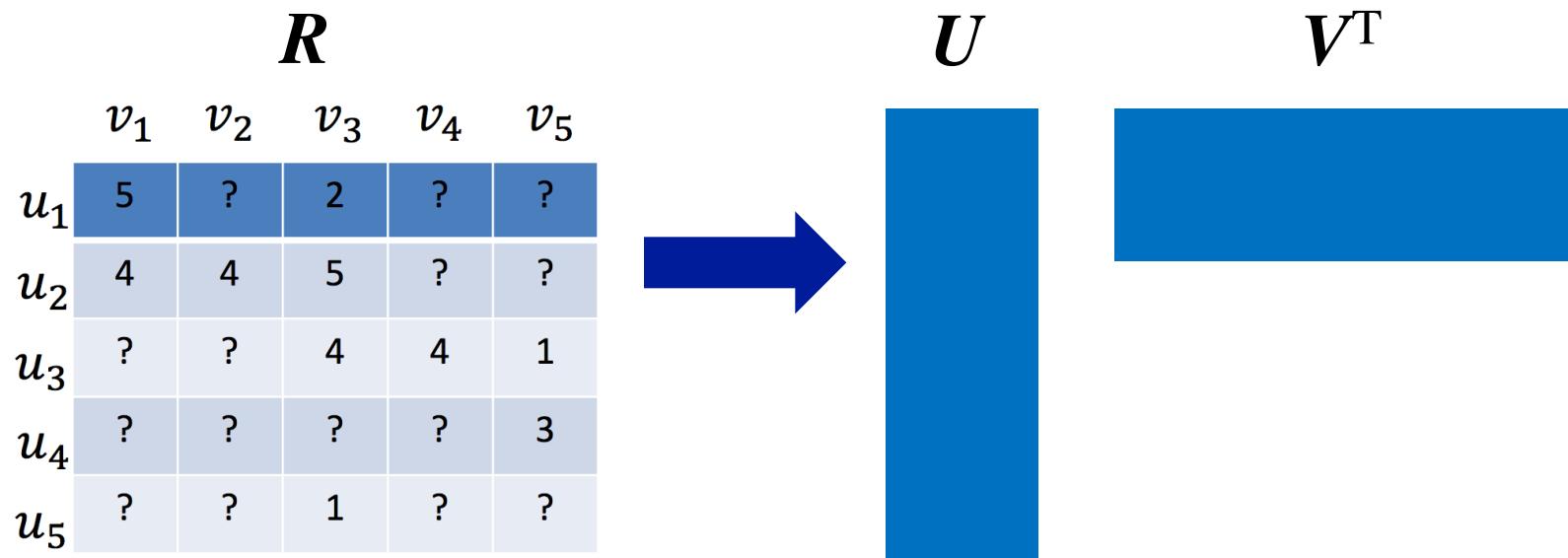


Traditional Recommender Systems

- ❑ Content-based recommender (e.g., TF-IDF)
 - ❑ For textual information (e.g., news, documents)
 - ❑ *Limitation: limited content analysis, over-specialization*
- ❑ Collaborative filtering based recommender
 - ❑ Memory-based CF (e.g., PCC, similarity)
 - ❑ Model-based CF (e.g., factorization based)
 - ❑ *Limitation: data sparsity, cold-start problem*
- ❑ Hybrid recommender system

Matrix Factorization (MF) based CF

- Low-rank MF on the user-item rating matrix R
- User preference vector U
- Item characteristic vector V



Matrix Factorization (MF) based CF

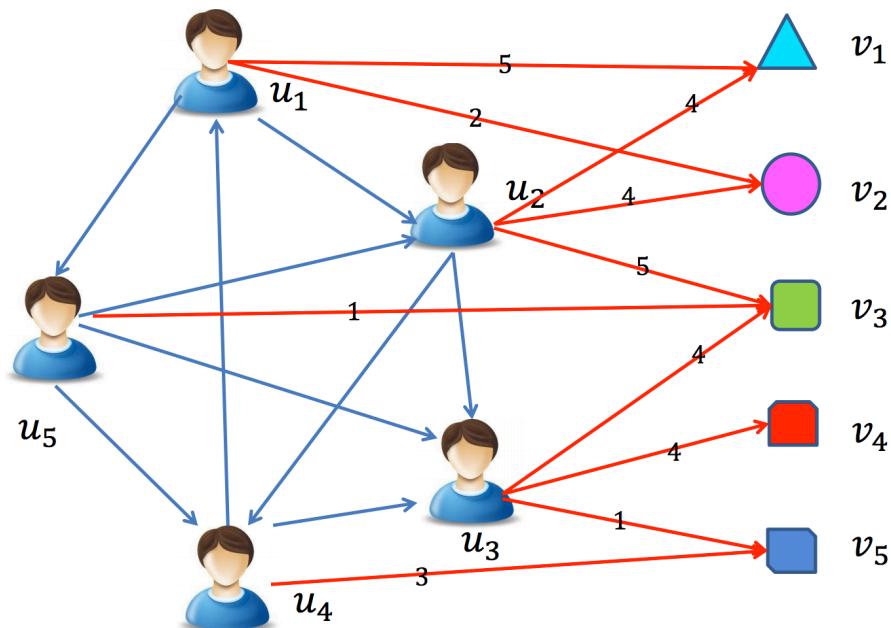
- Low-rank MF on the user-item rating matrix R
- User preference vector U
- Item characteristic vector V
- Observed weight matrix W

$$\min_{\mathbf{U}, \mathbf{V}} \sum_{i=1}^n \sum_{j=1}^m \boxed{\mathbf{W}_{ij}} (\mathbf{R}_{ij} - \mathbf{U}_i \mathbf{V}_j^\top)^2 + \boxed{\alpha(\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)}$$

avoid **over-fitting**,
controlled by the **parameter**

Social Recommendation

Social relations



	u_1	u_2	u_3	u_4	u_5
u_1	0	1	0	0	1
u_2	0	0	1	1	0
u_3	0	0	0	0	0
u_4	1	0	1	0	0
u_5	0	1	1	1	0

	v_1	v_2	v_3	v_4	v_5
u_1	5	?	2	?	?
u_2	4	4	5	?	?
u_3	?	?	4	4	1
u_4	?	?	?	?	3
u_5	?	?	1	?	?

Memory based Social Recommender

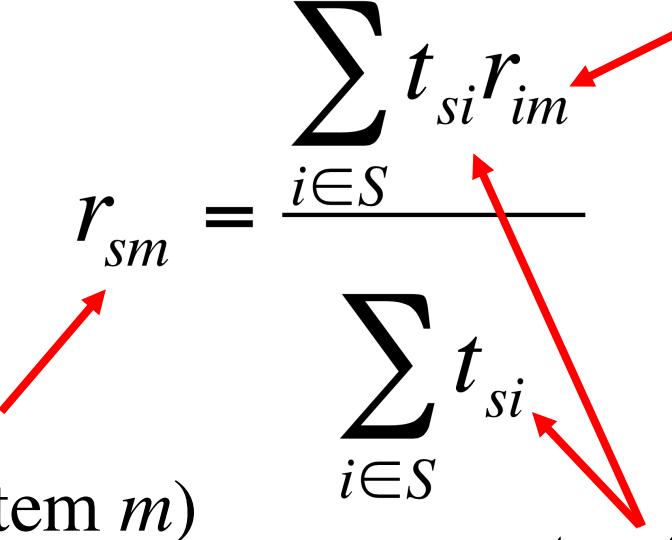
□ TidalTrust

$$r_{sm} = \frac{\sum_{i \in S} t_{si} r_{im}}{\sum_{i \in S} t_{si}}$$

rating (user i , item m)

rating (user s , item m)

trust from social relation (user s , user i)



Memory based Social Recommender

□ MoleTrust

$$p_{a,i} = \bar{r}_a + \frac{\sum_{u=1}^k w_{a,u}(r_{u,i} - \bar{r}_u)}{\sum_{u=1}^k w_{a,u}}$$

average rating (user a)

rating (user u , item i)

predicted rating (user a , item i)

trust from social relation (user a , user u)

average rating (user u)

The diagram illustrates the MoleTrust formula with red arrows indicating the flow of information from the labels to the terms in the equation. The labels are: 'average rating (user a)', 'rating (user u , item i)', 'predicted rating (user a , item i)', and 'trust from social relation (user a , user u)'. The first two labels point to the terms \bar{r}_a and $r_{u,i}$ respectively. The third label points to the entire fraction. The fourth label points to the term $w_{a,u}$.

Memory based Social Recommender

□ TrustWalker

probability of
user u 's random walk
from item i to item j

$$P(Y_{u,i} = j) = \frac{sim(i, j)}{\sum_{l \in RI_u} sim(i, l)}$$

similarity measure
(item i , item j)

Pearson correlation
of (item i , item j)

$$sim(i, j) = \frac{1}{1 + e^{-\frac{|UC_{i,j}|}{2}}} \times corr(i, j)$$

common user set
of (item i , item j)



Model based Social Recommender

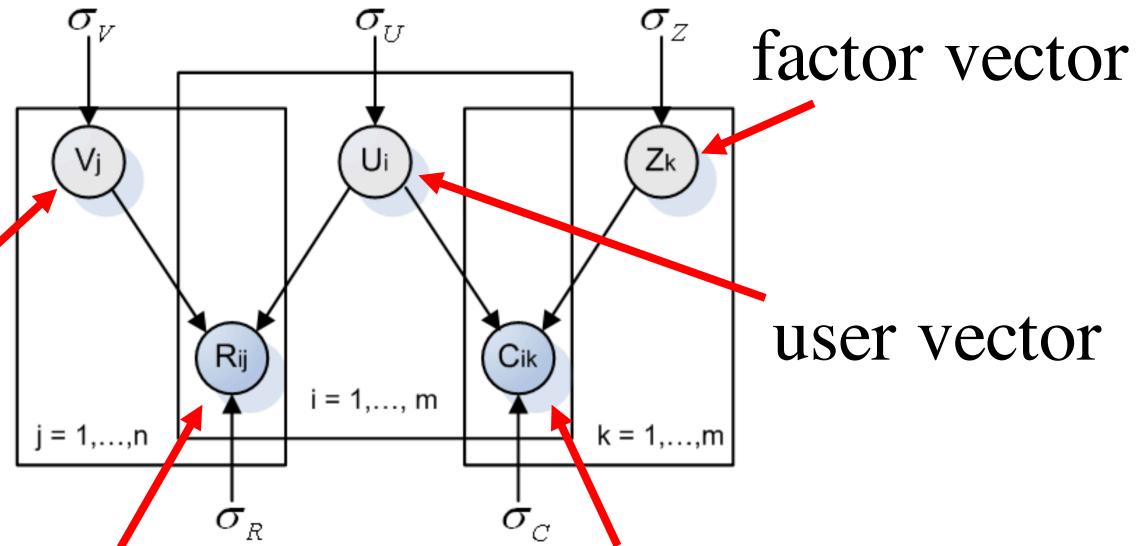
- Optimization methods such as gradient based methods can be applied to find a well-worked optimal solution.
- MF has a nice probabilistic interpretation with Gaussian noise.
- MF is very flexible and allows us to include prior knowledge.

$$\begin{aligned} & \textit{Social Recommendation CF} \\ &= \textit{Basic CF} + \textit{Social Information Model} \end{aligned}$$

Model based Social Recommender

□ SoRec

item vector



factor vector

user vector

R : user-item
rating matrix

	v_1	v_2	v_3	v_4	v_5
u_1	5	?	2	?	?
u_2	4	4	5	?	?
u_3	?	?	4	4	1
u_4	?	?	?	?	3
u_5	?	?	1	?	?

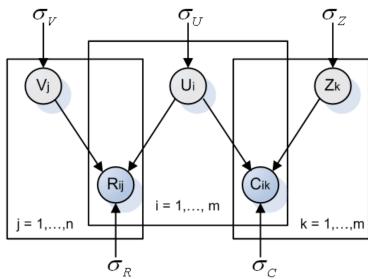
C : user-user
social matrix

	u_1	u_2	u_3	u_4	u_5
u_1	0	1	0	0	1
u_2	0	0	1	1	0
u_3	0	0	0	0	0
u_4	1	0	1	0	0
u_5	0	1	1	1	0

Model based Social Recommender

□ SoRec

$$p(\mathcal{C}|U, V, \sigma_R^2) = \prod_{i=1}^m \prod_{j=1}^n \mathcal{N} \left[\left(r_{ij} | g(U_i^T V_j), \sigma_R^2 \right) \right]^{I_{ij}^R}$$



Gaussian distribution

Logistic function Observed

$$p(C|U, Z, \sigma_C^2) = \prod_{i=1}^m \prod_{k=1}^m \mathcal{N} \left[\left(c_{ik} | g(U_i^T Z_k), \sigma_C^2 \right) \right]^{I_{ik}^C}$$

Model based Social Recommender

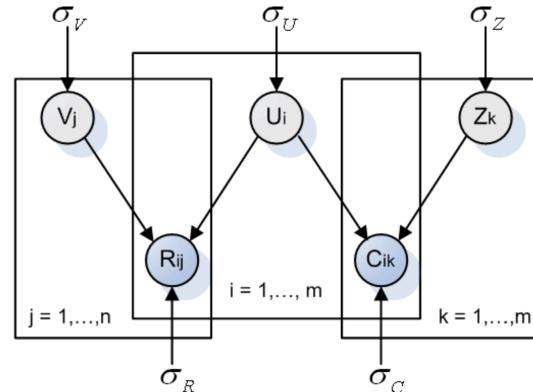
□ SoRec

behavioral term

$$\mathcal{L}(R, C, U, V, Z) =$$

$$\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n I_{ij}^R \underbrace{(r_{ij} - g(U_i^T V_j))^2}_{social\ term} + \frac{\lambda_C}{2} \sum_{i=1}^m \sum_{k=1}^m I_{ik}^C \underbrace{(c_{ik}^* - g(U_i^T Z_k))^2}_{regularization\ terms}$$

$$+ \frac{\lambda_U}{2} \|U\|_F^2 + \frac{\lambda_V}{2} \|V\|_F^2 + \frac{\lambda_Z}{2} \|Z\|_F^2, \quad (9)$$



Model based Social Recommender

□ SoRec

Gradient Descent Methods

$$\frac{\partial \mathcal{L}}{\partial U_i} = \sum_{j=1}^n I_{ij}^R g'(U_i^T V_j) \underline{(g(U_i^T V_j) - r_{ij}) V_j}$$

deviate of
Logistic
function

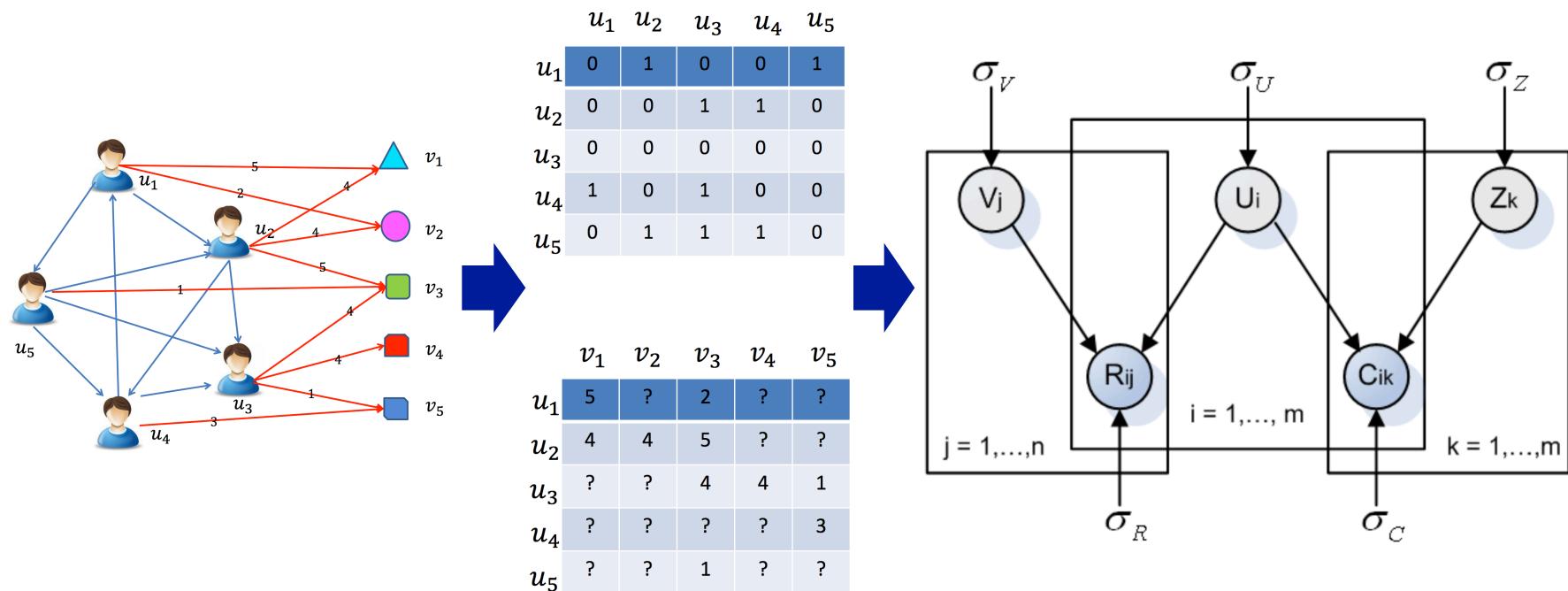
$$+ \lambda_C \sum_{k=1}^m I_{ik}^C g'(U_i^T Z_k) \underline{(g(U_i^T Z_k) - c_{ik}^*) Z_k} + \lambda_U U_i,$$

$$\frac{\partial \mathcal{L}}{\partial V_j} = \sum_{i=1}^m I_{ij}^R g'(U_i^T V_j) \underline{(g(U_i^T V_j) - r_{ij}) U_i} + \lambda_V V_j,$$

$$\frac{\partial \mathcal{L}}{\partial Z_k} = \lambda_C \sum_{i=1}^m I_{ik}^C g'(U_i^T Z_k) \underline{(g(U_i^T Z_k) - c_{ik}^*) U_i} + \lambda_Z Z_k, (10)$$

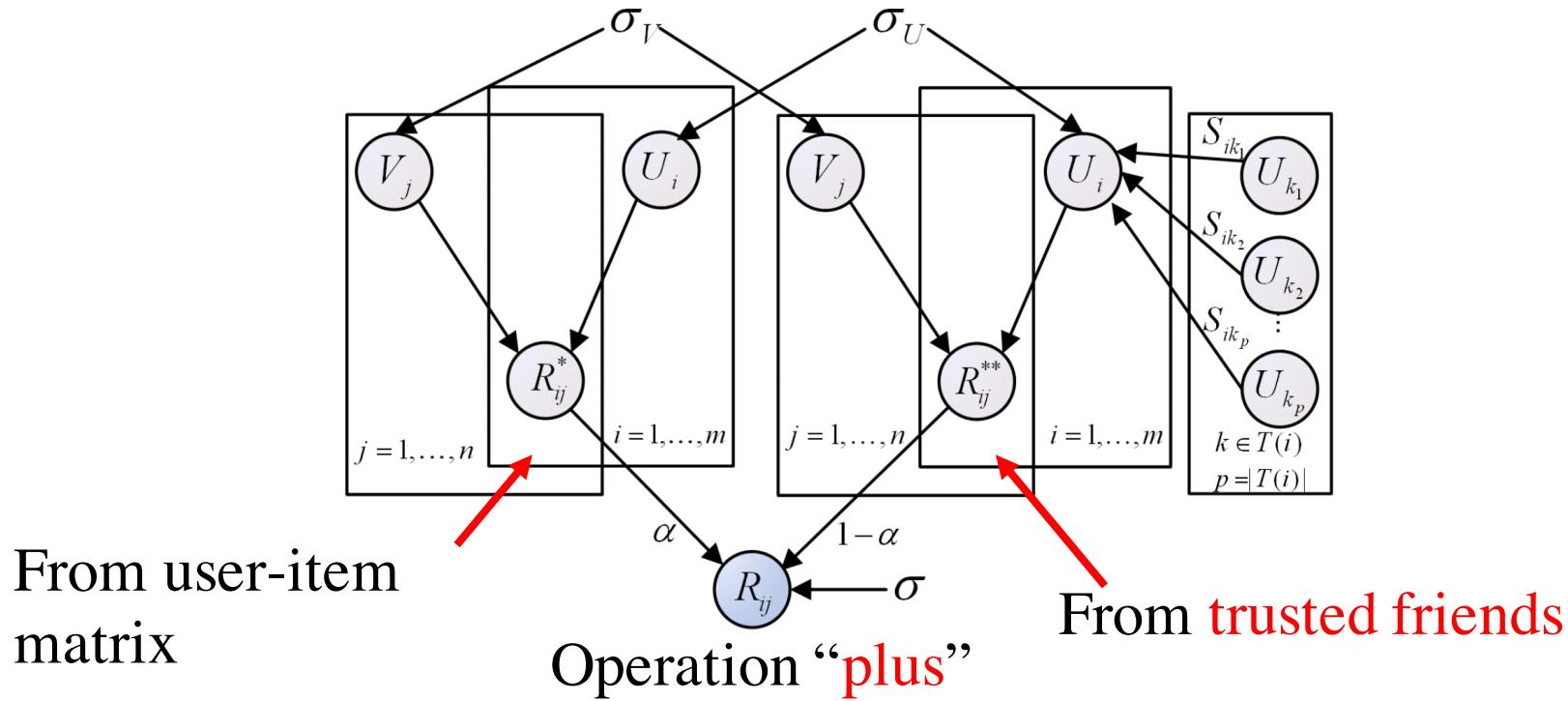
Model based Social Recommender

□ SoRec



Model based Social Recommender

- Replacing social with trust
- “Social Trust” Ensemble for Epinion data



Model based Social Recommender

□ “Social Trust” Ensemble

$$\begin{aligned} \mathcal{L}(R, S, U, V) & \quad \text{From user-item matrix} \\ &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n I_{ij}^R (R_{ij} - g(\underline{\alpha U_i^T V_j} + \underline{(1-\alpha) \sum_{k \in \mathcal{T}(i)} S_{ik} U_k^T V_j}))^2 \\ &+ \frac{\lambda_U}{2} \|U\|_F^2 + \frac{\lambda_V}{2} \|V\|_F^2, \end{aligned} \tag{13}$$

Model based Social Recommender

□ “Social Trust” Ensemble

*Gradient
Descent
Methods*

$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial U_i} = & \alpha \sum_{j=1}^n I_{ij}^R g'(\alpha U_i^T V_j + (1 - \alpha) \sum_{k \in \mathcal{T}(i)} S_{ik} U_k^T V_j) V_j \\
 & \times (g(\alpha U_i^T V_j + (1 - \alpha) \sum_{k \in \mathcal{T}(i)} S_{ik} U_k^T V_j) - R_{ij}) \\
 & + (1 - \alpha) \sum_{p \in \mathcal{B}(i)} \sum_{j=1}^n I_{pj}^R g'(\alpha U_p^T V_j + (1 - \alpha) \sum_{k \in \mathcal{T}(p)} S_{pk} U_k^T V_j) \\
 & \times (g(\alpha U_p^T V_j + (1 - \alpha) \sum_{k \in \mathcal{T}(p)} S_{pk} U_k^T V_j) - R_{pj}) S_{pi} V_j + \lambda_U U_i, \\
 \frac{\partial \mathcal{L}}{\partial V_j} = & \sum_{i=1}^m I_{ij}^R g'(\alpha U_i^T V_j + (1 - \alpha) \sum_{k \in \mathcal{T}(i)} S_{ik} U_k^T V_j) \\
 & \times (g(\alpha U_i^T V_j + (1 - \alpha) \sum_{k \in \mathcal{T}(i)} S_{ik} U_k^T V_j) - R_{ij}) \\
 & \times (\alpha U_i + (1 - \alpha) \sum_{k \in \mathcal{T}(i)} S_{ik} U_k^T) + \lambda_V V_j, \tag{14}
 \end{aligned}$$

Model based Social Recommender

□ SoReg

Average-based regularization:

Regularize with the average of friends' tastes

$$\min_{U, V} \mathcal{L}_1(R, U, V) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n I_{ij} (R_{ij} - U_i^T V_j)^2$$



$$+ \frac{\alpha}{2} \sum_{i=1}^m \|U_i - \frac{\sum_{f \in \mathcal{F}^+(i)} \text{Sim}(i, f) \times U_f}{\sum_{f \in \mathcal{F}^+(i)} \text{Sim}(i, f)}\|_F^2,$$

$$+ \frac{\lambda_1}{2} \|U\|_F^2 + \frac{\lambda_2}{2} \|V\|_F^2. \quad (8)$$

Information loss: Friends may have diverse tastes!!!

Model based Social Recommender

□ SoReg

Individual-based regularization:

Regularize with friends individually

$$\begin{aligned} \min_{U, V} \mathcal{L}_2(R, U, V) &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n I_{ij} (R_{ij} - U_i^T V_j)^2 \\ &+ \frac{\beta}{2} \sum_{i=1}^m \sum_{f \in \mathcal{F}^+(i)} Sim(i, f) \|U_i - U_f\|_F^2 \\ &+ \lambda_1 \|U\|_F^2 + \lambda_2 \|V\|_F^2. \end{aligned} \tag{11}$$


Related Work

	Behavior	Content	Social	Trust
Collaborative filtering (CF) [Herlocker <i>et al.</i> . TOIS; Koren KDD]	✓			
Content-based filtering with CF [Balabanovic <i>et al.</i> ; Liu <i>et al.</i> . CIKM;]	✓	✓		
SoRec [Ma <i>et al.</i> . CIKM, TIS] SoReg [Ma <i>et al.</i> . WSDM]	✓		✓	
Trust-based methods [Massa <i>et al.</i> . RecSys; Jamali <i>et al.</i> . KDD; Ma <i>et al.</i> . SIGIR, TIST]	✓			✓

❑ Q: What are the **factors** of users' decisions on retweeting?
Can we **observe** them from the data? How to **integrate** the information for accurate prediction?

Observation: Social Contextual Factors

- Will Michelle Obama share this message?
- Please list your reasons.



Barack Obama

Happy birthday, Michelle Obama!

[Like](#) · [Comment](#) · [Share](#) · January 18, 2013

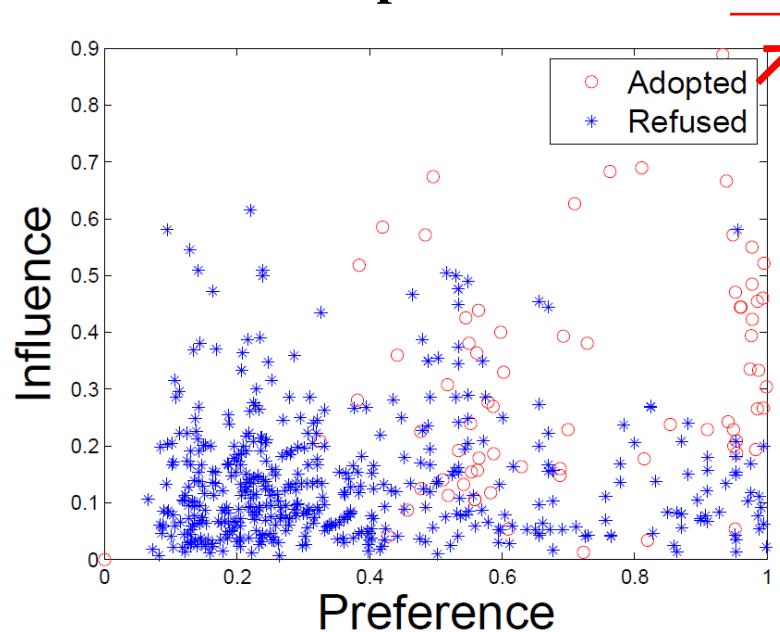
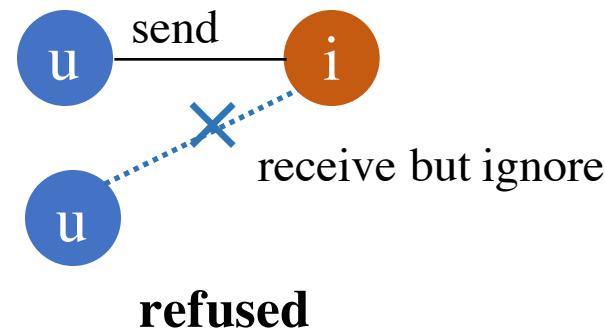
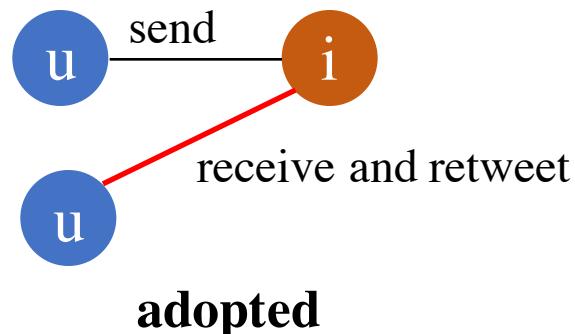


Michelle Obama shared Barack Obama's photo.

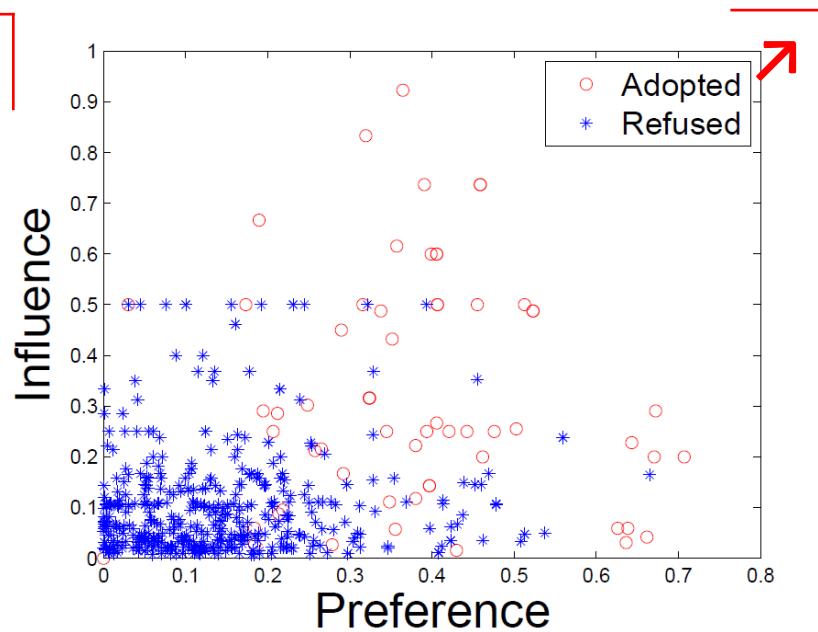
January 18, 2013 ·



Observation: Social Contextual Factors



China's Facebook: Renren



China's Twitter: Tencent Weibo

Representation: From Contextual Information to Contextual Factors

Content

Item-item similarity

Item latent features V

Behavior

User-item interaction

User latent features U

Social

User-user social relation

Item sender G

Interaction frequency

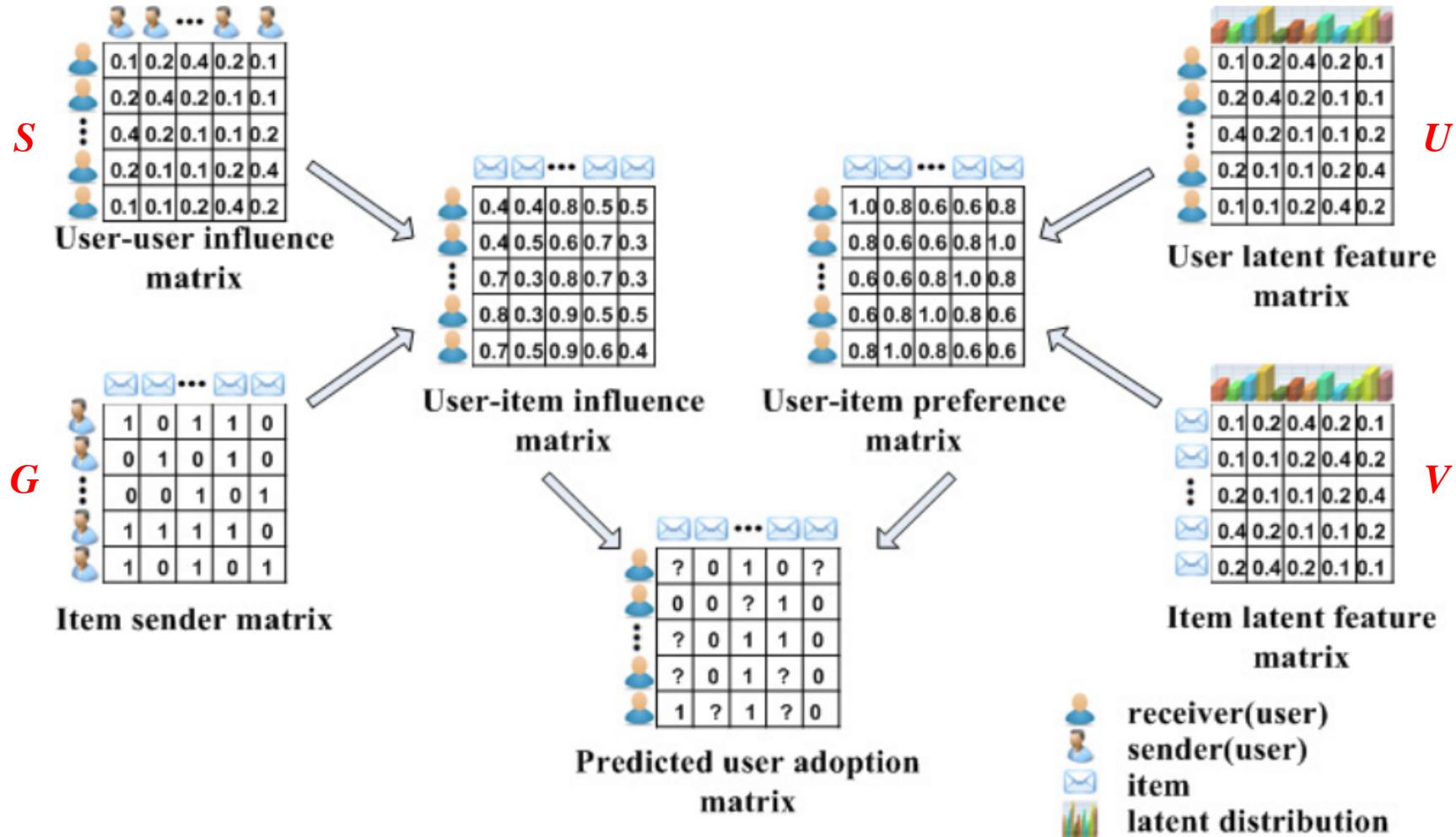
User-user interaction

User-user influence S

Personal preference
on the given item

Interpersonal influence
from the item's sender

Model: ContextMF



Model: ContextMF

behavior influence preference

$$P(\mathbf{R}|\mathbf{S}, \mathbf{U}, \mathbf{V}, \sigma_R^2) = \prod_{i=1}^M \prod_{j=1}^N \mathcal{N}(\underline{\mathbf{R}_{ij}} | \underline{\mathbf{S}_i \mathbf{G}_j^\top} \odot \underline{\mathbf{U}_i^\top \mathbf{V}_j}, \sigma_R^2)$$

behavior interaction frequency/trust

item content

$$\begin{aligned} \mathcal{J} = & ||\mathbf{R} - \mathbf{S}\mathbf{G}^\top \odot \mathbf{U}^\top \mathbf{V}||_F^2 + \alpha ||\mathbf{W} - \mathbf{U}^\top \mathbf{U}||_F^2 \\ & + \beta ||\mathbf{C} - \mathbf{V}^\top \mathbf{V}||_F^2 + \gamma ||\mathbf{S} - \mathbf{F}||_F^2 \\ & + \delta ||\mathbf{S}||_F^2 + \eta ||\mathbf{U}||_F^2 + \lambda ||\mathbf{V}||_F^2 \end{aligned}$$

social relation

Model: ContextMF

- Gradient descent method

$$\frac{\partial \mathcal{J}}{\partial \mathbf{S}} = 2 \left(-\mathbf{R}(\mathbf{G} \odot \mathbf{V}^\top \mathbf{U}) + (\mathbf{S}\mathbf{G}^\top \odot \mathbf{U}^\top \mathbf{V})\mathbf{G} \right. \\ \left. + \gamma(\mathbf{S} - \mathbf{F}) + \delta\mathbf{S} \right)$$

$$\frac{\partial \mathcal{J}}{\partial \mathbf{U}} = 2 \left(-\mathbf{V}\mathbf{R}^\top + \mathbf{V}(\mathbf{G}\mathbf{S}^\top \odot \mathbf{V}^\top \mathbf{U}) - 2\alpha\mathbf{U}\mathbf{W} \right. \\ \left. + 2\alpha\mathbf{U}\mathbf{U}^\top \mathbf{U} + \eta\mathbf{U} \right)$$

$$\frac{\partial \mathcal{J}}{\partial \mathbf{V}} = 2 \left(-\mathbf{U}\mathbf{R} + \mathbf{U}(\mathbf{S}\mathbf{G}^\top \odot \mathbf{U}^\top \mathbf{V}) - 2\beta\mathbf{V}\mathbf{C} \right. \\ \left. + 2\beta\mathbf{V}\mathbf{V}^\top \mathbf{V} + \lambda\mathbf{V} \right)$$



Experimental Results

Method	MAE	RMSE	$\hat{\tau}$	$\hat{\rho}$
Renren Dataset				
Content-based [1]	0.3842	0.4769	0.5409	0.5404
Item CF [25]	0.3601	0.4513	0.5896	0.5988
FeedbackTrust [22]	0.3764	0.4684	0.5433	0.5469
Influence-based [9]	0.3859	0.4686	0.5394	0.5446
SoRec [19]	0.3276	0.4127	0.6168	0.6204
SoReg [20]	0.2985	0.3537	0.7086	0.7140
Influence MF	0.3102	0.3771	0.6861	0.7006
Preference MF	0.3032	0.3762	0.6937	0.7036
Context MF	0.2416	0.3086	0.7782	0.7896

Tencent Weibo Dataset

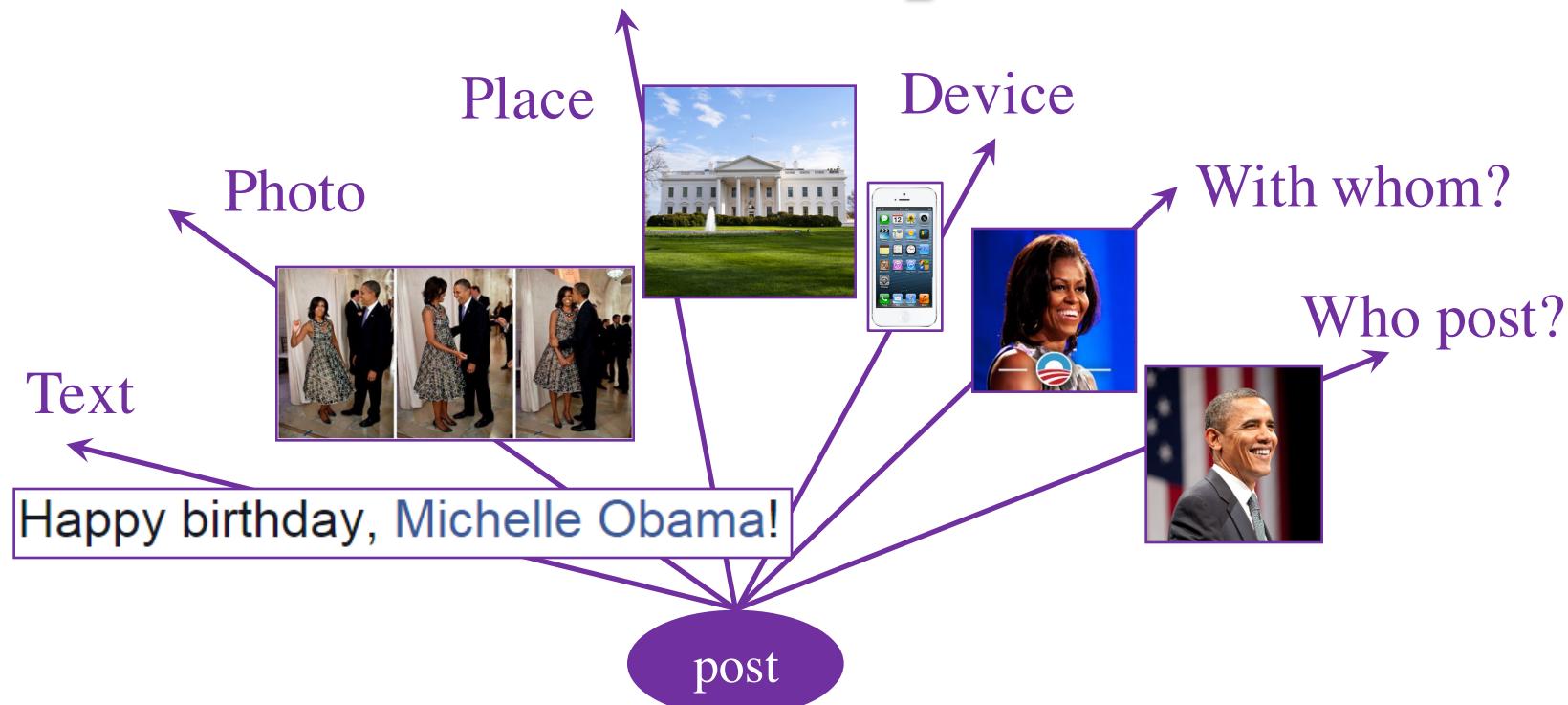
Content-based [1]	0.2576	0.3643	0.7728	0.7777
Item CF [25]	0.2375	0.3372	0.7867	0.8049
FeedbackTrust [22]	0.2830	0.3887	0.7094	0.7115
Influence-based [9]	0.2651	0.3813	0.7163	0.7275
SoRec [19]	0.2256	0.3325	0.7973	0.8064
SoReg [20]	0.1997	0.2962	0.8390	0.8423
Influence MF	0.2183	0.3206	0.8179	0.8258
Preference MF	0.2111	0.3088	0.8384	0.8453
Context MF	0.1514	0.2348	0.8570	0.8685

vs. SoReg [TIST'11]	Renren	Tencent Weibo
MAE	$\downarrow 19.1\%$	$\downarrow 24.2\%$
RMSE	$\downarrow 12.8\%$	$\downarrow 20.7\%$
Kendall's	$\uparrow 9.82\%$	$\uparrow 2.1\%$
Spearman's	$\uparrow 10.6\%$	$\uparrow 3.1\%$

□ **Deployed in Weibo News Feed.** Improved conversion rate from 5.78% to 8.27% (relatively **43%**).

□ #citations = **149**

Observation: Spatial Context



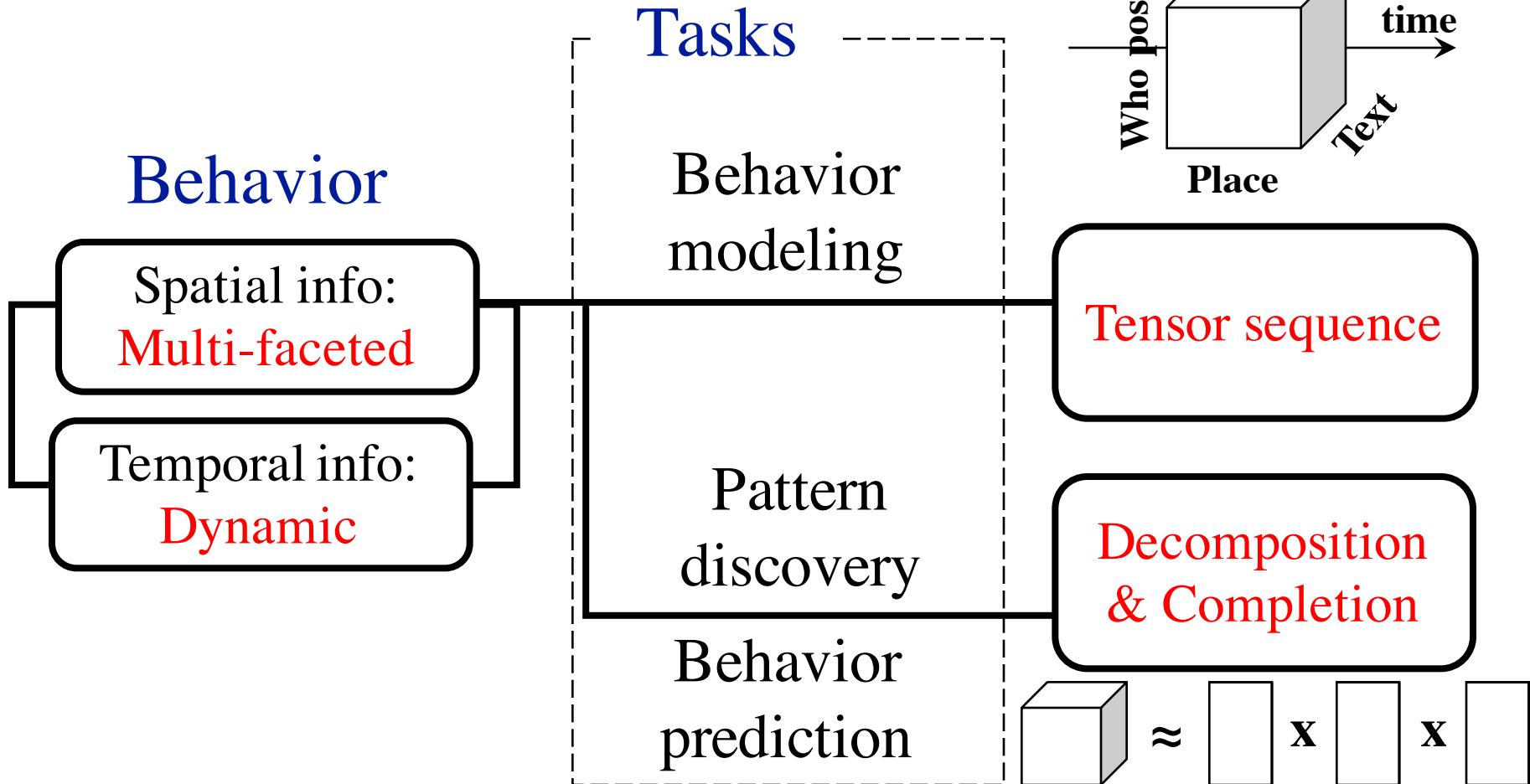
Jan. 18
Birthday party
@ White house



Observation: Temporal Context

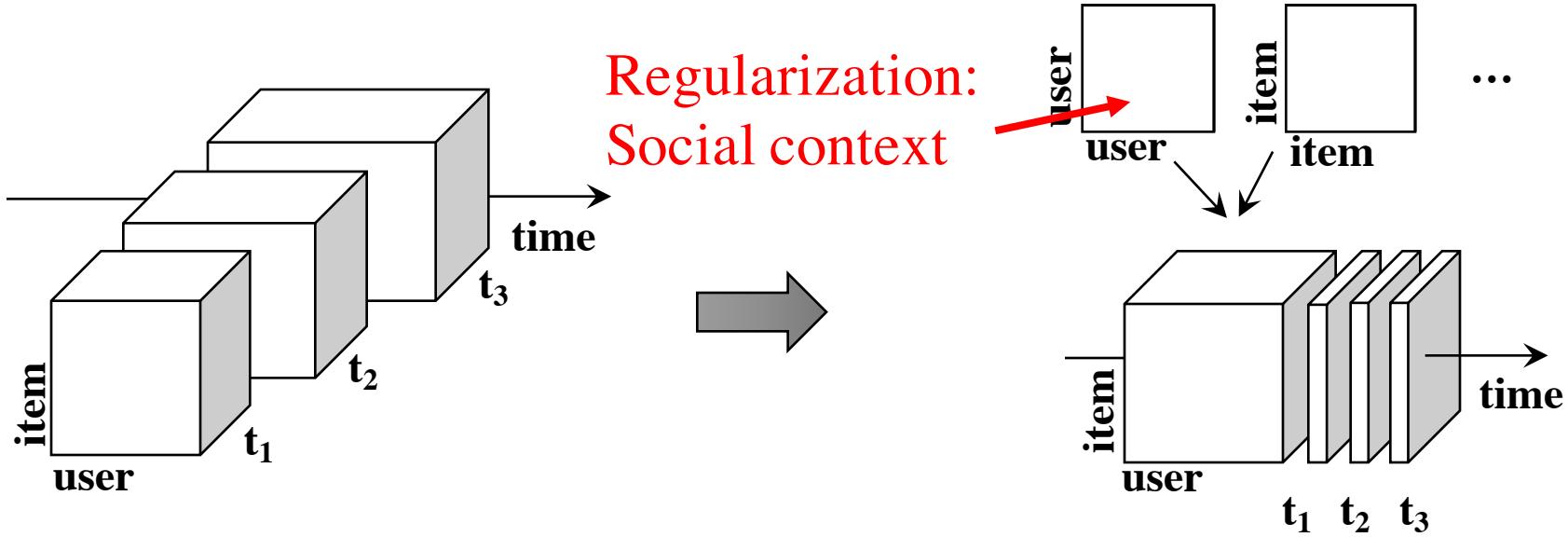


Representation: Tensor Sequence



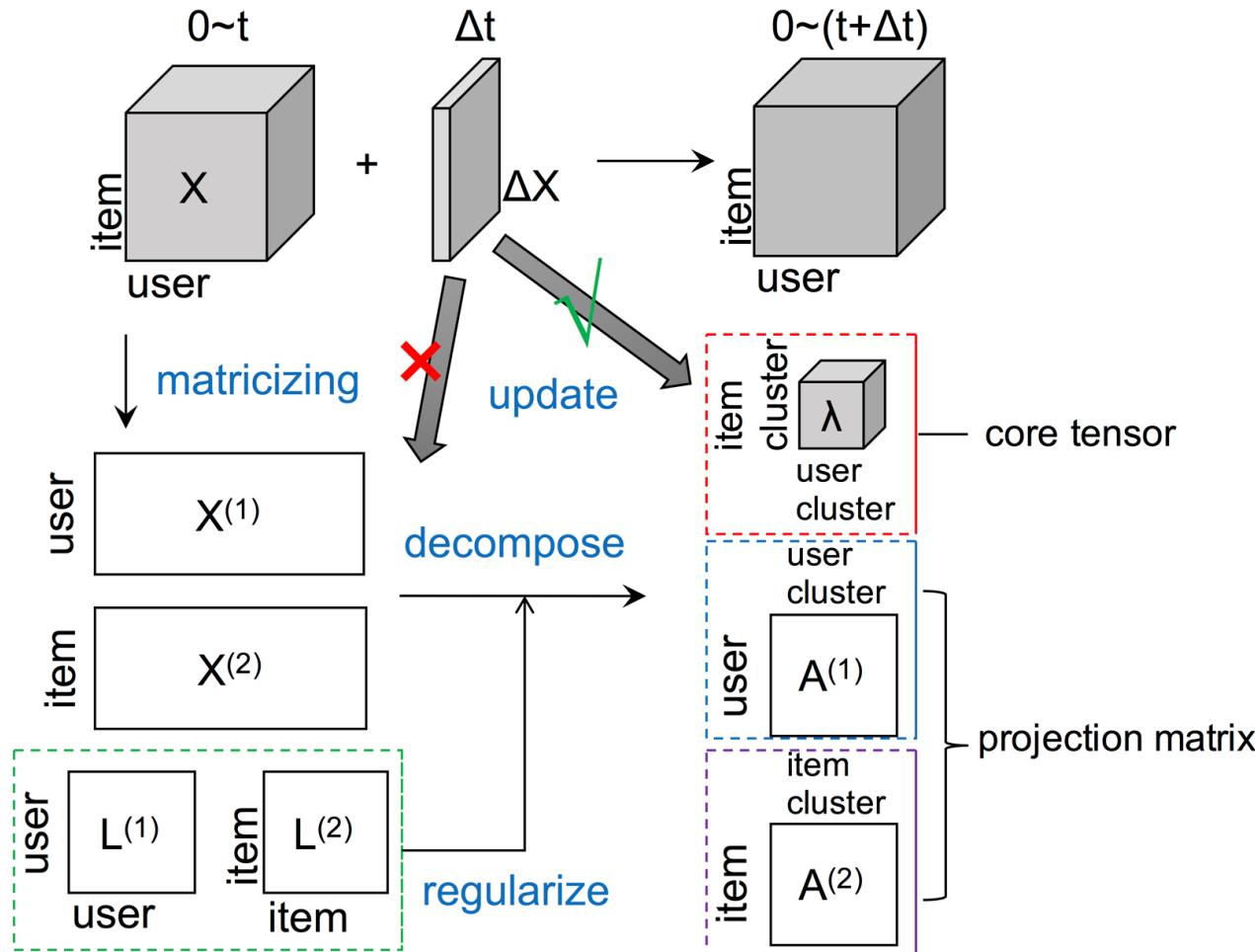
Challenges: Sparsity and Complexity

- Addressing **sparsity**: *Flexible regularization with auxiliary data*
- Addressing **high complexity**: *Incremental updates for projection matrix*



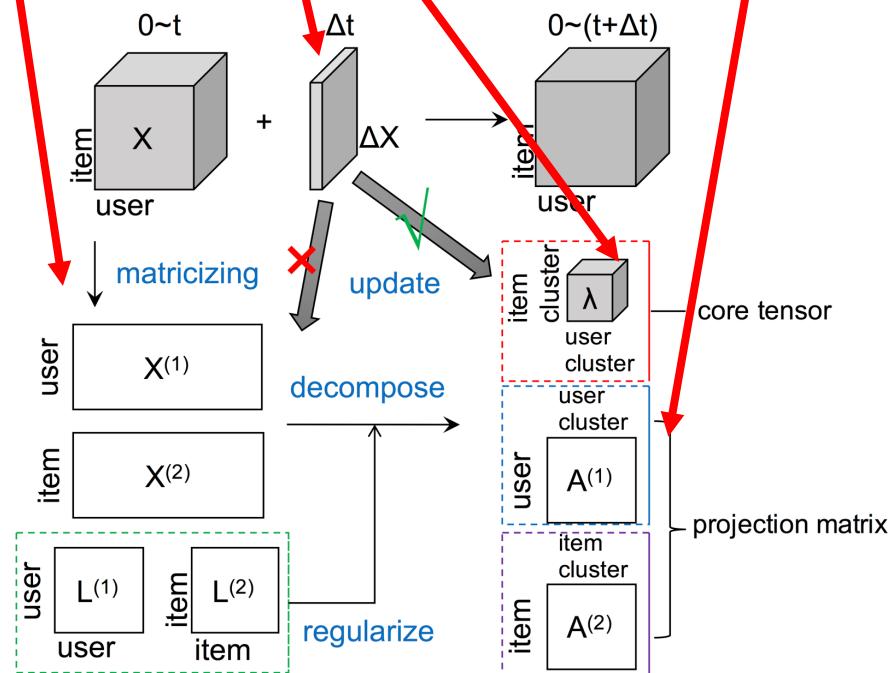
Model: FEMA

Flexible Evolutionary Multi-faceted Analysis



Tensor Perturbation Theory

$$[(\mathbf{X}^{(m)} + \Delta\mathbf{X}^{(m)})(\mathbf{X}^{(m)} + \Delta\mathbf{X}^{(m)})^\top + \mu^{(m)} \mathbf{L}^{(m)}] \cdot (\mathbf{a}_i^{(m)} + \Delta\mathbf{a}_i^{(m)}) = (\lambda_i^{(m)} + \Delta\lambda_i^{(m)}) (\mathbf{a}_i^{(m)} + \Delta\mathbf{a}_i^{(m)})$$



Algorithm: FEMA

Approximation

Require: $\mathcal{X}_t, \Delta\mathcal{X}_t, \mathbf{A}_t^{(m)}|_{m=1}^M, \lambda_t^{(m)}|_{m=1}^M$

for $m = 1, \dots, M$ **do**

for $i = 1, \dots, r^{(m)}$ **do**

 Compute $\Delta\lambda_{t,i}^{(m)}$ using

$$\Delta\lambda_i^{(m)} = \mathbf{a}_i^{(m)\top} (\mathbf{X}^{(m)} \Delta\mathbf{X}^{(m)\top} + \Delta\mathbf{X}^{(m)} \mathbf{X}^{(m)\top}) \mathbf{a}_i^{(m)}$$

 and compute

$$\lambda_{t+1,i}^{(m)} = \lambda_{t,i}^{(m)} + \Delta\lambda_{t,i}^{(m)};$$

 Compute $\Delta\mathbf{a}_{t,i}^{(m)}$ using

$$\Delta\mathbf{a}_i^{(m)} = \sum_{j \neq i} \frac{\mathbf{a}_j^{(m)\top} (\mathbf{X}^{(m)} \Delta\mathbf{X}^{(m)\top} + \Delta\mathbf{X}^{(m)} \mathbf{X}^{(m)\top}) \mathbf{a}_i^{(m)}}{\lambda_i^{(m)} - \lambda_j^{(m)}} \mathbf{a}_j^{(m)}$$

 and compute

$$\mathbf{a}_{t+1,i}^{(m)} = \mathbf{a}_{t,i}^{(m)} + \Delta\mathbf{a}_{t,i}^{(m)} \text{ and } \mathbf{A}_{t+1}^{(m)} = \{\mathbf{a}_{t+1,i}^{(m)}\};$$

end for

end for

$$\mathcal{Y}_{t+1} = (\mathcal{X}_t + \Delta\mathcal{X}_t) \prod_{m=1}^M \times_{(m)} \mathbf{A}_{t+1}^{(m)\top};$$

return $\mathbf{A}_{t+1}^{(m)}|_{m=1}^M, \lambda_{t+1}^{(m)}|_{m=1}^M, \mathcal{Y}_{t+1}$

Bound Guarantee

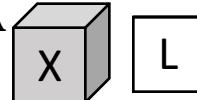
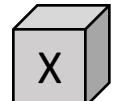
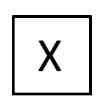
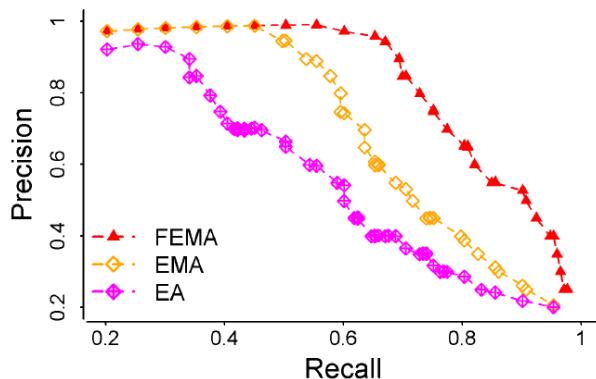
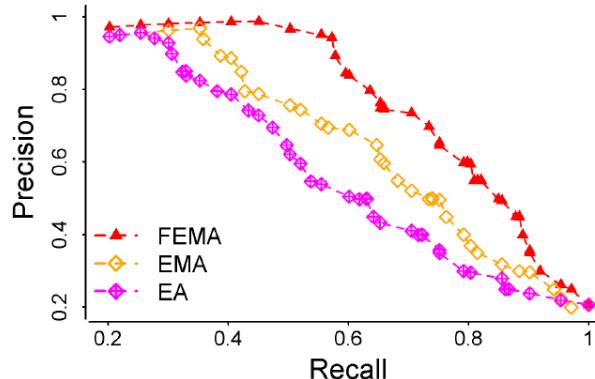
core tensor

$$|\Delta\lambda_i^{(m)}| \leq 2(\lambda_{\mathbf{X}^{(m)\top} \mathbf{X}^{(m)}}^{\max})^{\frac{1}{2}} \|\Delta\mathbf{X}^{(m)}\|_2$$

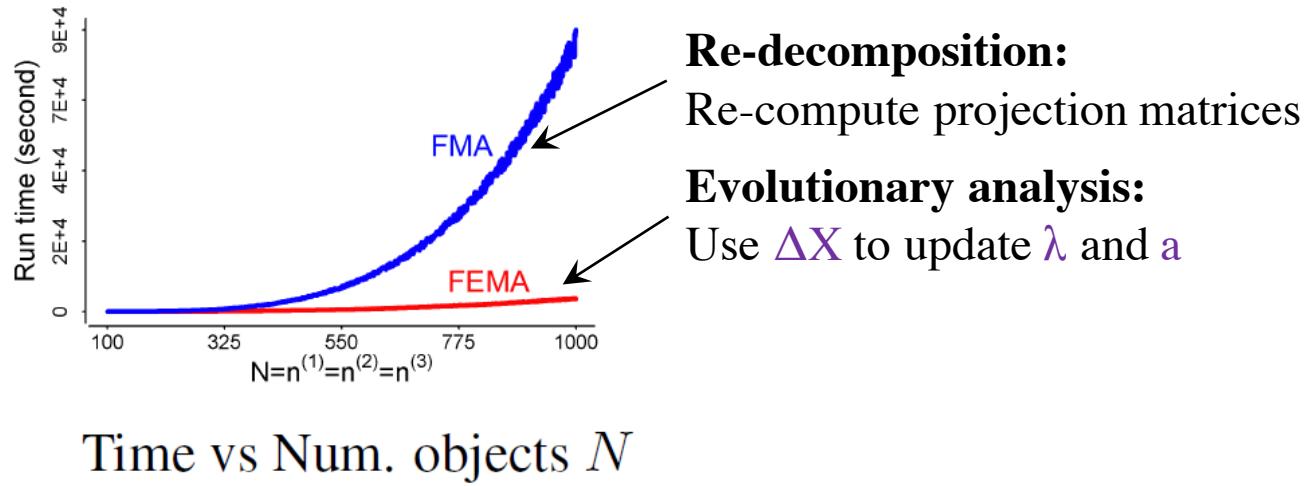
$$|\Delta\mathbf{a}_i^{(m)}| \leq 2\|\Delta\mathbf{X}^{(m)}\|_2 \sum_{j \neq i} \frac{(\lambda_{\mathbf{X}^{(m)\top} \mathbf{X}^{(m)}}^{\max})^{\frac{1}{2}}}{|\lambda_i^{(m)} - \lambda_j^{(m)}|}$$

projection matrix

Results: FEMA > EMA > EA

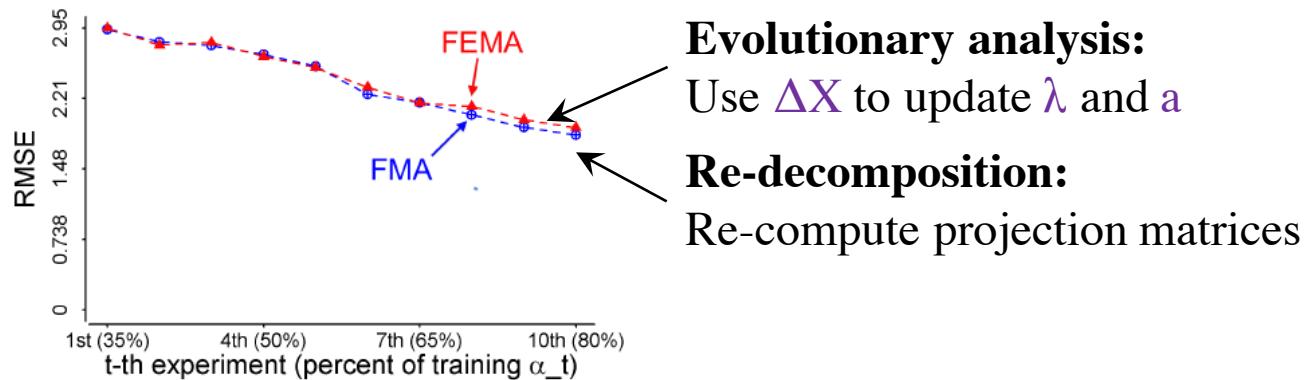
	Microsoft Academic Search		Tencent Weibo mentions “@”	
	MAE	RMSE	MAE	RMSE
FEMA 	0.735	0.944	0.894	1.312
EMA 	0.794	1.130	0.932	1.556
EA 	0.979	1.364	1.120	1.873
Precision vs Recall				

Results: Efficiency



Re-decomposition:
Re-compute projection matrices

Evolutionary analysis:
Use ΔX to update λ and a



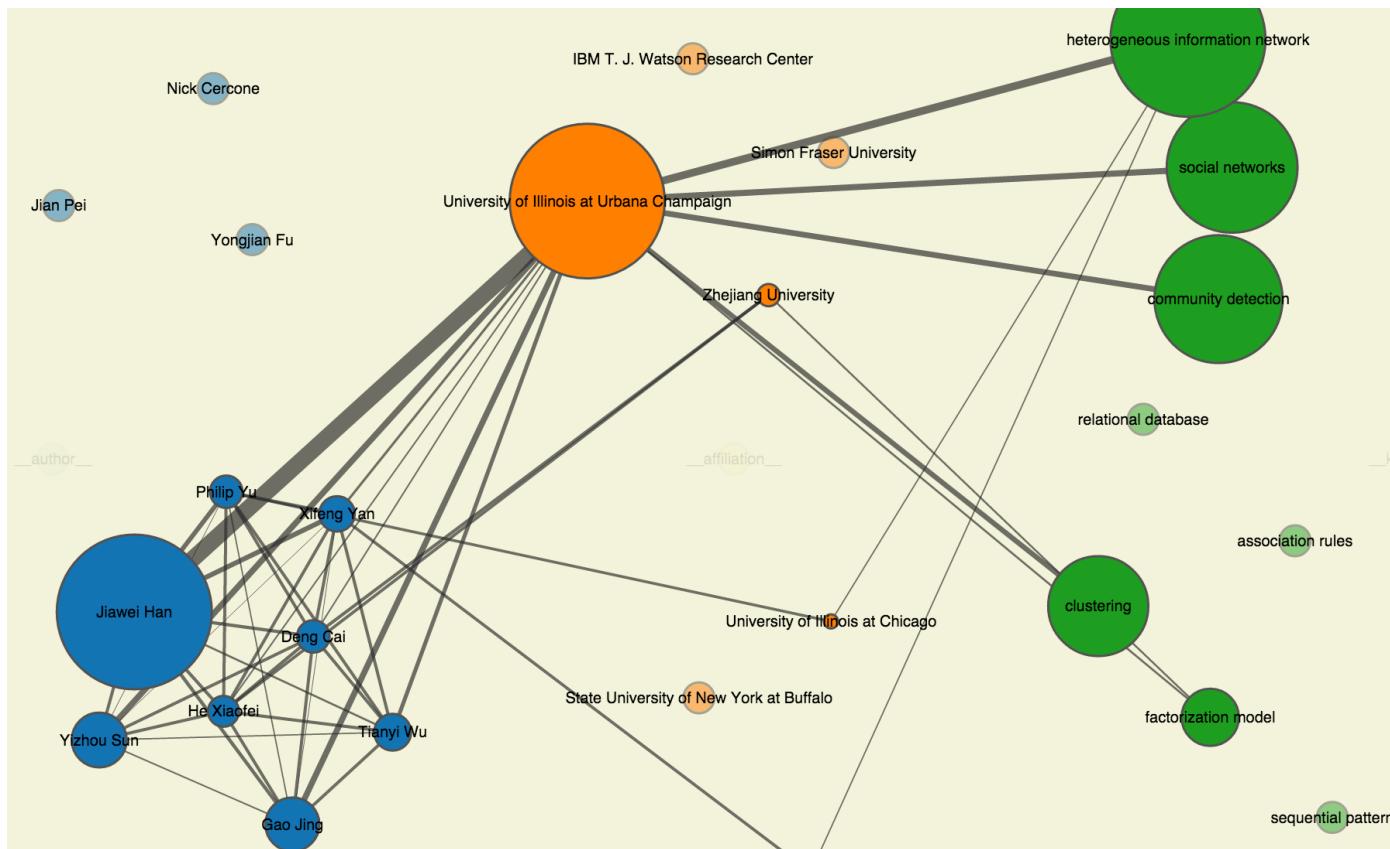
Evolutionary analysis:
Use ΔX to update λ and a

Re-decomposition:
Re-compute projection matrices

The loss is small.

Demo: Author@Affiliation#Keyword

<http://www.meng-jiang.com/demos/fema/mas/>



Observation: Multiple Domains



Osmar Zaiane

20 hrs · Twitter · 

#DataScientists need ability to tell the story about #data and convey #business value <https://t.co/VNN2rXaLuV> #BigData #datascience #dataviz

 Like  Comment  Share

The Globe and Mail shared Globe Politics's video.
19 hrs · 

Watch highlights from Stephen Harper's concession speech





Philip Bohannon shared a link.
5 hrs · 



British Library offers over 1 million free vintage images for download

9#
Closed Group

Joined  Share  ...

Discussion Members Events Photos Files Search this group

Write Post Add Photo / Video Ask Question Add File

Write something...

RECENT ACTIVITY

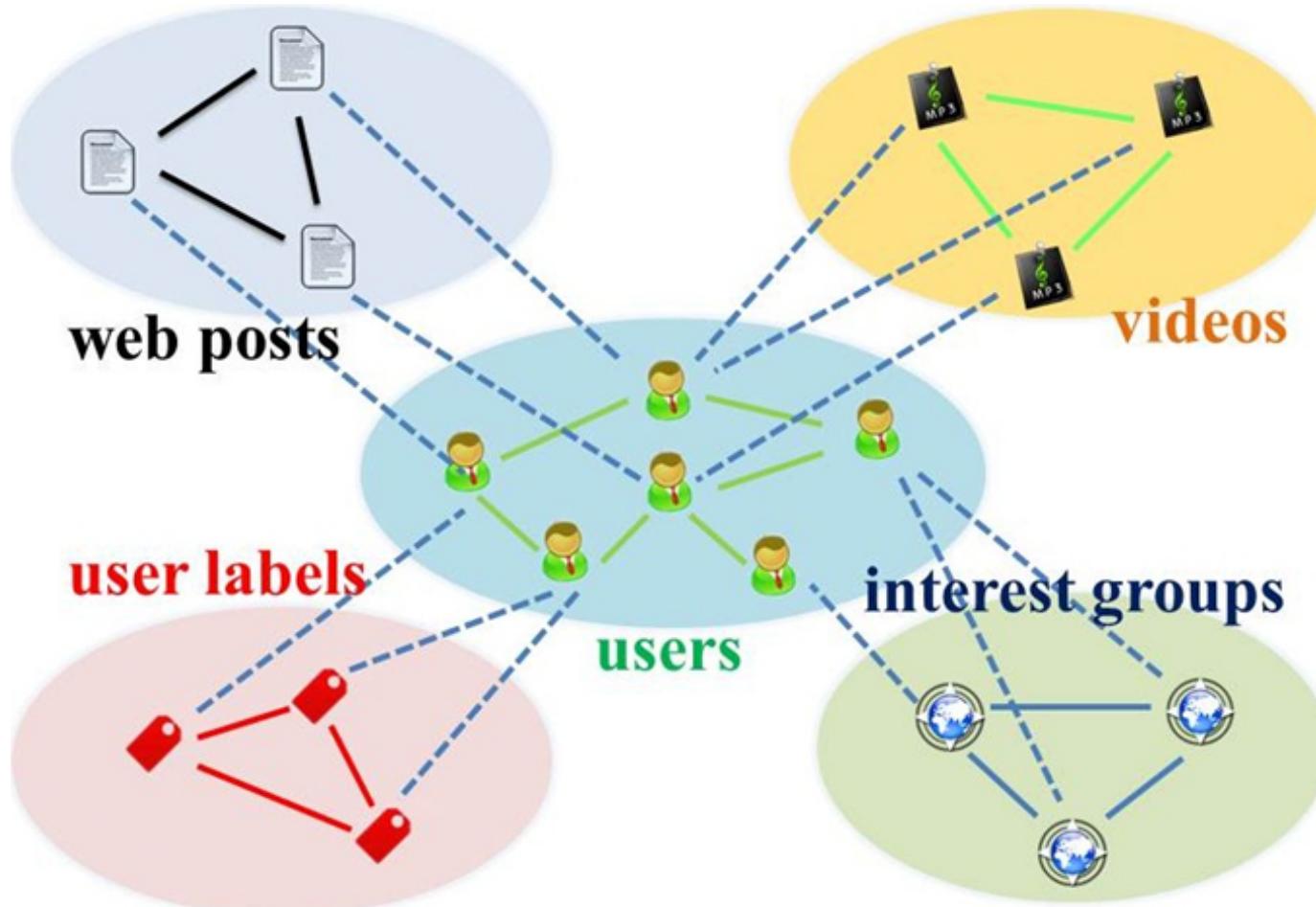
MEMBERS 1,049 Members (4 new)
+ Add People to Group



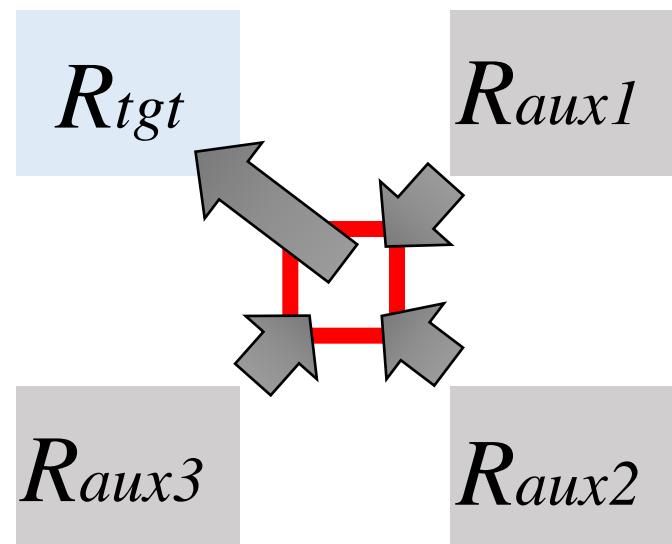
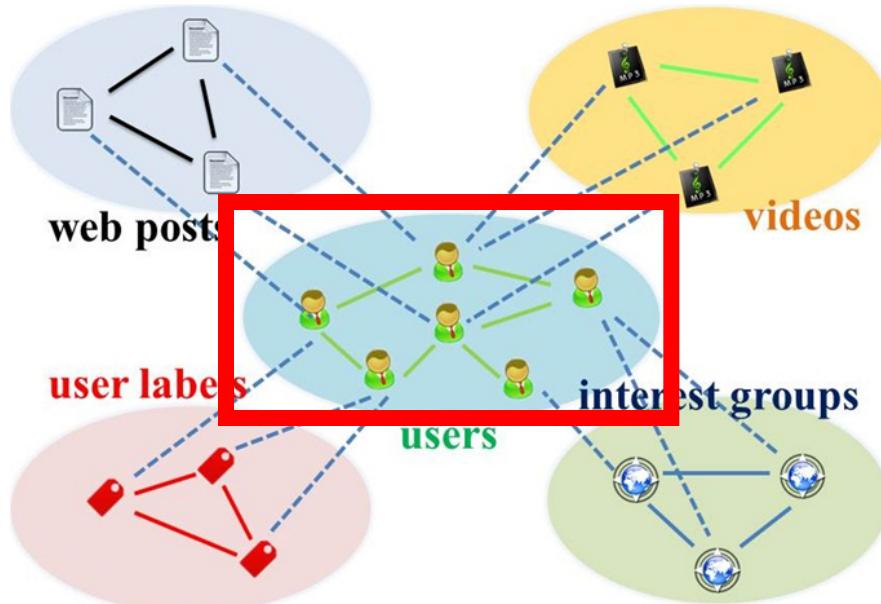
Invite by Email

Religious Views	Christian
Interests	Basketball, writing, spending time w/ kids
Favorite Music	Miles Davis, John Coltrane, Bob Dylan, Stevie Wonder, Johann Sebastian Bach (cello suites), and The Fugees
Favorite Movies	Casablanca, Godfather I & II, Lawrence of Arabia and One Flew Over the Cuckoo's Nest
Favorite TV Shows	Sportscenter
Favorite Quotations	"The Arc of the moral universe is long, but it bends towards justice." (MLK)

Representation: Star-Structured Graph



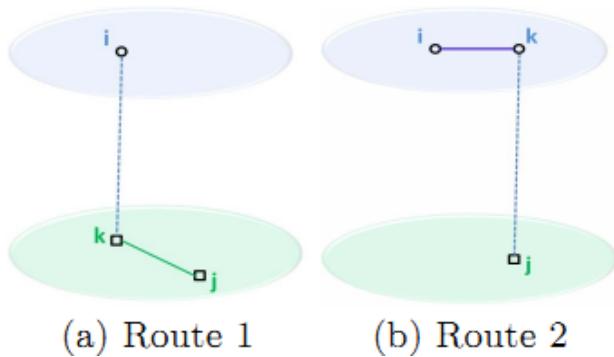
Representation: Social Bridge



Bridge: Tie strength

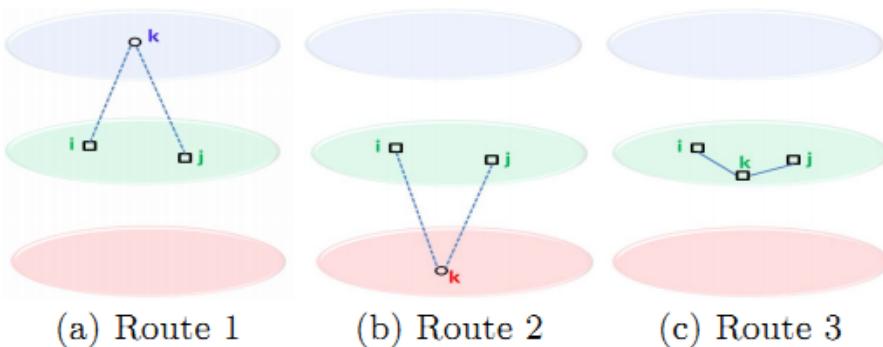
Algorithm: Hybrid Random Walk

□ Updating cross-domain links



$$\begin{aligned}
 p_{ij}^{(\mathcal{UP})+} &= \delta \sum_{u_k \in \mathcal{U}} r_{ik}^{(\mathcal{U})} p_{kj}^{(\mathcal{UP})+} + (1 - \delta) \sum_{p_k \in \mathcal{P}} p_{ik}^{(\mathcal{UP})+} r_{kj}^{(\mathcal{P})} \\
 p_{ij}^{(\mathcal{UP})-} &= \delta \sum_{u_k \in \mathcal{U}} r_{ik}^{(\mathcal{U})} p_{kj}^{(\mathcal{UP})-} + (1 - \delta) \sum_{p_k \in \mathcal{P}} p_{ik}^{(\mathcal{UP})-} r_{kj}^{(\mathcal{P})} \\
 p_{ij}^{(\mathcal{UT})+} &= \eta \sum_{u_k \in \mathcal{U}} r_{ik}^{(\mathcal{U})} p_{kj}^{(\mathcal{UT})+} + (1 - \eta) \sum_{t_k \in \mathcal{T}} p_{ik}^{(\mathcal{UT})+} r_{kj}^{(\mathcal{T})} \\
 \mathbf{P}^{(\mathcal{UP})+}(t+1) &= \delta \mathbf{R}^{(\mathcal{U})}(t) \mathbf{P}^{(\mathcal{UP})+}(t) + (1 - \delta) \mathbf{P}^{(\mathcal{UP})+}(t) \mathbf{R}^{(\mathcal{P})} \\
 \mathbf{P}^{(\mathcal{UP})-}(t+1) &= \delta \mathbf{R}^{(\mathcal{U})}(t) \mathbf{P}^{(\mathcal{UP})-}(t) + (1 - \delta) \mathbf{P}^{(\mathcal{UP})-}(t) \mathbf{R}^{(\mathcal{P})} \\
 \mathbf{P}^{(\mathcal{UT})+}(t+1) &= \eta \mathbf{R}^{(\mathcal{U})}(t) \mathbf{P}^{(\mathcal{UT})+}(t) + (1 - \eta) \mathbf{P}^{(\mathcal{UT})+}(t) \mathbf{R}^{(\mathcal{T})}
 \end{aligned}$$

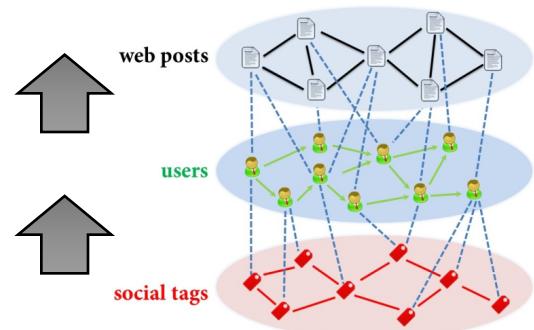
□ Updating within-domain links



$$\begin{aligned}
 r_{ij}^{(\mathcal{U})} &= \tau^{(\mathcal{P})} (\mu \sum_{p_k \in \mathcal{P}} p_{ik}^{(\mathcal{UP})+} p_{jk}^{(\mathcal{UP})+} + (1 - \mu) \sum_{p_k \in \mathcal{P}} p_{ik}^{(\mathcal{UP})-} p_{jk}^{(\mathcal{UP})-}) \\
 &\quad + \tau^{(\mathcal{T})} \sum_{t_k \in \mathcal{T}} p_{ik}^{(\mathcal{UT})+} p_{jk}^{(\mathcal{UT})+} + \tau^{(\mathcal{U})} \sum_{u_k \in \mathcal{U}} r_{ik}^{(\mathcal{U})} r_{kj}^{(\mathcal{U})}
 \end{aligned} \tag{12}$$

$$\begin{aligned}
 \mathbf{R}^{(\mathcal{U})}(t+1) &= \\
 &\quad \tau^{(\mathcal{P})} (\mu \mathbf{P}^{(\mathcal{UP})+}(t) \mathbf{P}^{(\mathcal{UP})+}(t)^T + (1 - \mu) \mathbf{P}^{(\mathcal{UP})-}(t) \mathbf{P}^{(\mathcal{UP})-}(t)^T) \\
 &\quad + \tau^{(\mathcal{T})} \mathbf{P}^{(\mathcal{UT})+}(t) \mathbf{P}^{(\mathcal{UT})+}(t)^T + \tau^{(\mathcal{U})} \mathbf{R}^{(\mathcal{U})}(t) \mathbf{R}^{(\mathcal{U})}(t)^T
 \end{aligned} \tag{13}$$

Results



Comparing with Random Walk with Restarts Models

Algorithm	MAE	Precision	Recall	F1	Kendall's $\hat{\tau}$
HRW	0.227±1.5e-3	0.711±1.3e-3	0.921±1.4e-3	0.802±1.1e-3	0.792±2.5e-3
BRW- R_U -P (TrustWalker)	0.276±1.1e-3	0.657±7.6e-4	0.935±9.8e-4	0.772±7.6e-4	0.774±1.6e-3
BRW- R_U	0.282±5.3e-3	0.655±4.0e-3	0.921±1.2e-2	0.765±7.7e-3	0.725±2.8e-3
BRW- W_U -P	0.292±1.1e-3	0.666±7.0e-4	0.900±5.2e-4	0.765±6.6e-4	0.725±8.5e-4
BRW- W_U (ItemRank)	0.318±1.4e-3	0.671±1.5e-3	0.713±2.4e-3	0.691±1.2e-3	0.661±2.2e-3
BRW-P	0.438±2.6e-4	0.571±3.4e-4	0.499±4.2e-4	0.532±3.2e-4	0.606±2.3e-4

Comparing with Social Recommendation Baselines

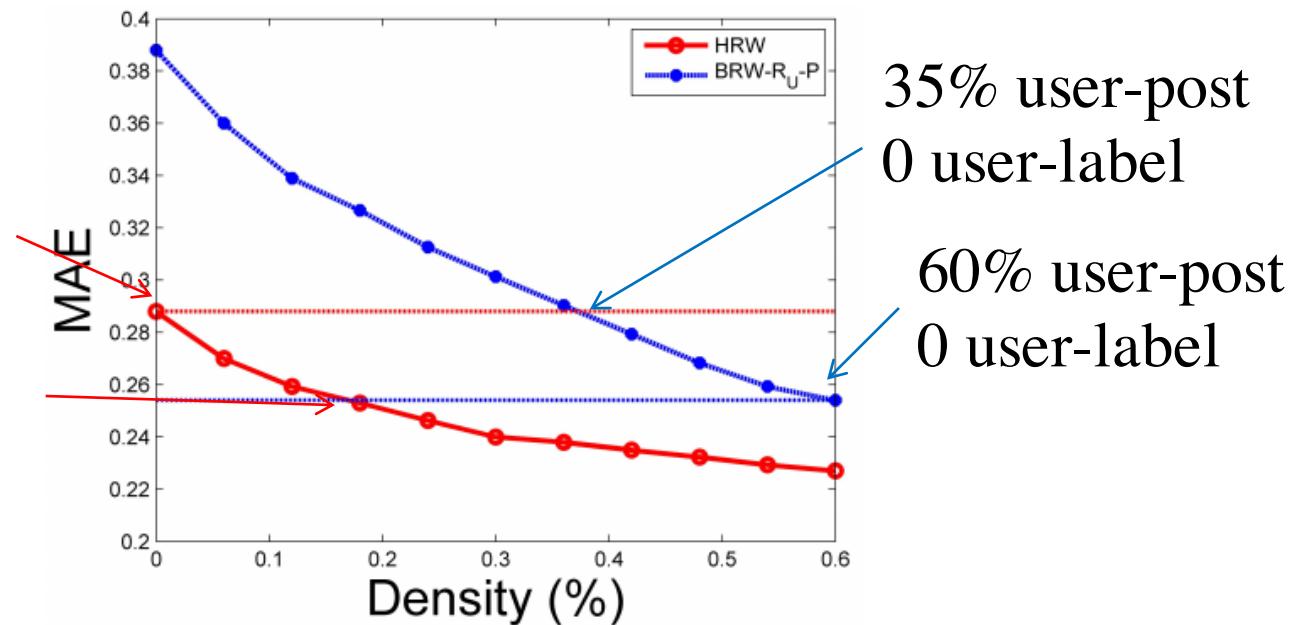
Algorithm	MAE	Precision	Recall	F1	Kendall's $\hat{\tau}$
HRW	0.227±1.5e-3	0.711±1.3e-3	0.921±1.4e-3	0.802±1.1e-3	0.792±2.5e-3
BRW- R_U -P (TrustWalker) [10]	0.276±1.1e-3	0.657±7.6e-4	0.935±9.8e-4	0.772±7.6e-4	0.774±1.6e-3
BRW- W_U (ItemRank) [8]	0.318±1.4e-3	0.671±1.5e-3	0.713±2.4e-3	0.691±1.2e-3	0.661±2.2e-3
MCF [5]	0.352±2.3e-4	0.592±1.8e-3	0.951±6.0e-4	0.730±1.3e-3	0.582±4.3e-4
CF [22]	0.506±3.4e-4	0.552±1.5e-3	0.589±7.2e-4	0.570±1.0e-3	0.540±5.2e-4

Results: Insight

- ❑ Knowledge transfer from auxiliary domains improves cold-start users' behavior prediction
 - ❑ Using aux. (label) data, saving **60-70%** tgt. (post) data

0 user-post
100% user-label

18% user-post
100% user-label



Observation: Multiple Platforms



Observation: Cross-Platform

Add Facebook Login to Your App or Website

Facebook Login for Apps is a secure, fast and convenient way for people to log into your app or website.



iOS



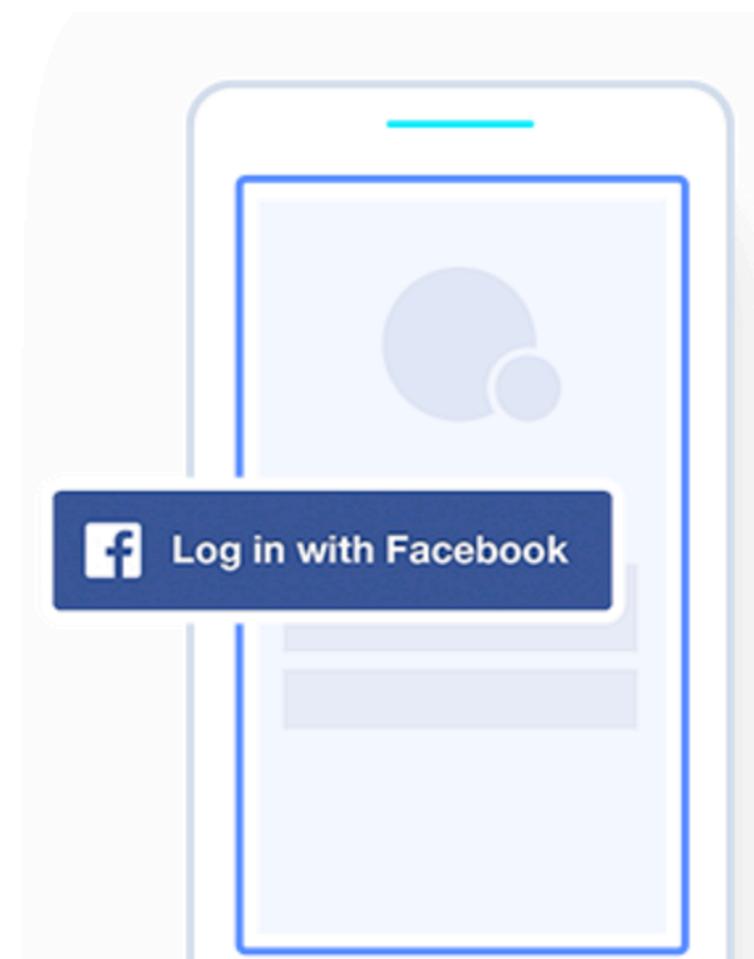
Android



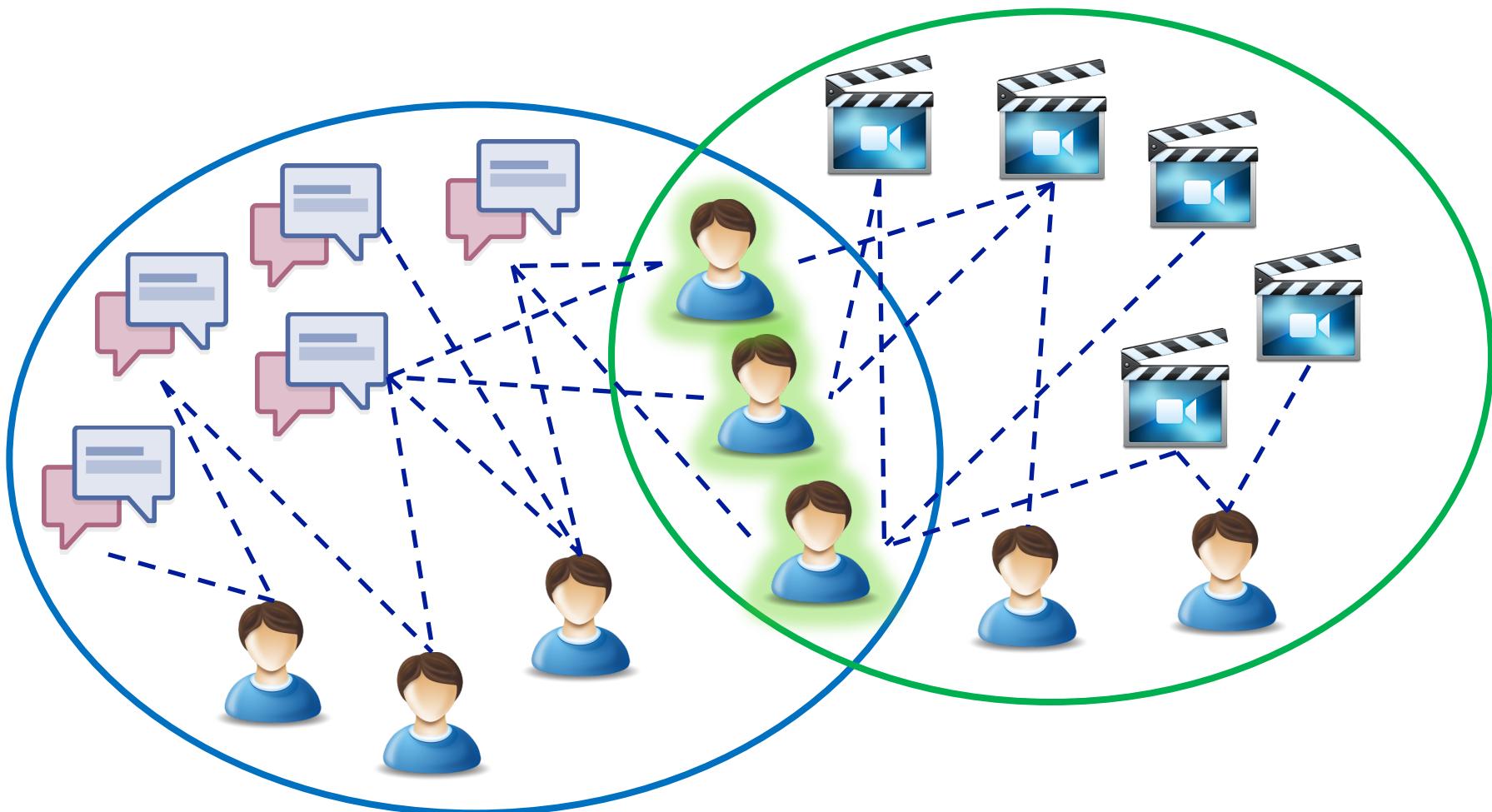
Websites or mobile websites



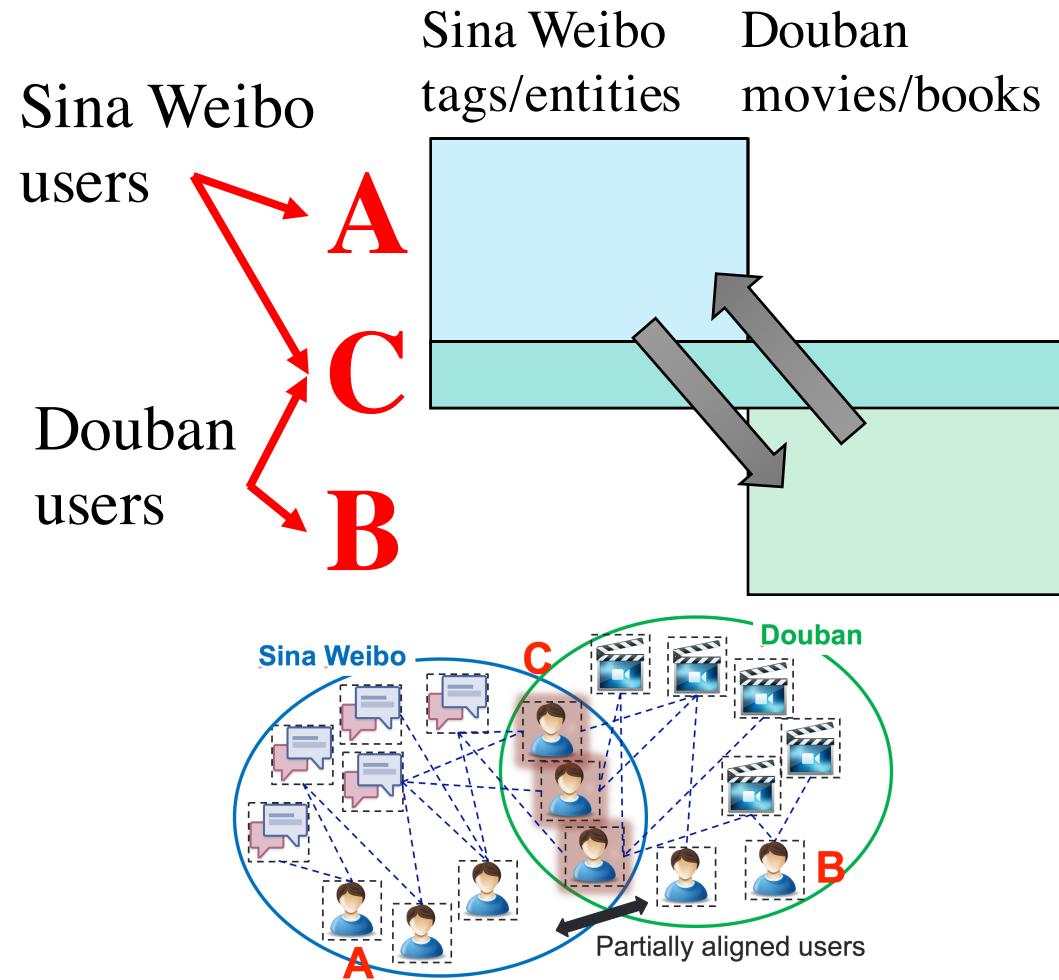
More platforms



Observation: Partially Overlapped Crowds



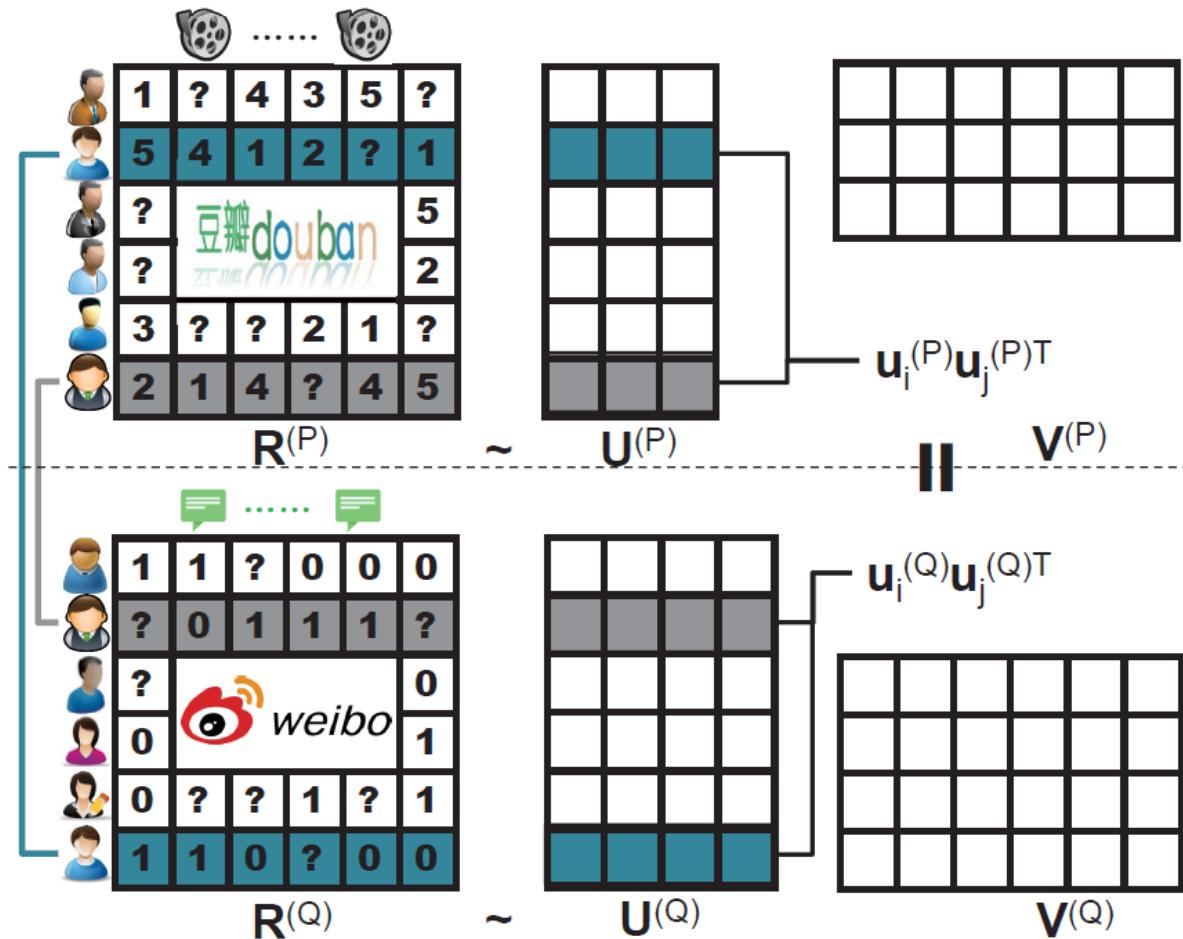
Representation: When NO Transfer



User set	Weibo tweet entity to Douban movie	RMSE	MAP
A	Auxiliary platform data!		
C	0.779	0.805	
B	1.439		0.640

User set	Douban book to Weibo social tag	RMSE	MAP
A		0.429	0.464
C	0.267	0.666	
B	Auxiliary platform data!		

Algorithm: XPTTrans



Algorithm: XPTTrans

□ Input

- Tgt./Aux. platform P/Q;
- Behavior data R(P)/R(Q);
- Observation W(P)/W(Q);
- Overlapping indicator W(P,Q),

□ Output

- User latent representation U(P)/U(Q);
- Item latent representation V(P)/V(Q);
- Missing values in R(P)

□ Objective function

Target platform Auxiliary platform

$$\mathcal{J} = \sum_{i,j} W_{i,j}^{(P)} \left(R_{i,j}^{(P)} - \sum_r U_{i,r}^{(P)} V_{r,j}^{(P)} \right)^2 + \lambda \sum_{i,j} W_{i,j}^{(Q)} \left(R_{i,j}^{(Q)} - \sum_r U_{i,r}^{(Q)} V_{r,j}^{(Q)} \right)^2 + \mu \sum_{i_1,j_1,i_2,j_2} W_{i_1,j_1}^{(P,Q)} W_{i_2,j_2}^{(P,Q)} \left(A_{i_1,i_2}^{(P)} - A_{j_1,j_2}^{(Q)} \right)^2$$

Overlapping user similarity
(Pair-wise regularization)

Results: Leveraging Auxiliary Platform Data

NO Transfer

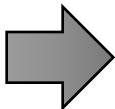
User set	Weibo tweet entity to Douban movie	
	RMSE	MAP
A	Auxiliary platform data!	
C	0.779	0.805
B	1.439	0.640

User set	Douban book to Weibo social tag	
	RMSE	MAP
A	0.429	0.464
C	0.267	0.666
B	Auxiliary platform data!	

Transfer via the Same Latent Space

User set	Weibo tweet entity to Douban movie	
	RMSE	MAP
A		
C	0.757	0.811
B	1.164 (-19%)	0.702 (+9.7%)

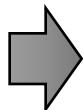
User set	Douban book to Weibo social tag	
	RMSE	MAP
A	0.411 (-4.2%)	0.487 (+5.0%)
C	0.256	0.681
B		



Results: Leveraging Different Latent Spaces

Transfer via the Same Latent Space

User set	Weibo tweet entity to Douban movie	
	RMSE	MAP
A		
C	0.757	0.811
B	1.164	0.702
User set	Douban book to Weibo social tag	
	RMSE	MAP
A	0.411	0.487
C	0.256	0.681
B		



Transfer via Different Latent Spaces

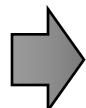
User set	Weibo tweet entity to Douban movie	
	RMSE	MAP
A		
C	0.715	0.821
B	0.722 (-38%)	0.820 (+17%)
User set	Douban book to Weibo social tag	
	RMSE	MAP
A	0.374 (-11 %)	0.533 (+12 %)
C	0.236	0.705
B		

Results: Where Amazing Happens

NO Transfer

User set	Weibo tweet entity to Douban movie	
	RMSE	MAP
A	Auxiliary platform data!	
C	0.779	0.805
B	1.439	0.640

User set	Douban book to Weibo social tag	
	RMSE	MAP
A	0.429	0.464
C	0.267	0.666
B	Auxiliary platform data!	

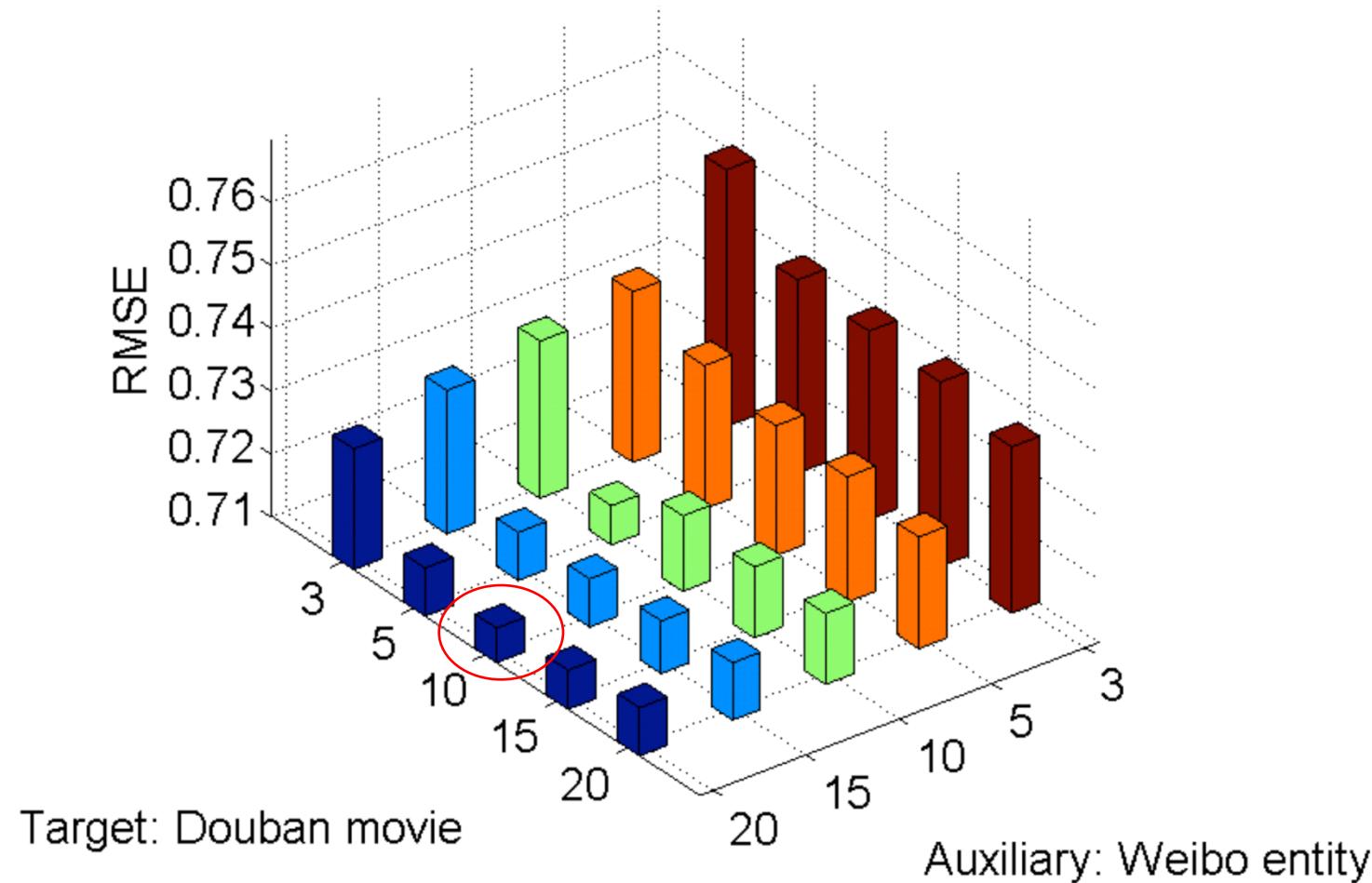


Transfer via Different Latent Spaces

User set	Weibo tweet entity to Douban movie	
	RMSE	MAP
A		
C	0.715	0.821
B	0.722	0.820

User set	Douban book to Weibo social tag	
	RMSE	MAP
A	0.374	0.533
C	0.236	0.705
B		

Results: Different Sizes of Latent Spaces





Summary

- ❑ Like, Reply, Share, Retweet, Favorite, Comment ...
- ❑ Memory based social recommenders
 - ❑ TidalTrust, MoleTrust, TrustWalker
- ❑ Model based social recommenders
 - ❑ SoRec, “Social Trust” Ensemble, SoReg
- ❑ **Observations, Representations, Models**
 - ❑ **ContextMF**: Social contexts (preference & influence)
 - ❑ **FEMA**: Spatiotemporal contexts (multidimensional)
 - ❑ **HybridRW**: Cross-domain behavior modeling
 - ❑ **XPTrans**: Cross-platform behavior modeling



I. Mining behavior networks with social and spatiotemporal contexts

I.2. Suspicious behavior detection



Ill-gotten Facebook Likes

25,000 Facebook Likes	50,000 Facebook Likes	100,000 Facebook Likes	200,000 Facebook Likes
\$265	\$525	\$1,000	\$1,750
Lifetime Replacement Warranty	Lifetime Replacement Warranty	Lifetime Replacement Warranty	Lifetime Replacement Warranty
Dedicated 24/7 Customer Service			
100% Risk Free, Try Us Today			
Order starts within 24 - 48 hours			
Order completed within 22 days	Order completed within 35 days	Order completed within 35 days	Order completed within 35 days

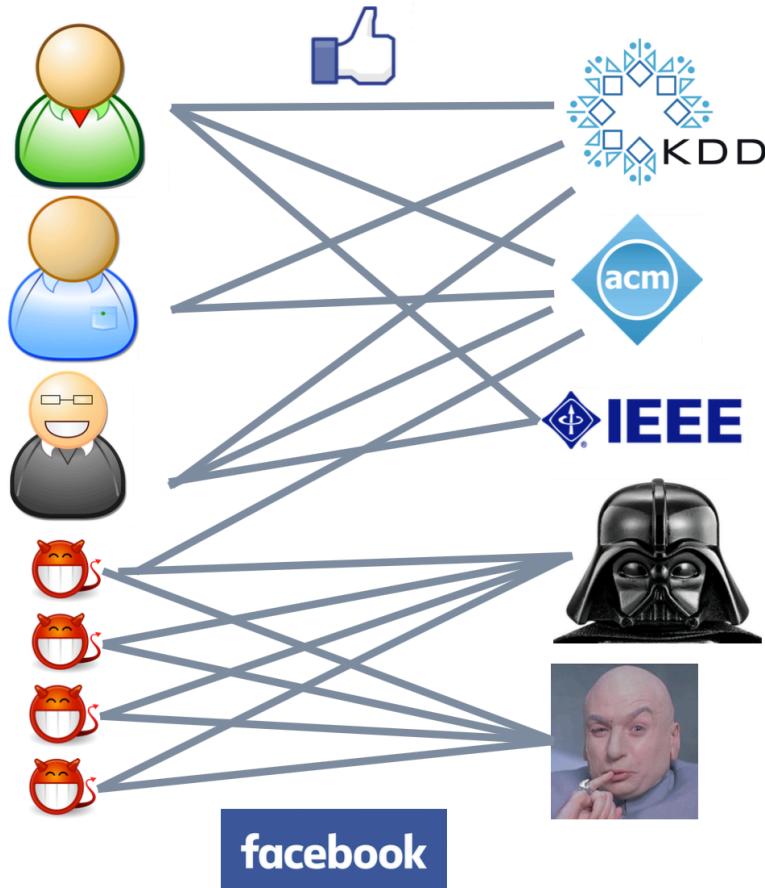
Suspicious Behavior Detection



Meng Jiang, Peng Cui and Christos Faloutsos.

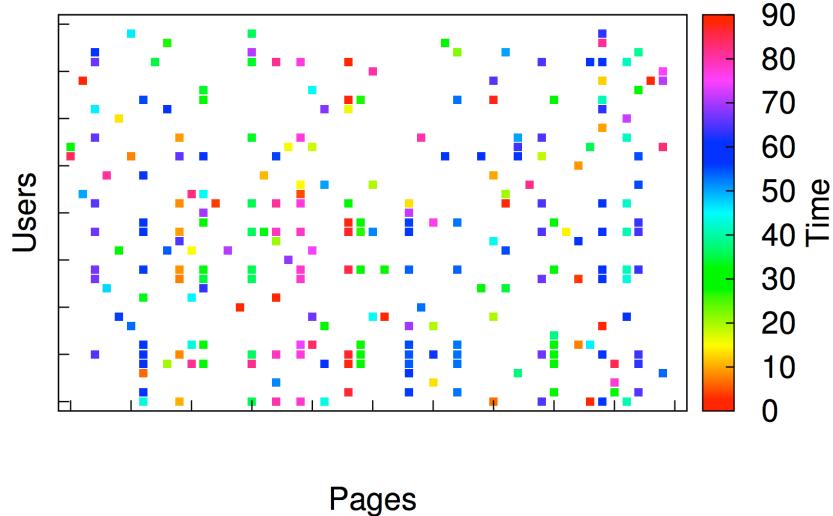
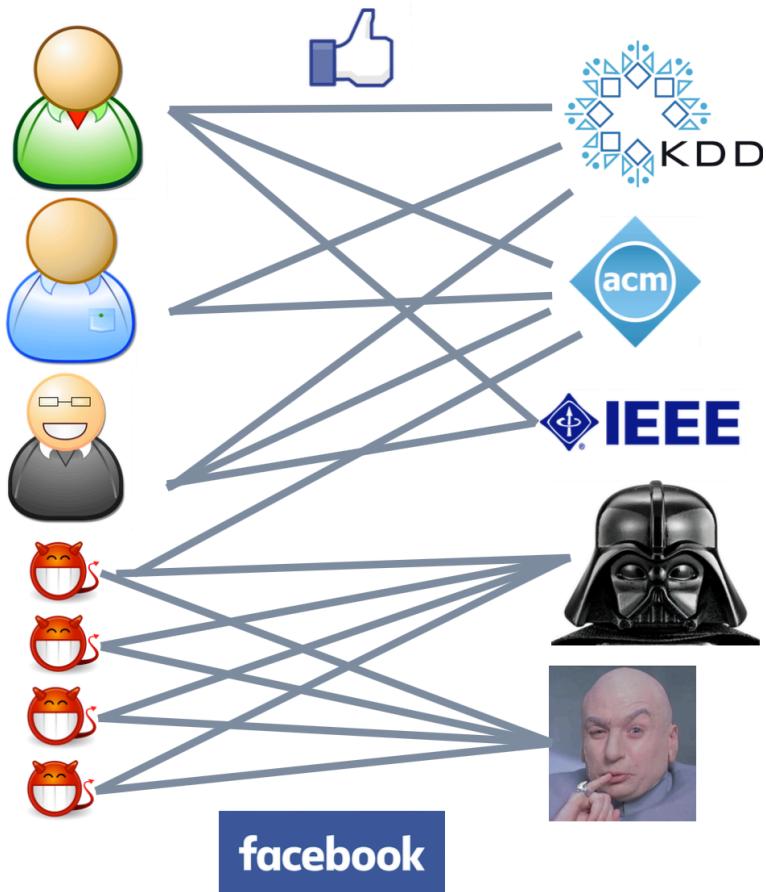
Suspicious Behavior Detection: Current Trends and Future Directions.
IEEE Intelligent Systems (ISSI), 2016.

Ill-gotten Facebook Likes

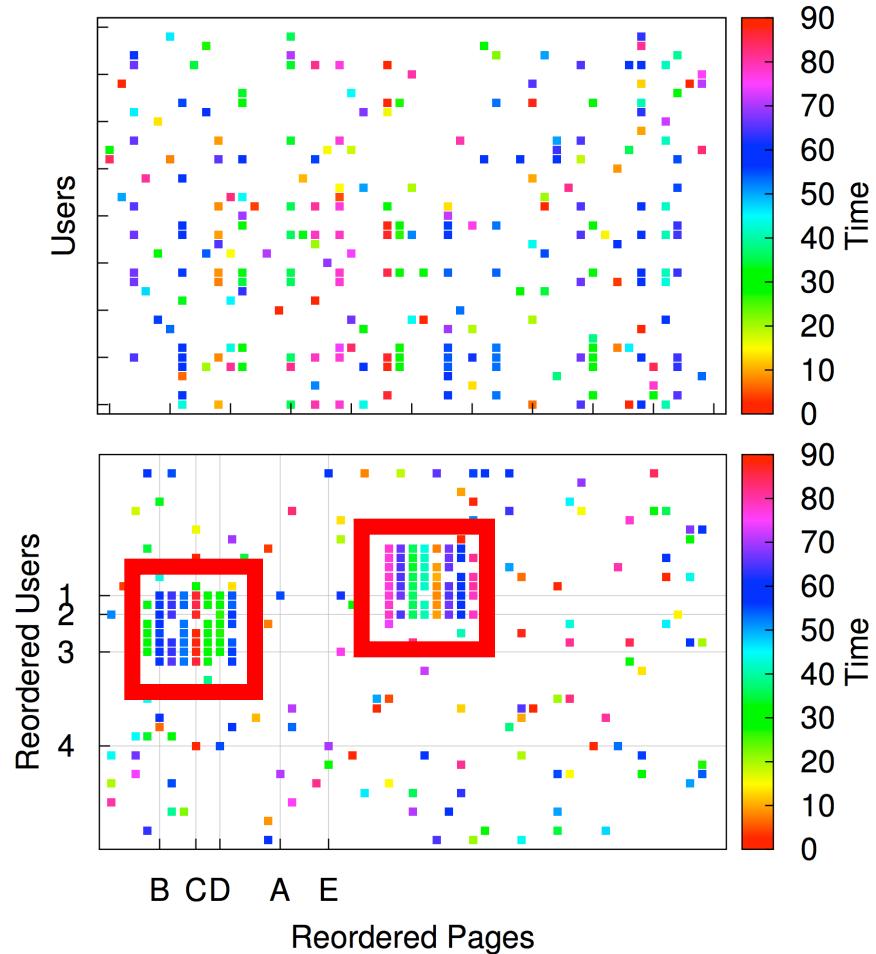
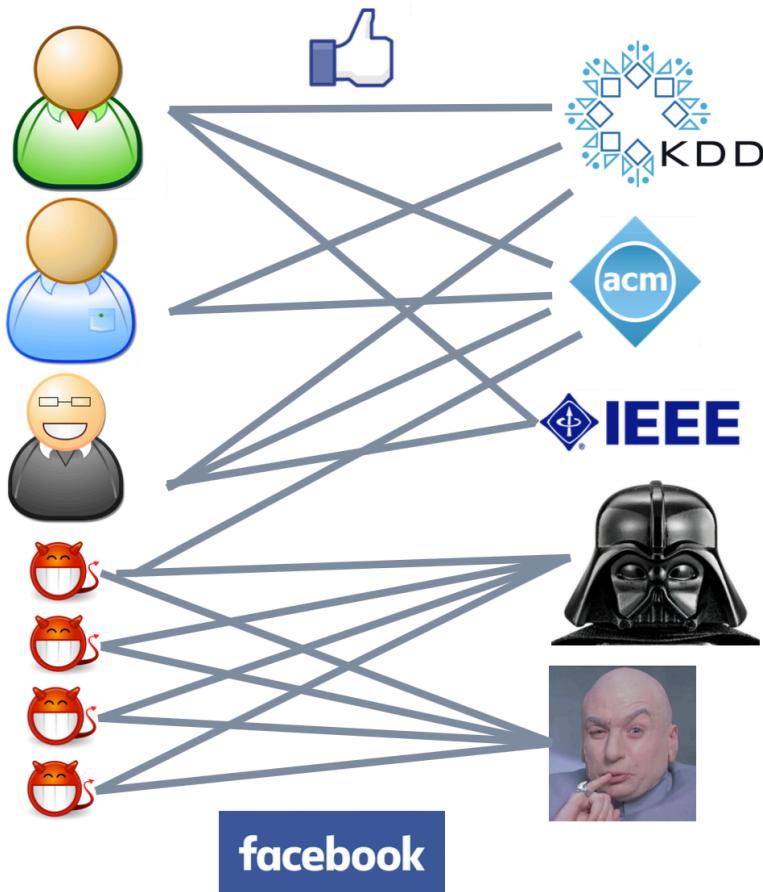


Beutel et al. **CopyCatch: Stopping Group Attacks by Spotting Lockstep Behavior in Social Networks.** WWW, 2013.

Observation: Graphical View



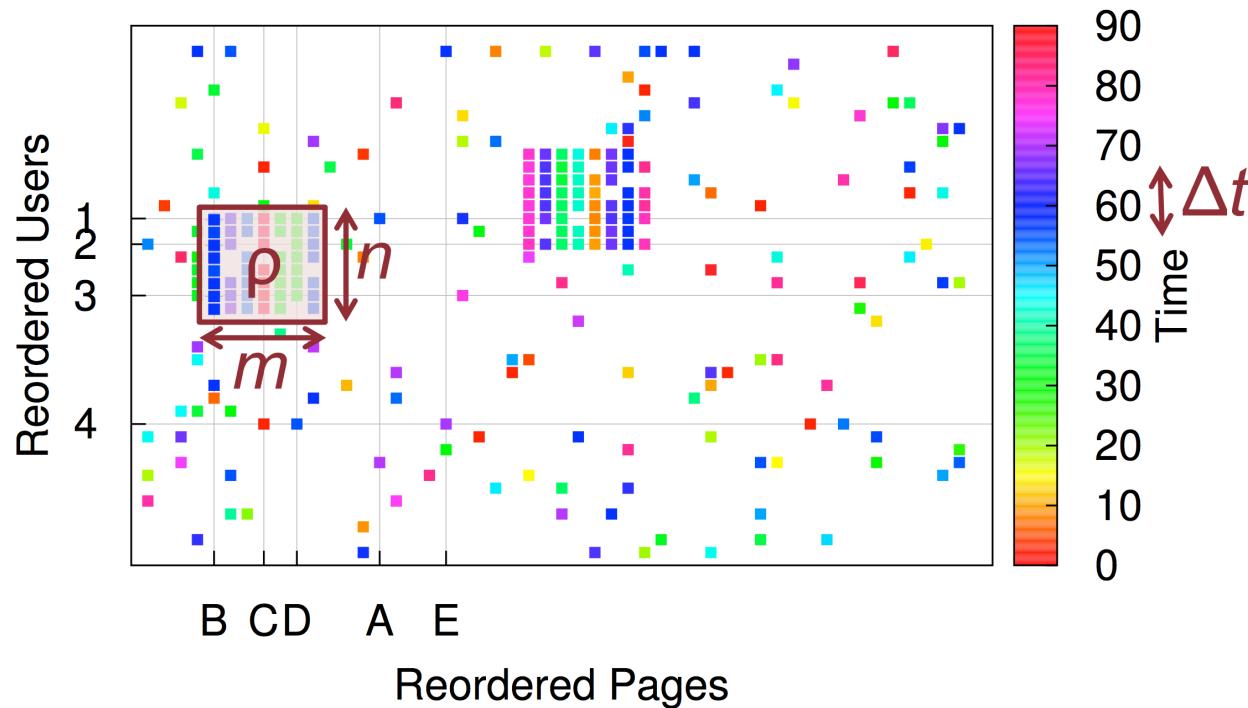
Observation: Reorder Matrix



Algorithm: Seed + Search

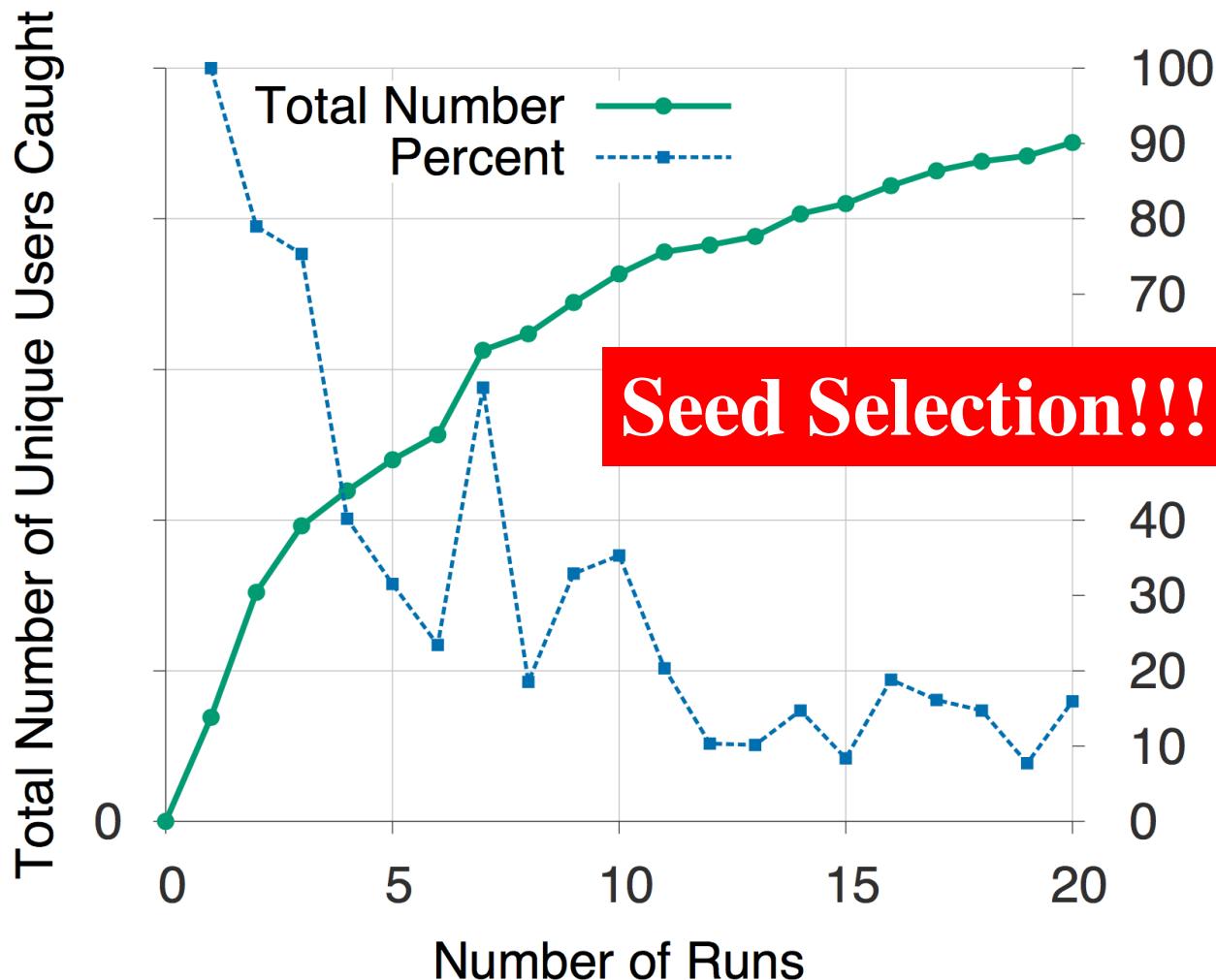
□ CopyCatch

□ “Near Bipartite Core”: n users, m Pages, Q , Δt





Experimental Result



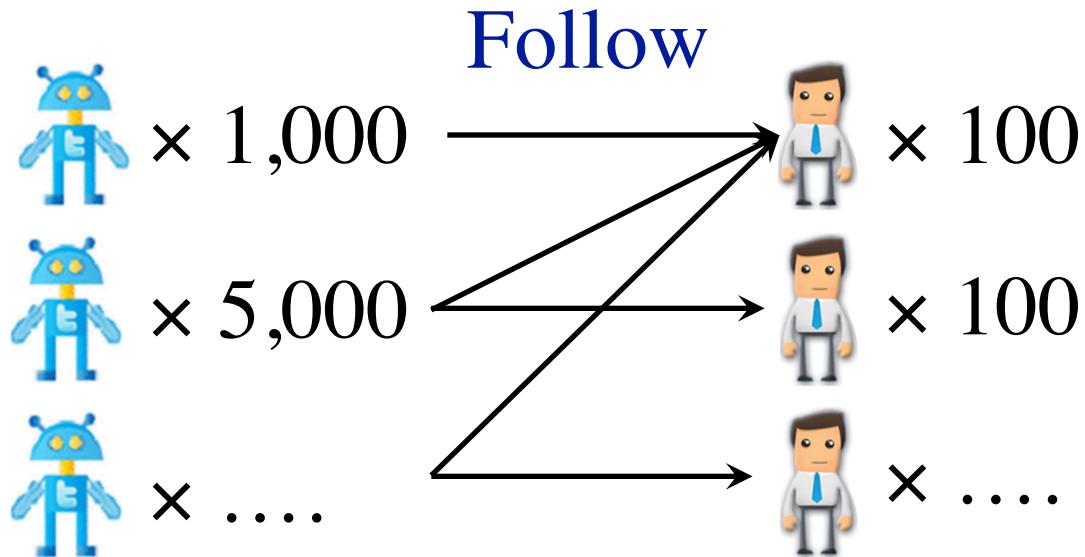
Serious Problem in Weibo



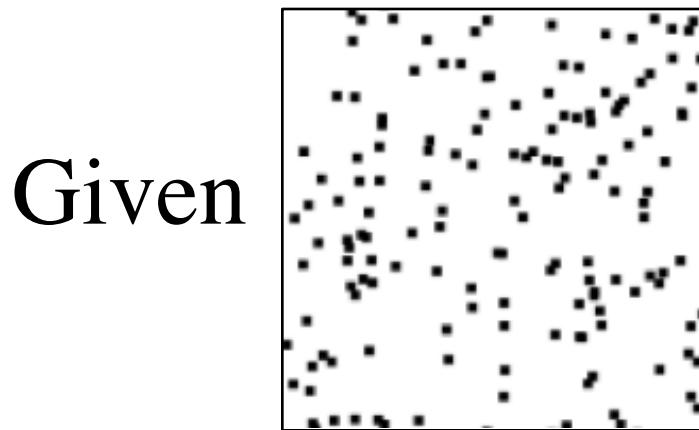
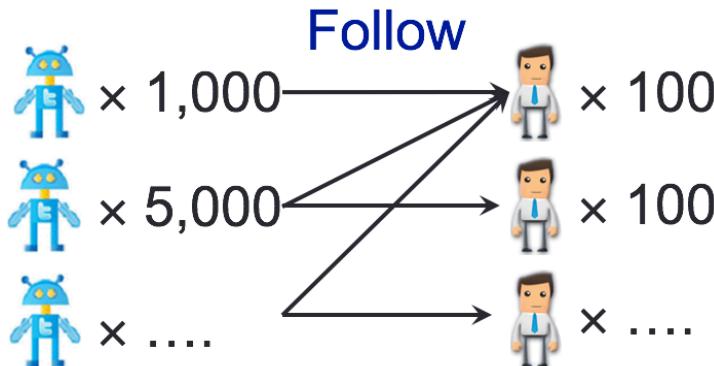
Experience-driven approaches:
features of #followees, #hashtags, #URLs...



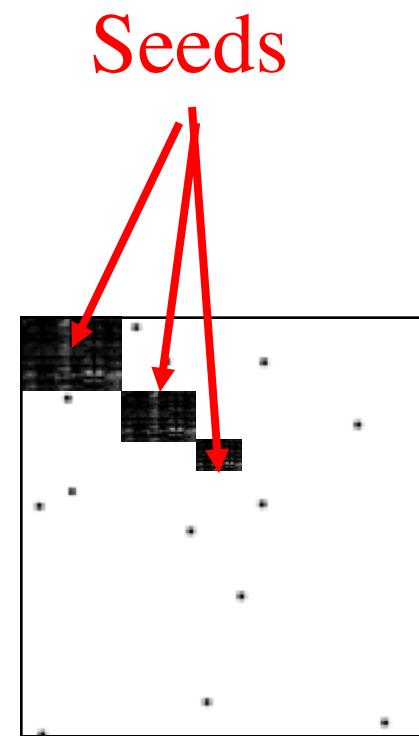
Zombie Followers



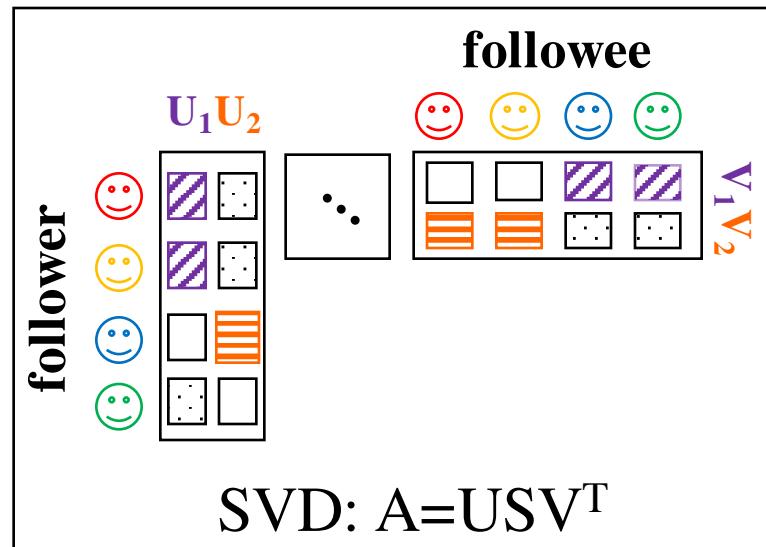
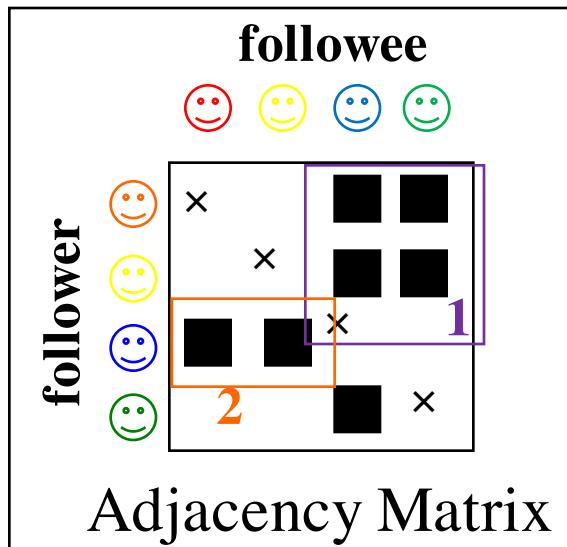
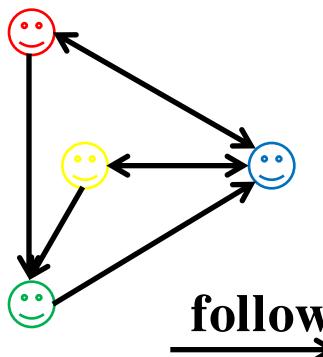
Observation: Reorder Matrix



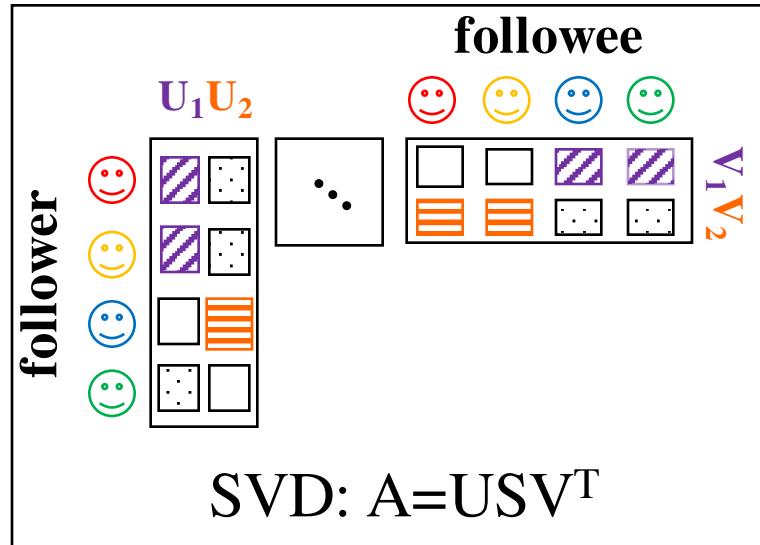
Reorder



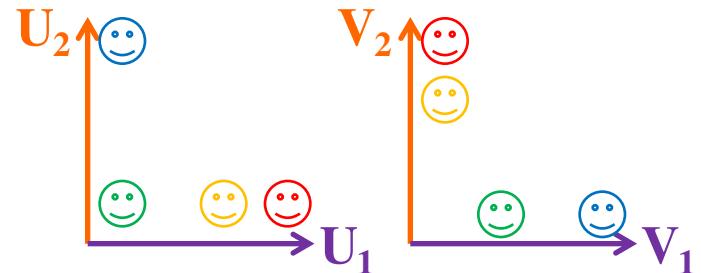
Representation: SVD Reminder



Representation: Spectral Subspace



Pairs of singular vectors:

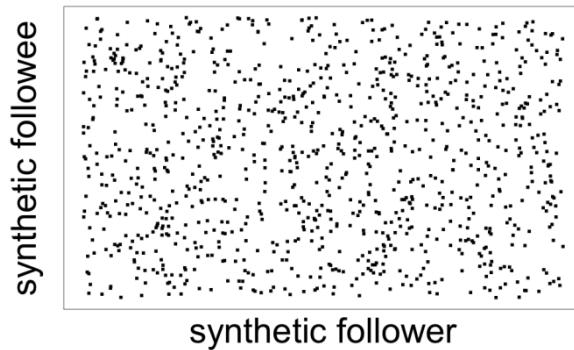


“Spectral Subspace Plot”

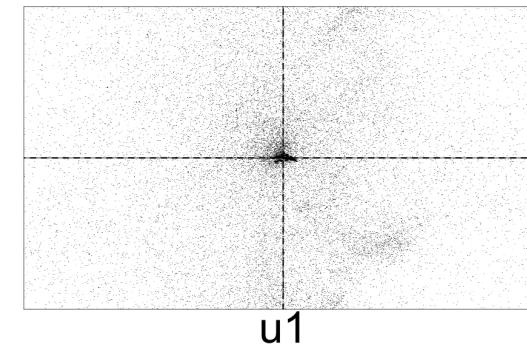
Spectral Subspace Plot: Case #0

- NO lockstep behavior: Scatter

Adjacency Matrix



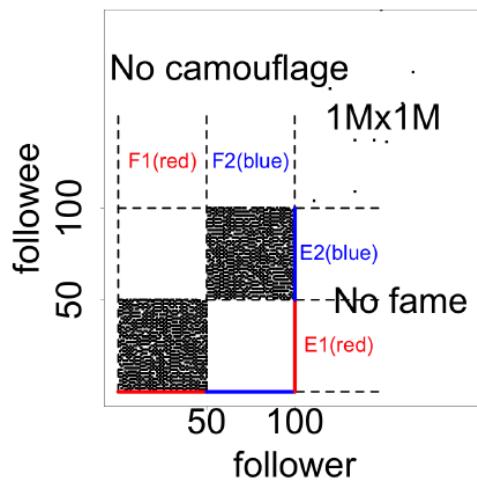
Spectral Subspace Plot



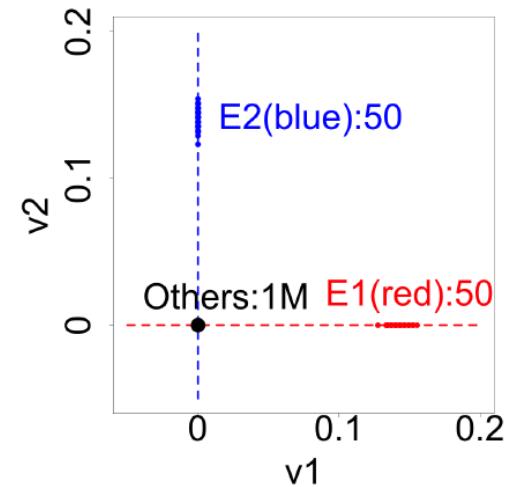
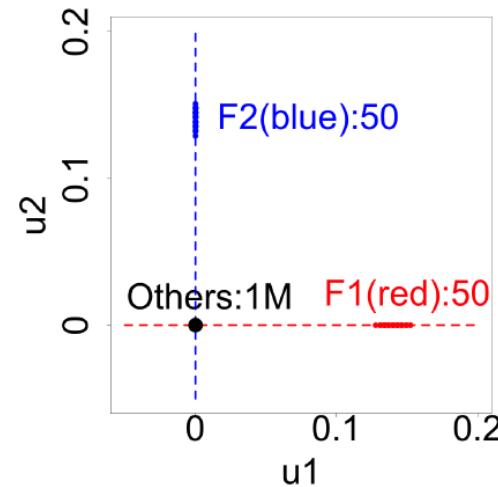
Spectral Subspace Plot: Case #1

- Non-overlapping lockstep: “Rays”

Adjacency Matrix



Spectral Subspace Plot

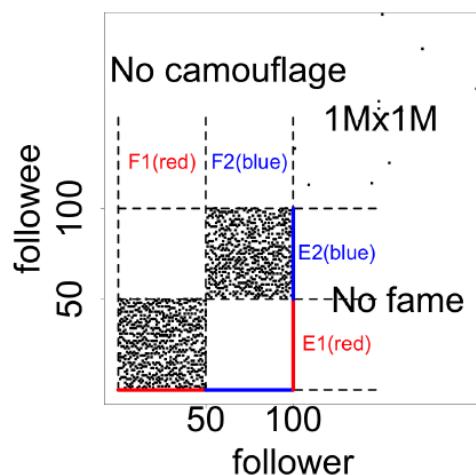


Rule 1 (short “rays”): two blocks, high density (90%), no “camouflage”, no “fame”

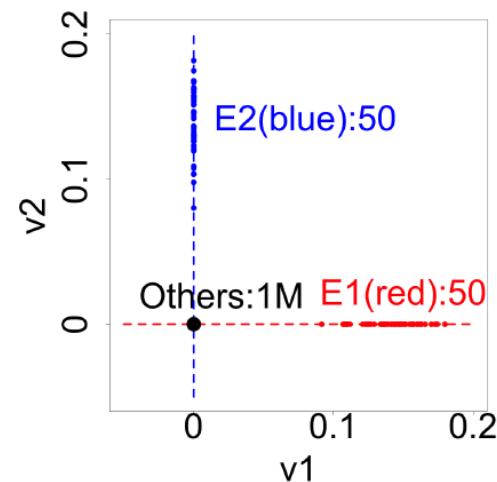
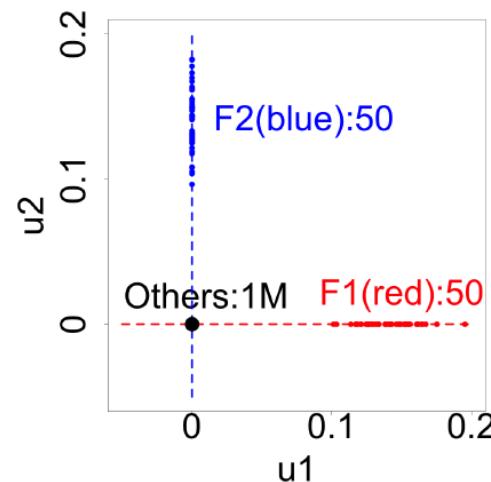
Spectral Subspace Plot: Case #2

- Non-overlapping: Low density, Elongation

Adjacency Matrix



Spectral Subspace Plot

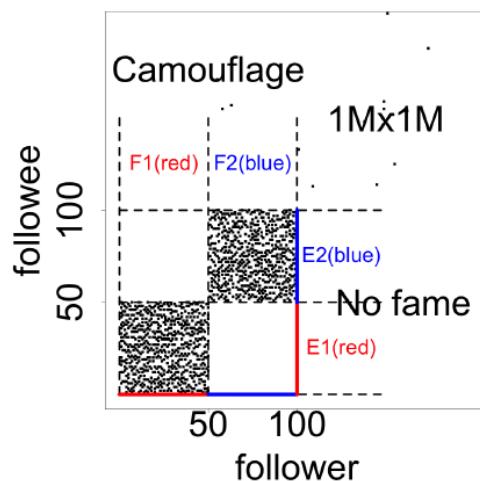


Rule 2 (long “rays”): two blocks, low density (50%), no “camouflage”, no “fame”

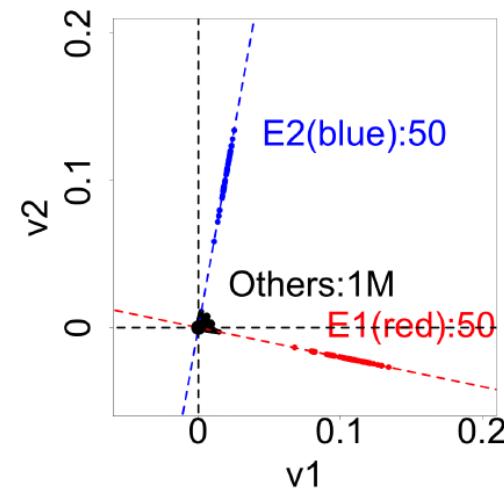
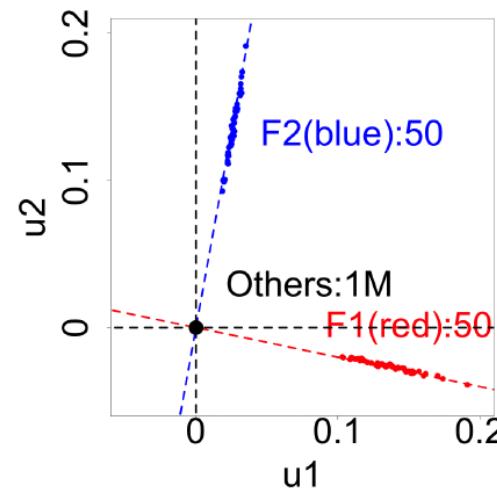
Spectral Subspace Plot: Case #3

- Non-overlapping: Camouflage/Fame, Tilting

Adjacency Matrix



Spectral Subspace Plot

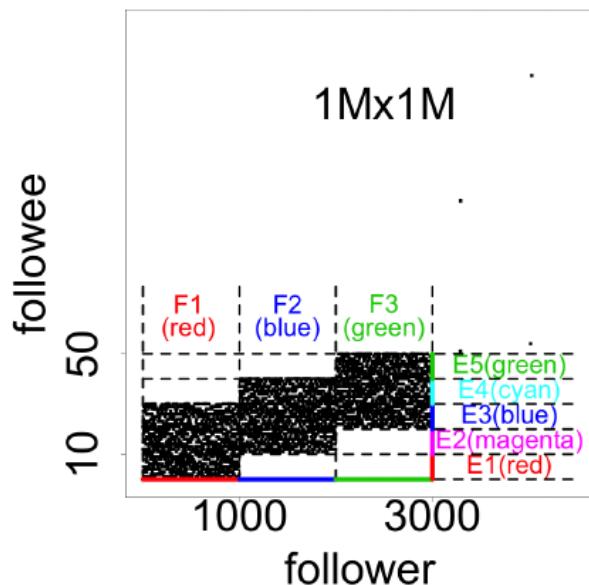


Rule 3 (tilting “rays”): two blocks, with “camouflage”, no “fame”

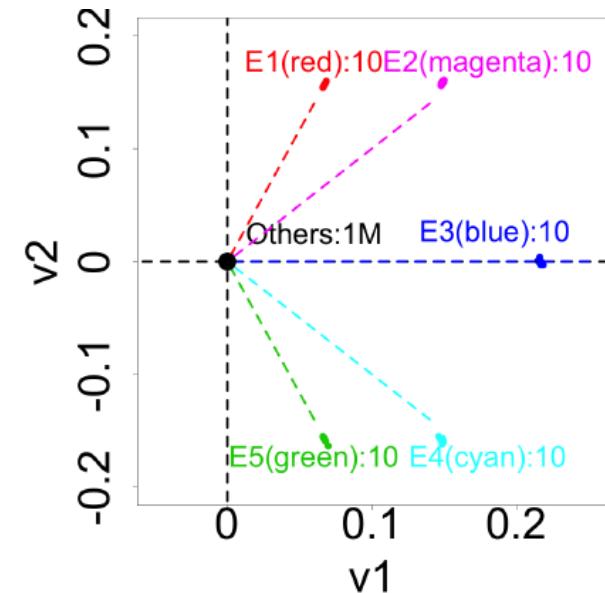
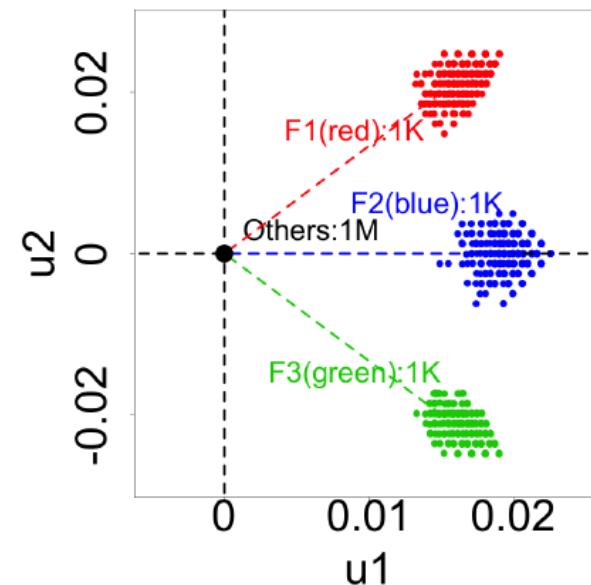
Spectral Subspace Plot: Case #4

- Overlapping: “Staircase”, “Pearls”

Adjacency Matrix



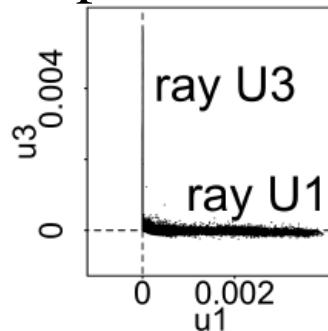
Spectral Subspace Plot



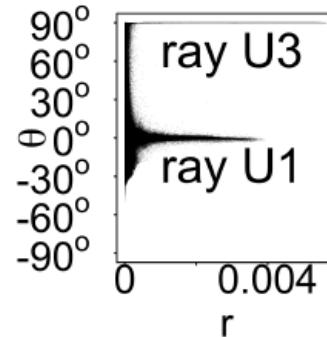
Rule 4 (“pearls”): a “staircase” of three partially overlapping blocks.

Algorithm: Reading & LockInfer

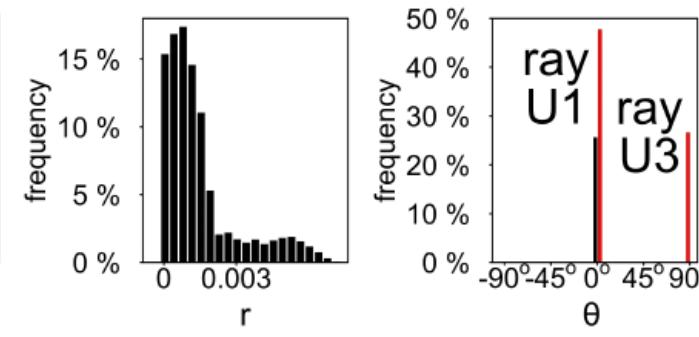
Spectral
Subspace Plot



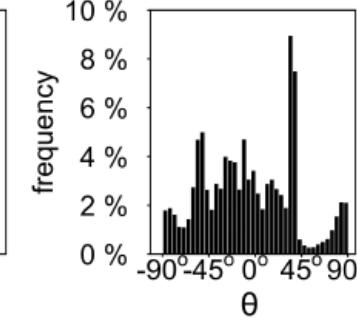
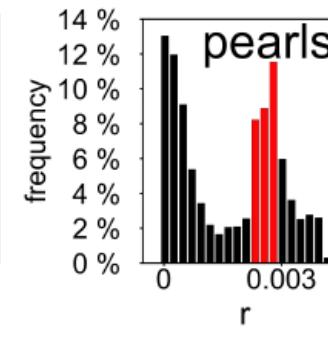
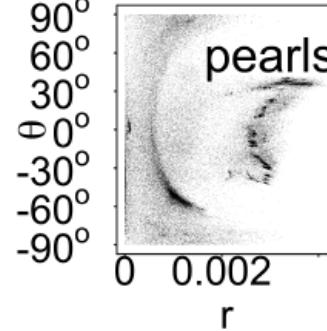
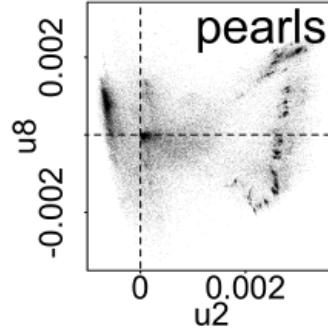
Polar Coordinate
Transform



Histograms



"rays" show two apparent spikes on θ frequency at 0° and 90°

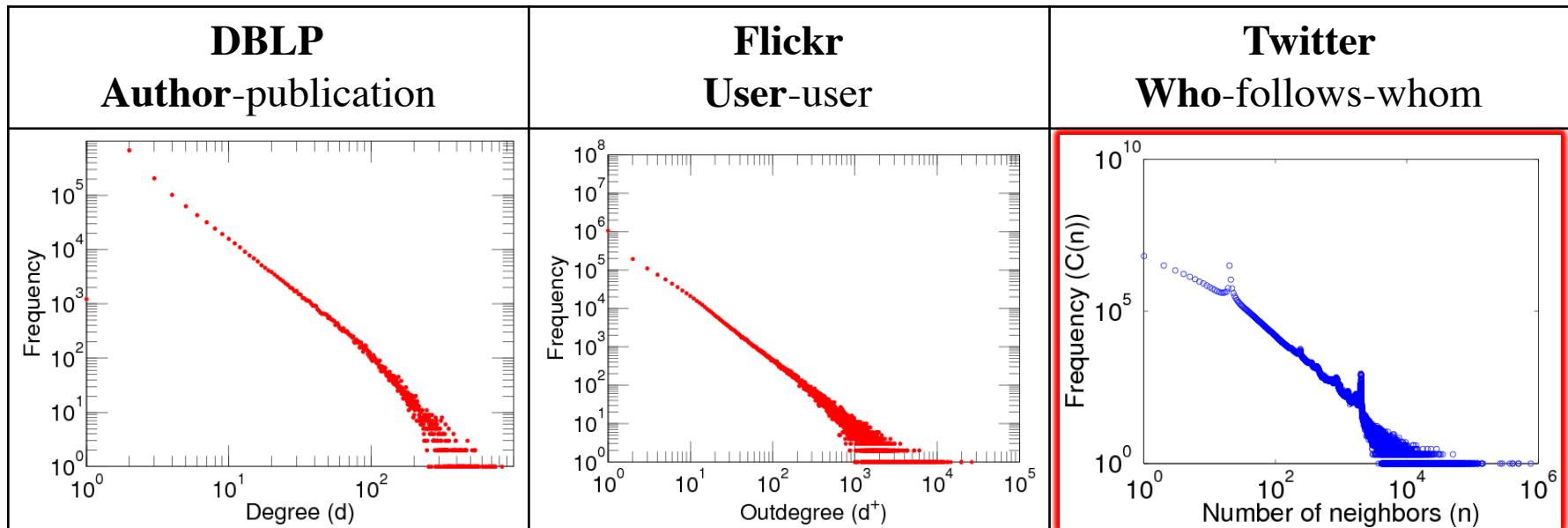


"pearls" show a spike on r frequency at a much-greater-than-zero value

High precision but low recall!!!

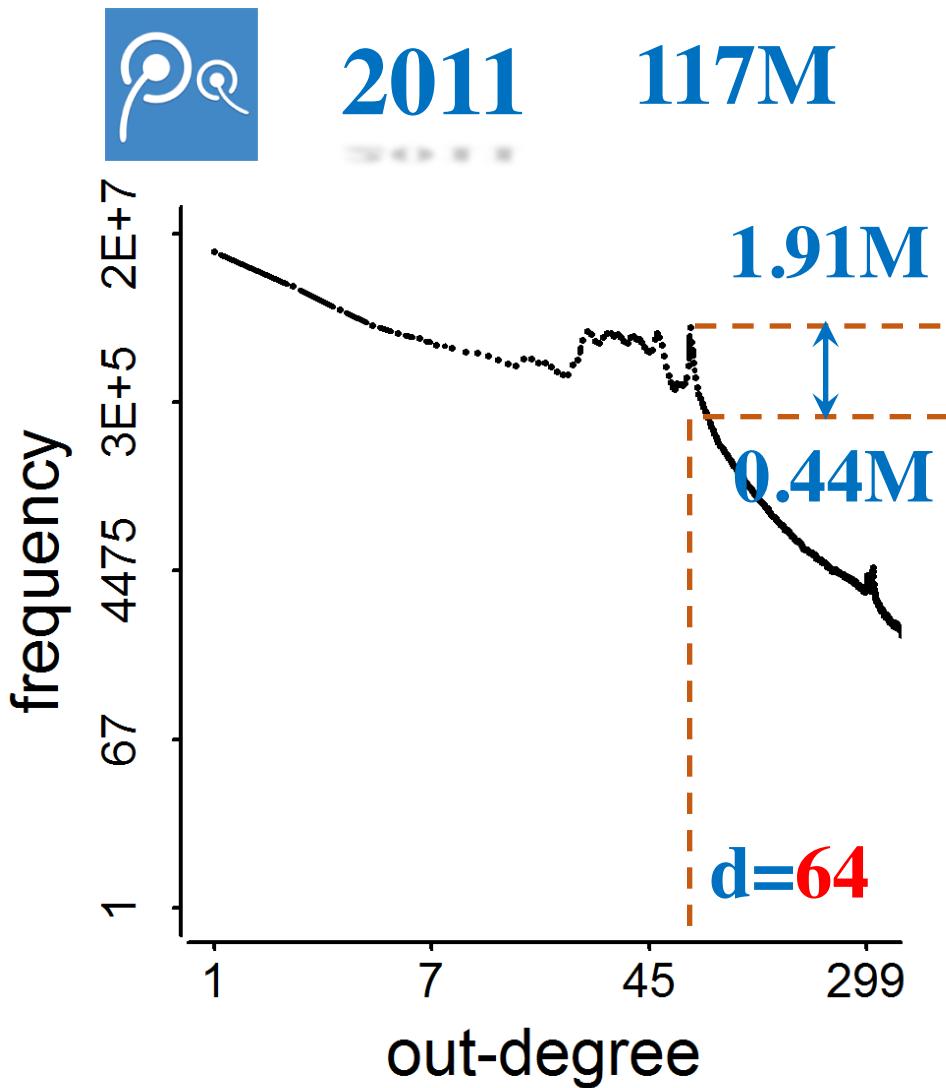
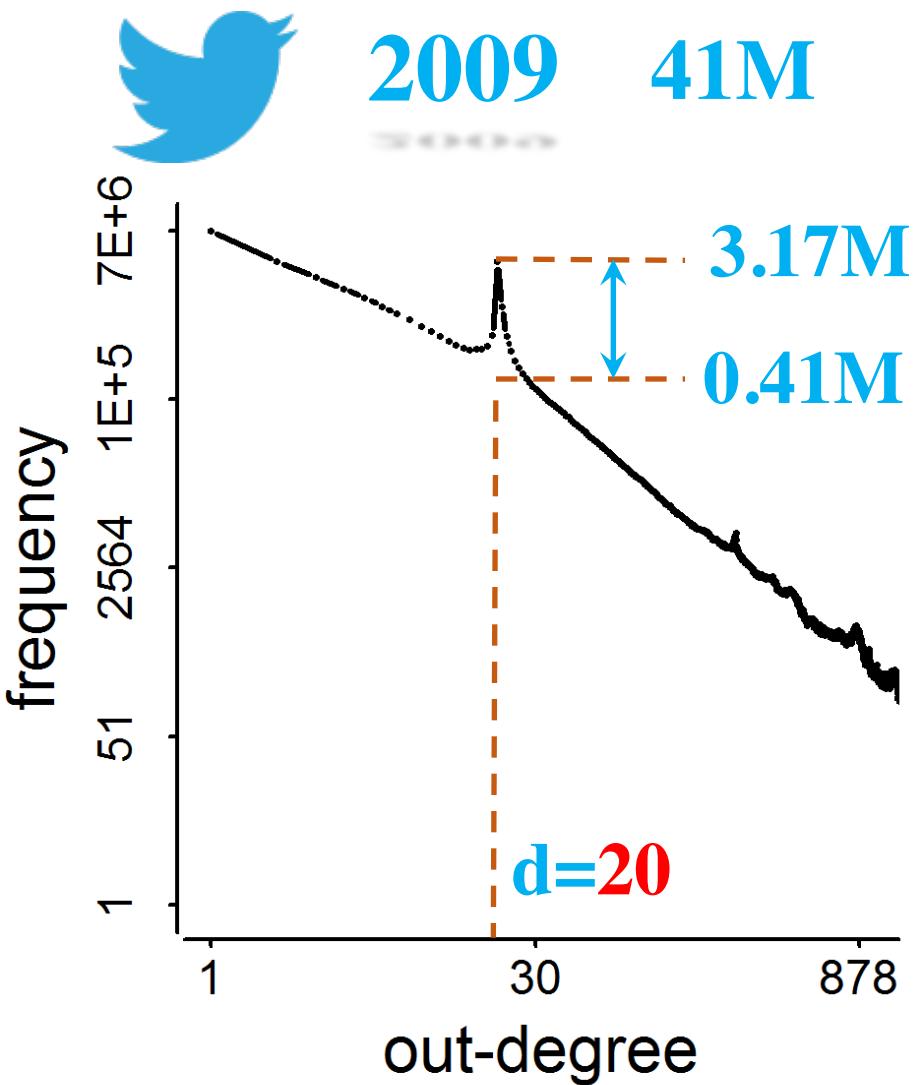
Out-Degree Distributions

- Power-law distribution [Faloutsos *et al.* SIGCOMM; Broder *et al.* Computer Networks; Chung *et al.* PNAS]



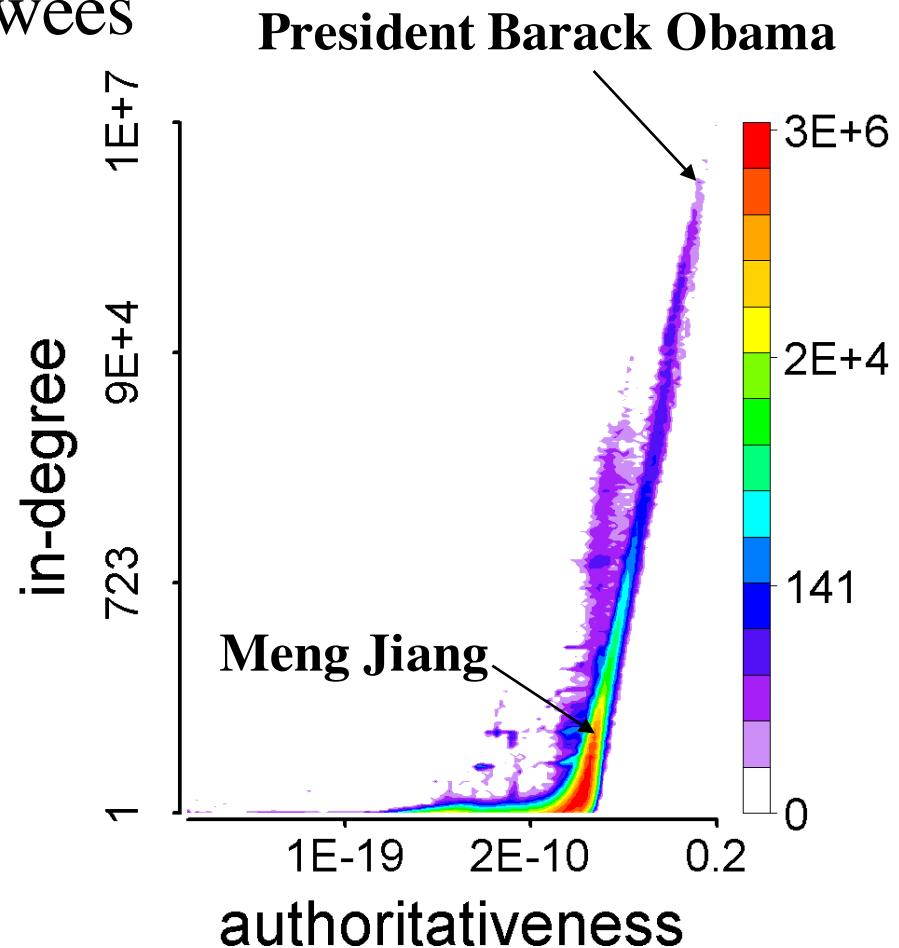
[konect.uni-koblenz.de/networks/]

Spikes!



Observation: How They Behave

- Feature space of followees [Kleinberg. JACM]



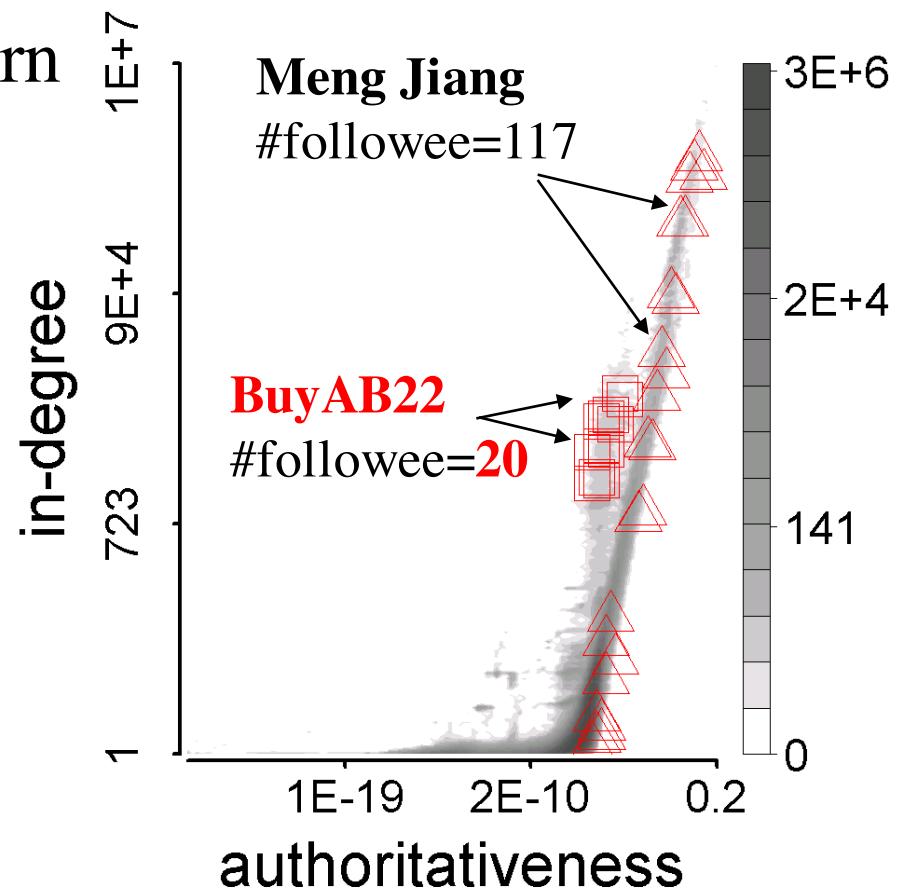
Observation: How They Behave

- Who are their followees?
- Their behavioral pattern
 - Synchronized

Similar with each other

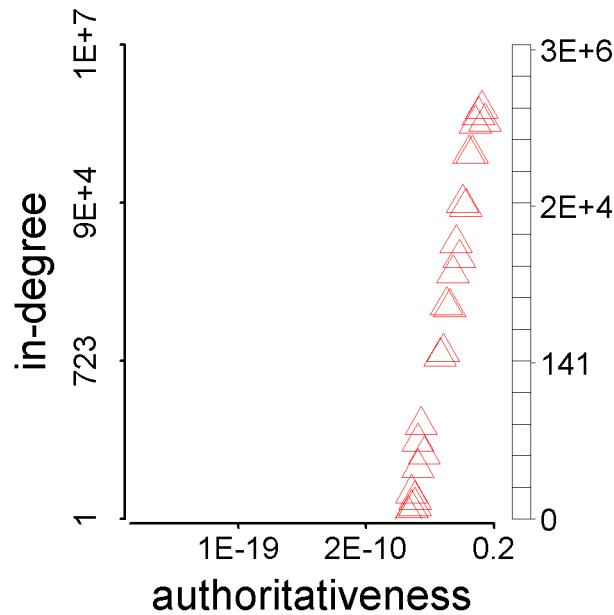
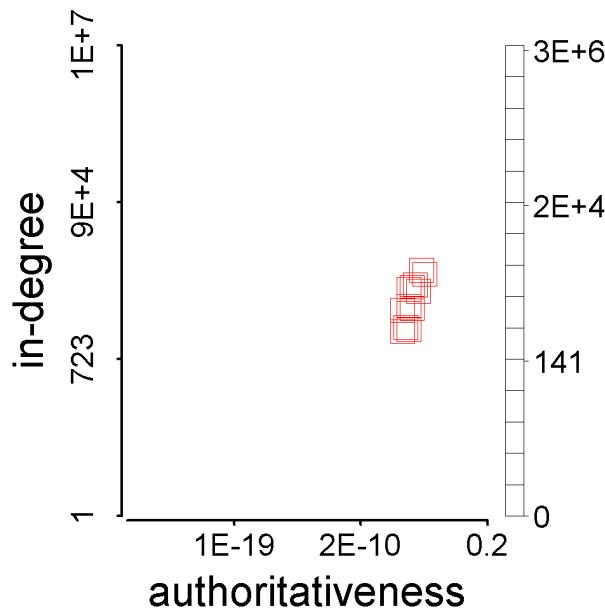
- Abnormal

Different from the majority



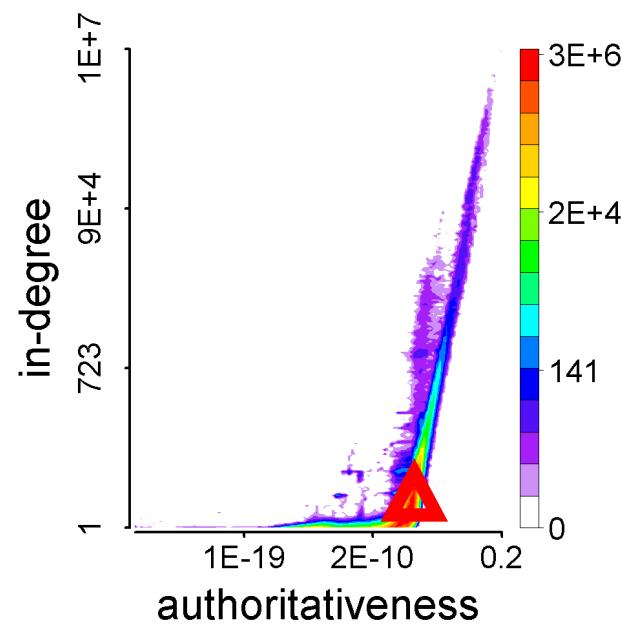
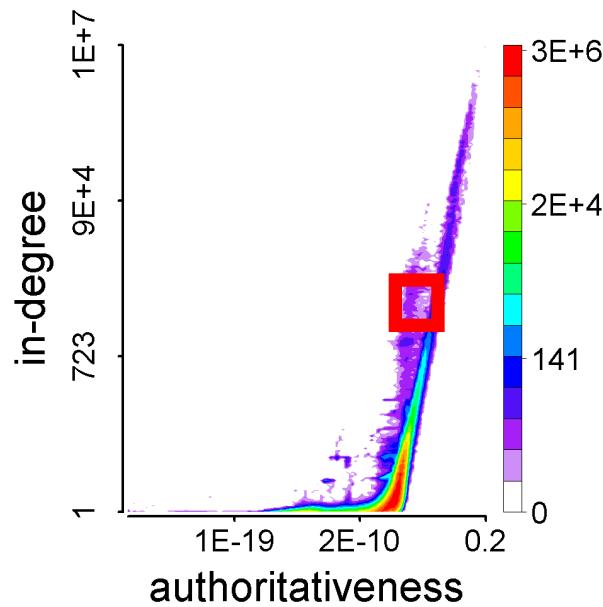
Represent Synchronicity

$$sync(u) = \frac{\sum_{(v, v') \in \mathcal{F}(u) \times \mathcal{F}(u)} \mathbf{p}(v) \cdot \mathbf{p}(v')}{d(u) \times d(u)}$$



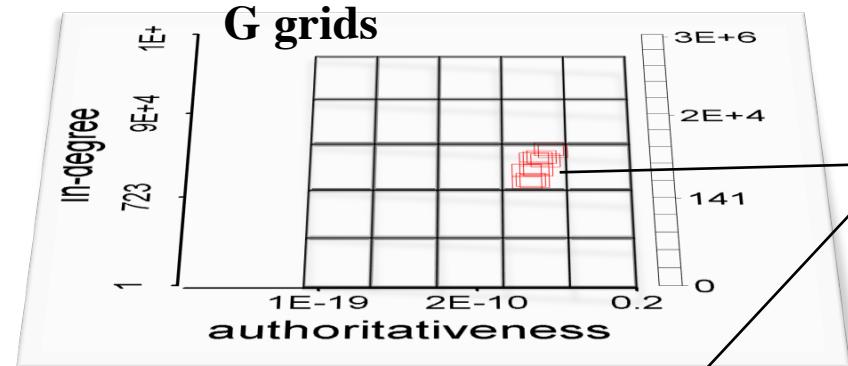
Represent Normality

$$\text{norm}(u) = \frac{\sum_{(v,v') \in \mathcal{F}(u) \times \mathcal{U}} \mathbf{p}(v) \cdot \mathbf{p}(v')}{d(u) \times N}$$

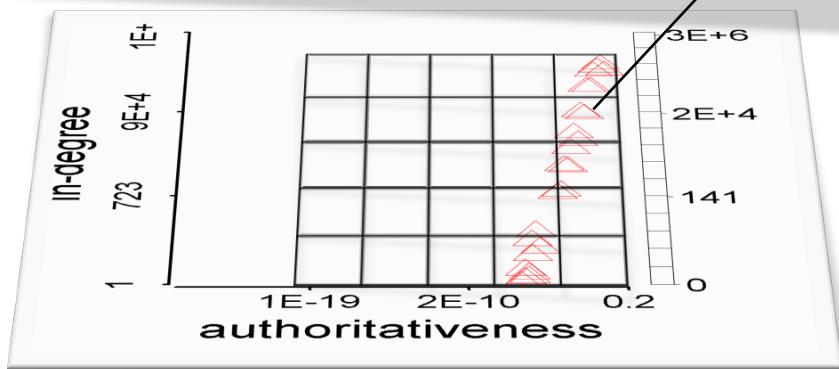




Theorem: Synchronicity vs. Normality



fp_g : #foreground points in grid g
 $\sum fp_g = F = d(u)$ (#followees of u)



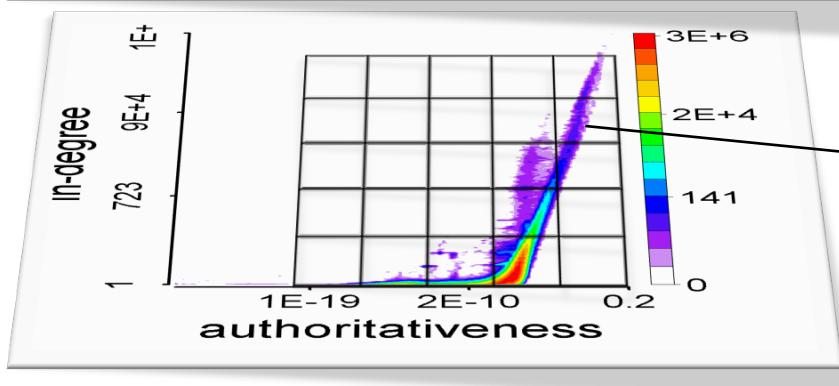
Given normality

$n = \sum(fp_g/F)(bp_g/B) = \sum f_g b_g$,
find minimal synchronicity

$$s = \sum(fp_g/F)(fp_g/F) = \sum f_g^2$$

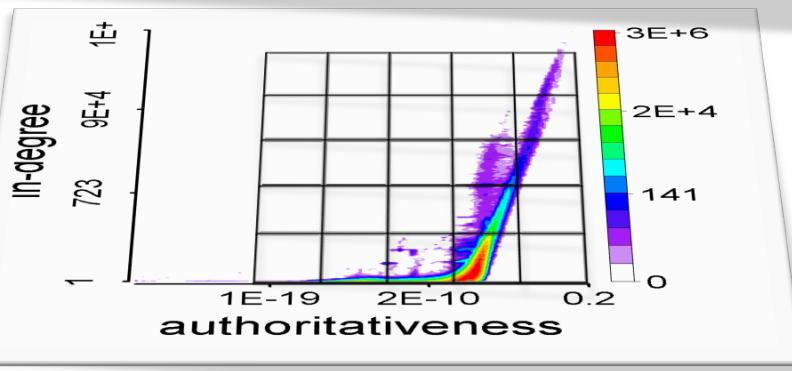
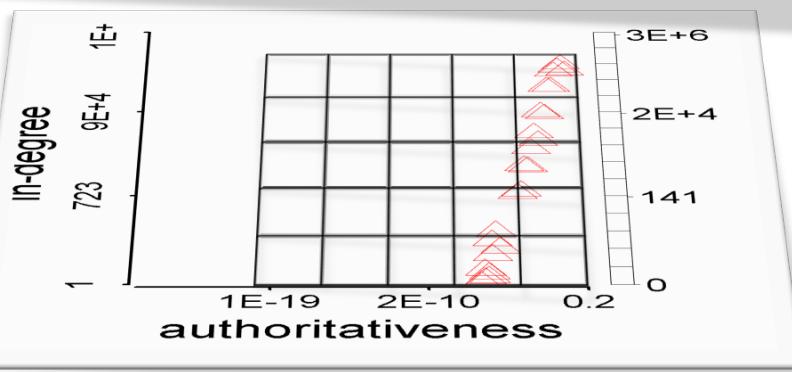
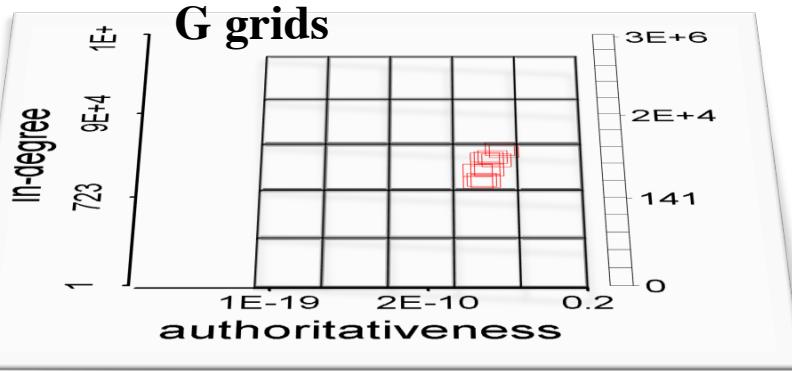
where

$$\sum f_g = 1, \sum b_g = 1$$



bp_g : #background points in grid g
 $\sum bp_g = B = N$ (#all users)

Theorem: Synchronicity vs. Normality



Solution.

Lagrange multiplier:

$$\text{minimize } s(f_g) = \sum f_g^2$$

$$\text{subject to } \sum f_g = 1, \sum f_g b_g = n$$

Lagrange function:

$$F(f_g, \lambda, \mu) = (\sum f_g^2) + \lambda(\sum f_g - 1) + \mu(\sum f_g b_g - n)$$

Gradients:

$$\begin{cases} \nabla_{f_g} F = 2 f_g + \lambda + \mu b_g = 0 \\ \nabla_{\lambda} F = \sum f_g - 1 = 0 \\ \nabla_{\mu} F = \sum f_g b_g - n = 0 \end{cases}$$

$$\begin{cases} 2 + \lambda G + \mu = 0 \\ 2 n + \lambda + \mu s_b = 0 \\ 2 s_{\min} + \lambda + \mu n = 0 \end{cases}$$

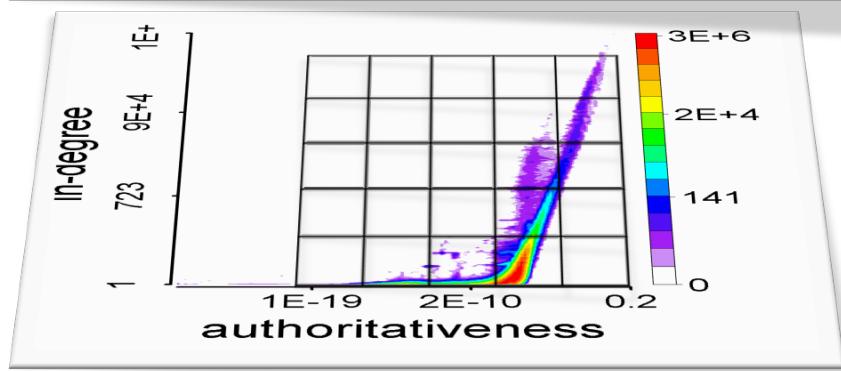
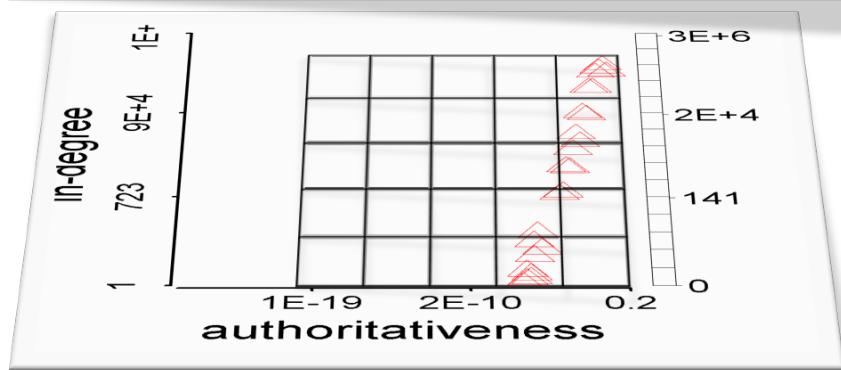
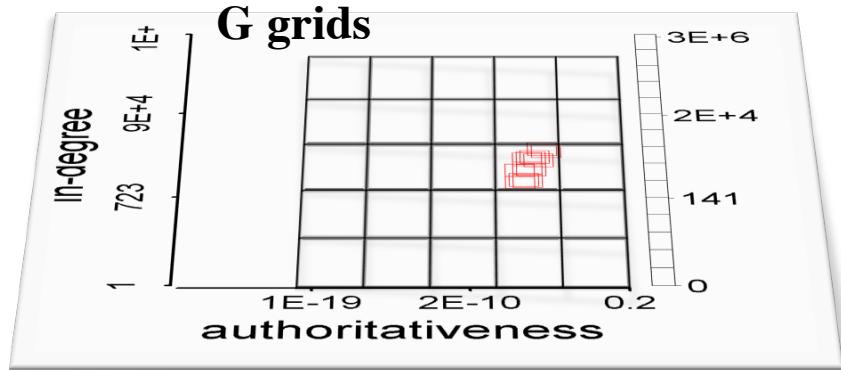
Σ $\times b_g \Sigma$ $\times f_g \Sigma$

where $s_b = \sum b_g^2$.

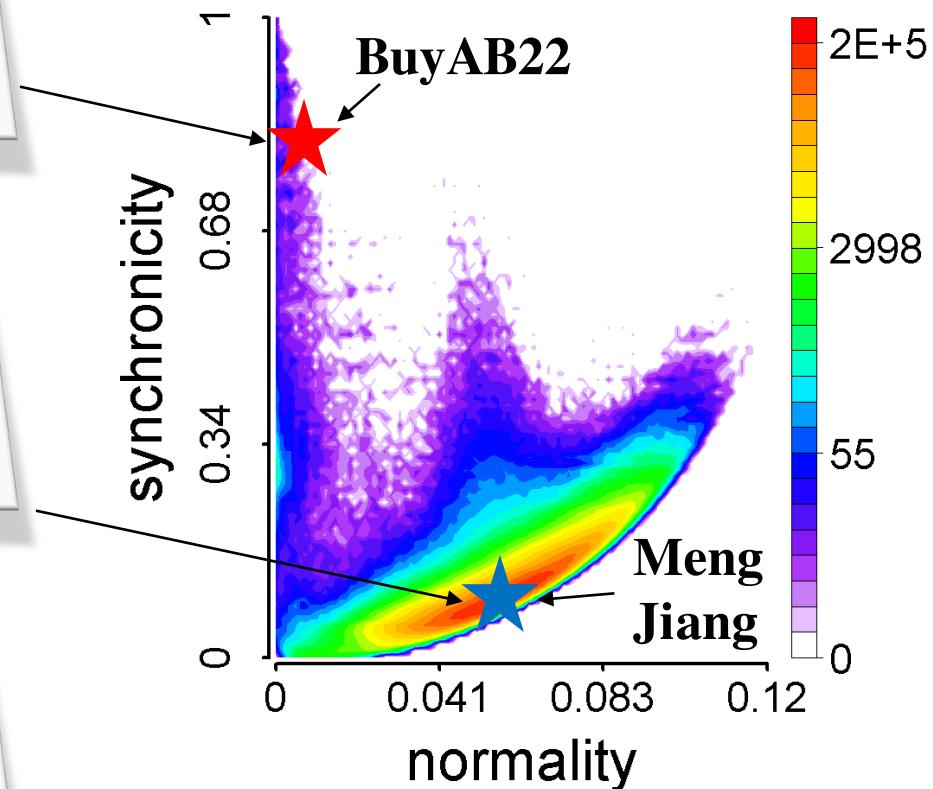
Therefore,

$$s_{\min} = \frac{-G n^2 + 2 n - s_b}{1 - G s_b}$$

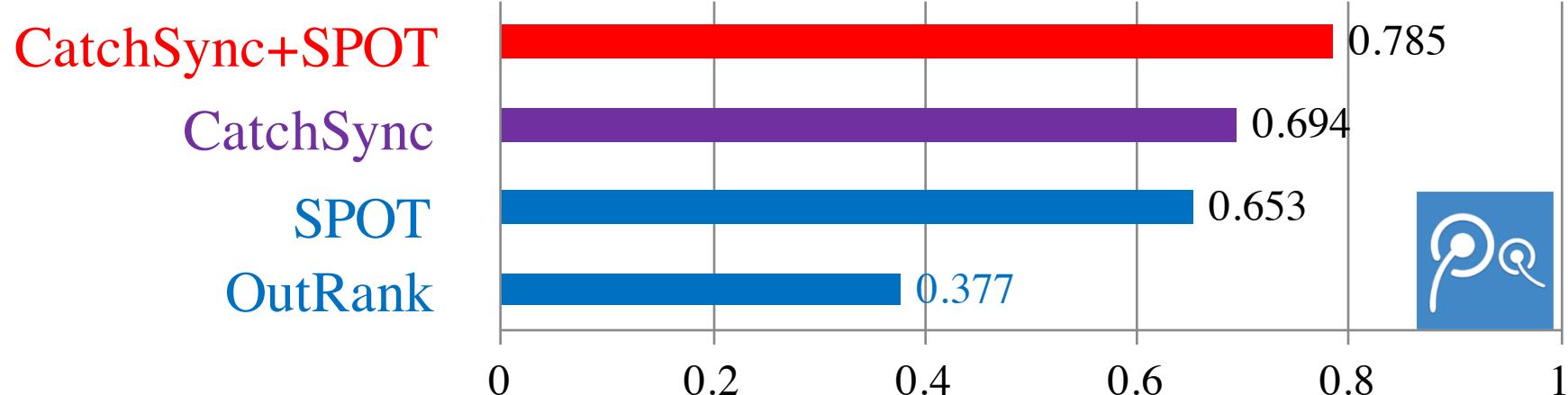
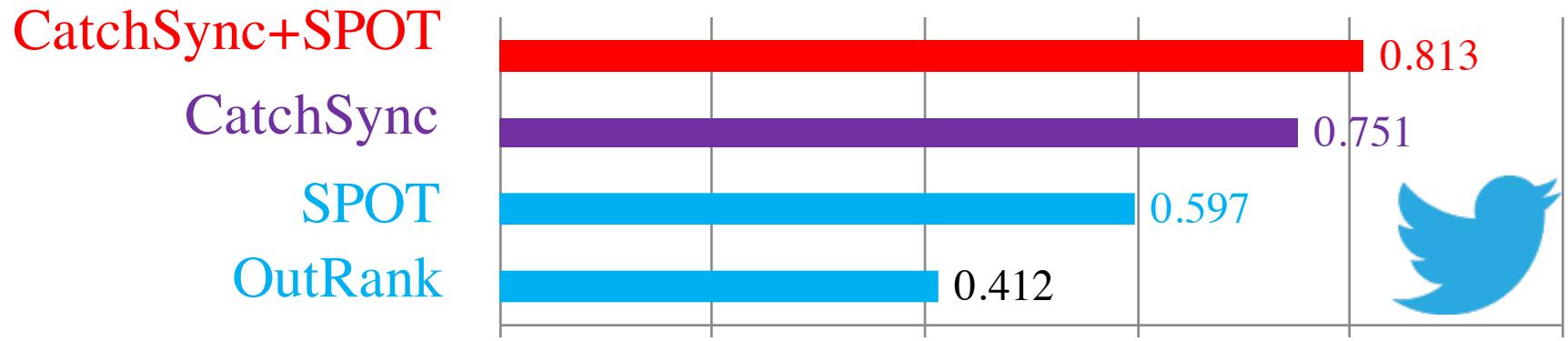
CatchSync Algorithm



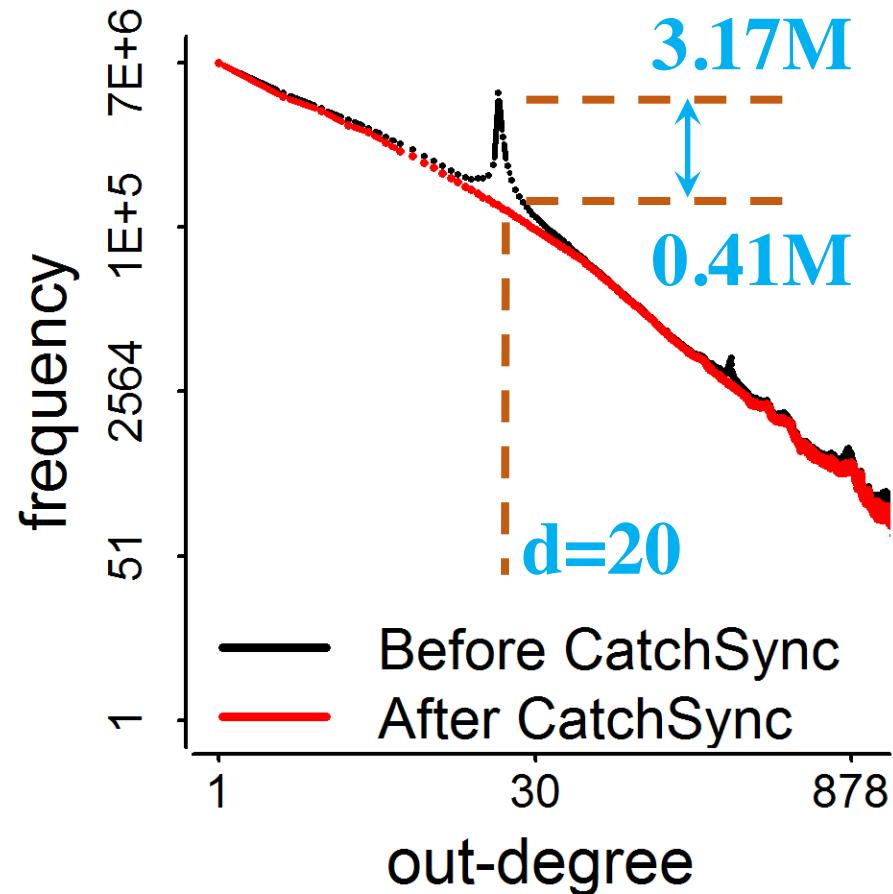
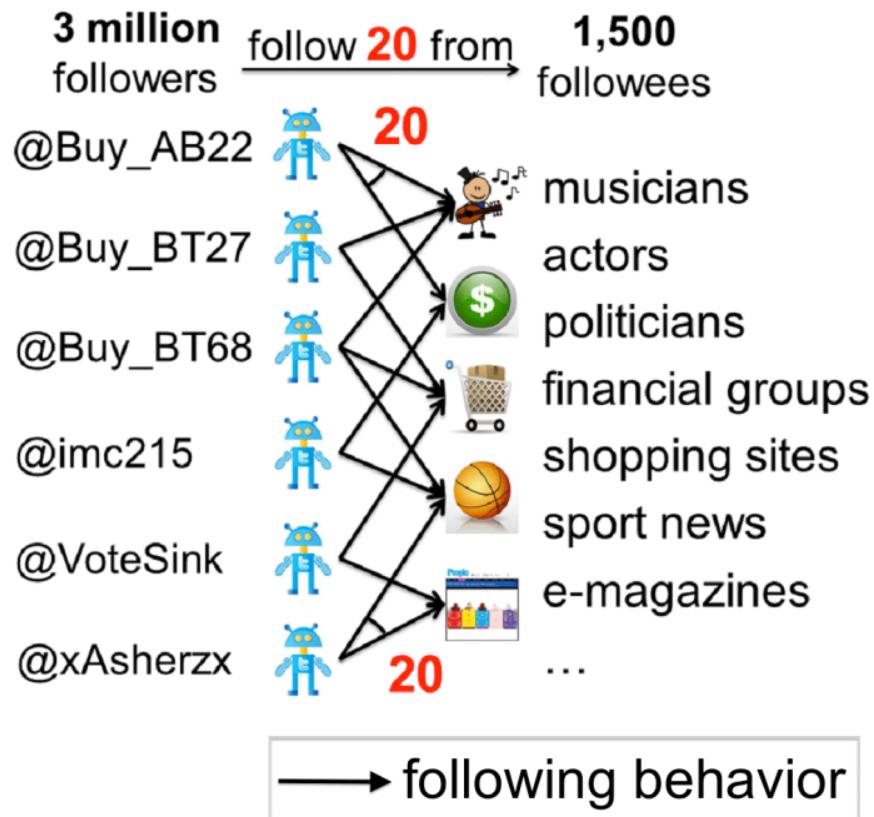
$$s_{\min} = \frac{-G n^2 + 2 n - s_b}{1 - G s_b}$$



Experimental Results



Experimental Results





Impact

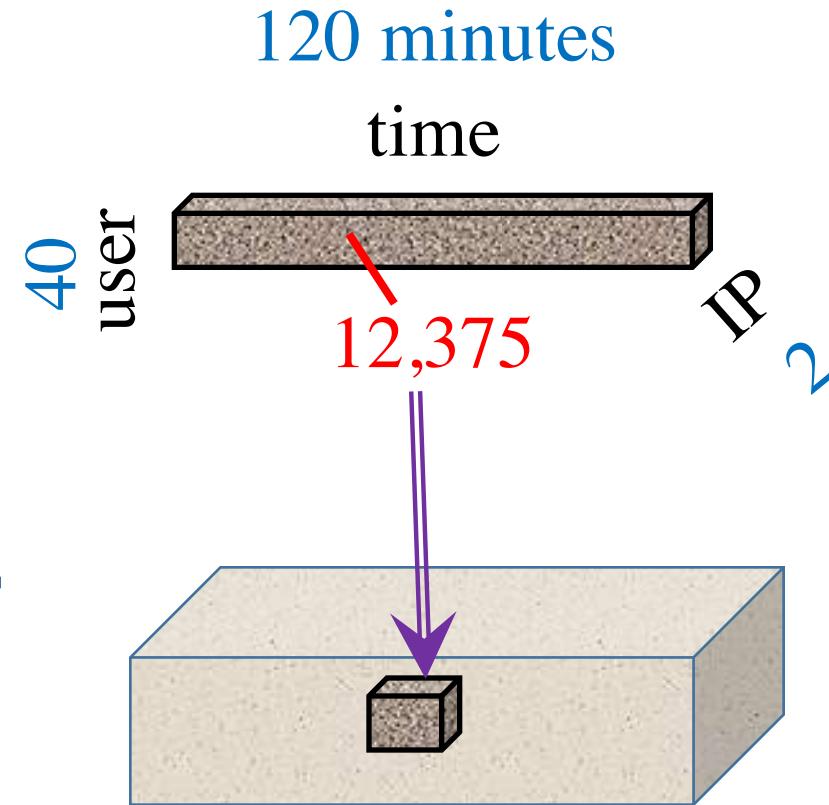
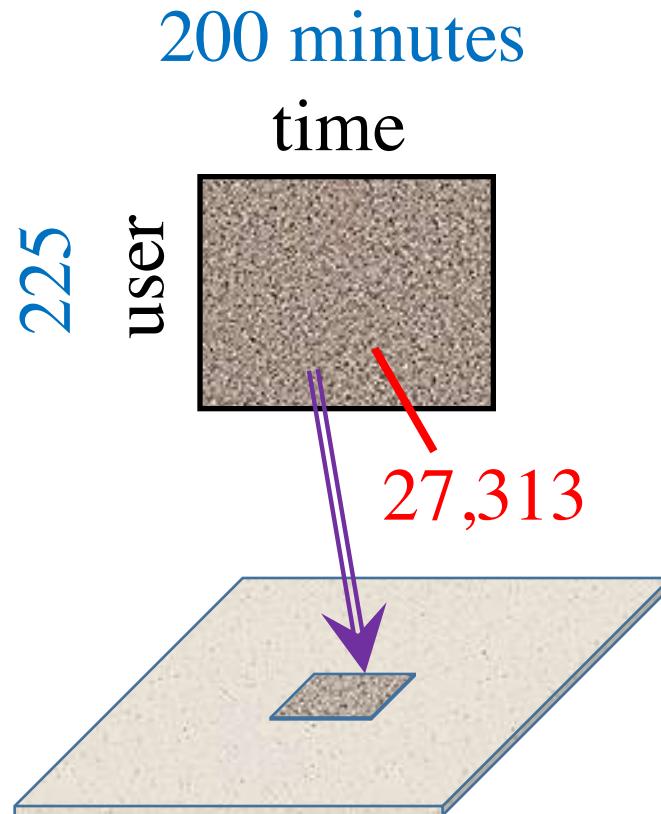
- ❑ M. Jiang, P. Cui, A. Beutel, C. Faloutsos and S. Yang.
“CatchSync: Catching Synchronized Behavior in Large
Directed Graphs” in **KDD’14 Best Paper Finalist**, Aug
2014. (#citations = **36**)
- ❑ Taught in
 - ❑ CMU 15-826: [Multimedia Databases and Data Mining](#)
 - ❑ UMich EECS 598: [Graph Mining and Exploration at Scale](#)
 - ❑ ASONAM’16 Tutorial: “[Identifying Malicious Actors on Social
Media](#)” by S. Kumar, F. Spezzano, V.S. Subrahmanian
- ❑ Deployed in Weibo? Unfortunately, in July 2014...



Observation: Spatiotemporal Contexts

Dataset	Dimension/Mode				Mass
Weibo's Retweeting	User	Root ID	IP	Time (min)	#retweet
	29.5M	19.8M	27.8M	56.9K	211.7M
Weibo's Trending (Hashtag)	User	Hashtag	IP	Time (min)	#tweet
	81.2M	1.6M	47.7M	56.9K	276.9M
Network attacks (LBNL)	Src-IP	Dest-IP	Port	Time (sec)	#packet
	2,345	2,355	6,055	3,610	230,836

Dense Block Indicates Suspiciousness

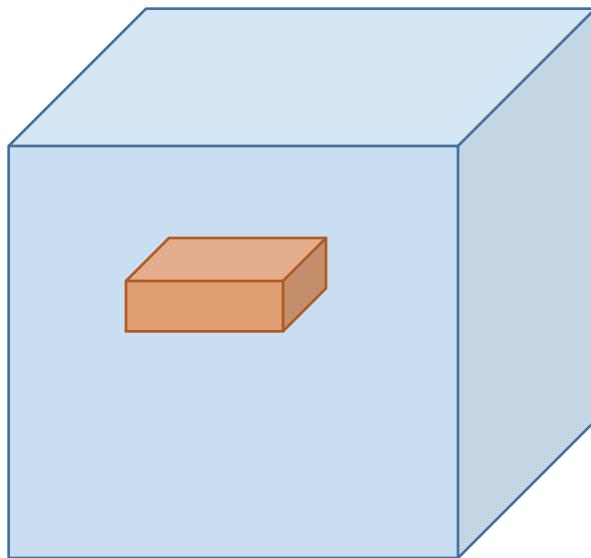


Q: Which is more suspicious?

We need a metric to evaluate the suspiciousness.

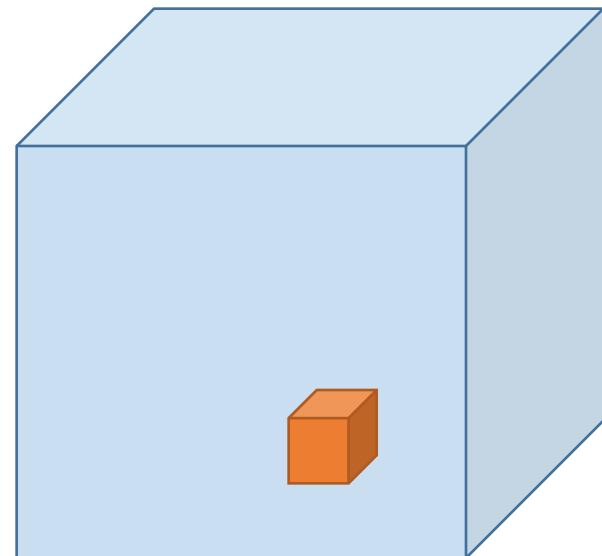
Criteria for Suspiciousness Metric

What properties are required of a good metric?



$$N_1 \times N_2 \times N_3$$

Count data with
total “mass” C



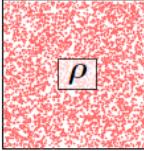
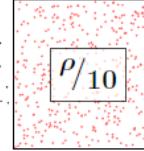
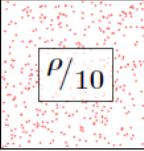
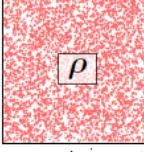
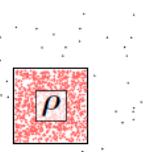
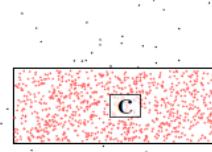
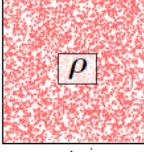
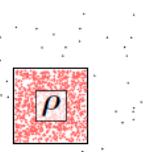
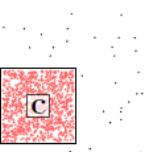
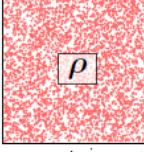
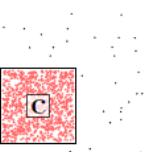
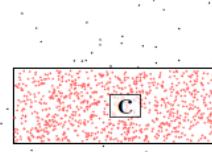
$$f\left(\begin{array}{c} n_1 \times n_2 \times n_3 \\ \text{mass } c \\ \text{density } \rho \end{array}\right)$$

VS

$$f\left(\begin{array}{c} n'_1 \times n'_2 \times n'_3 \\ \text{mass } c' \\ \text{density } \rho' \end{array}\right)$$

Axioms: 1 to 4

$$c_1 > c_2 \iff f(\mathbf{n}, c_1, \mathbf{N}, C) > f(\mathbf{n}, c_2, \mathbf{N}, C)$$

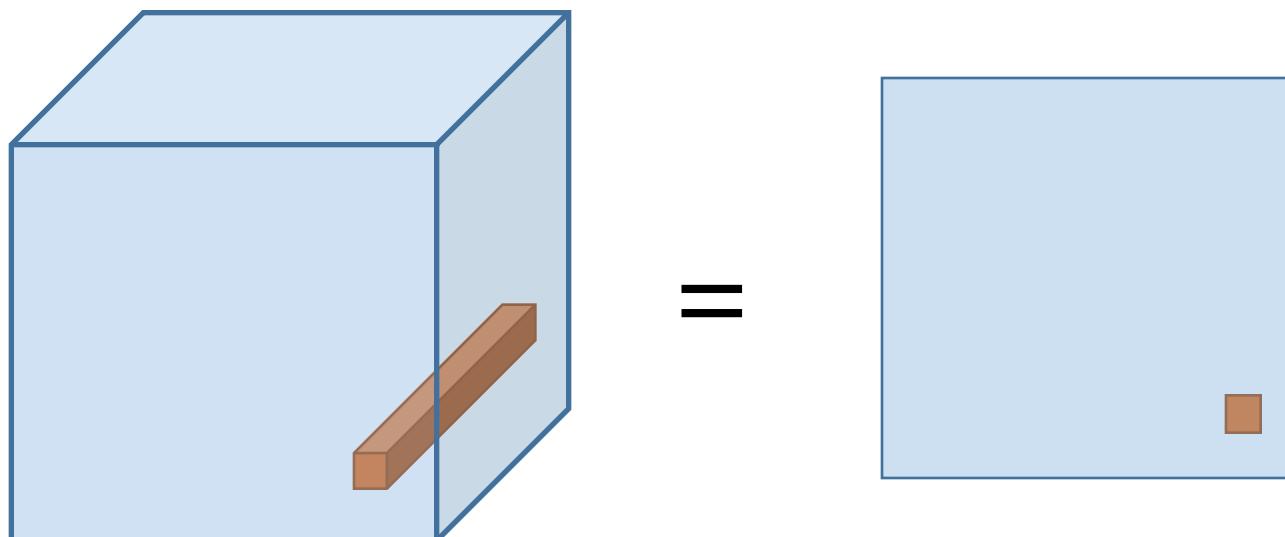
Density Axiom		Contrast Axiom	
	>		
	>		
Size Axiom		Concentration Axiom	
	>		
	>		

$$p_1 < p_2 \iff \hat{f}(\mathbf{n}, \rho, \mathbf{N}, p_1) > \hat{f}(\mathbf{n}, \rho, \mathbf{N}, p_2)$$

Axiom 5: Cross Dimensions

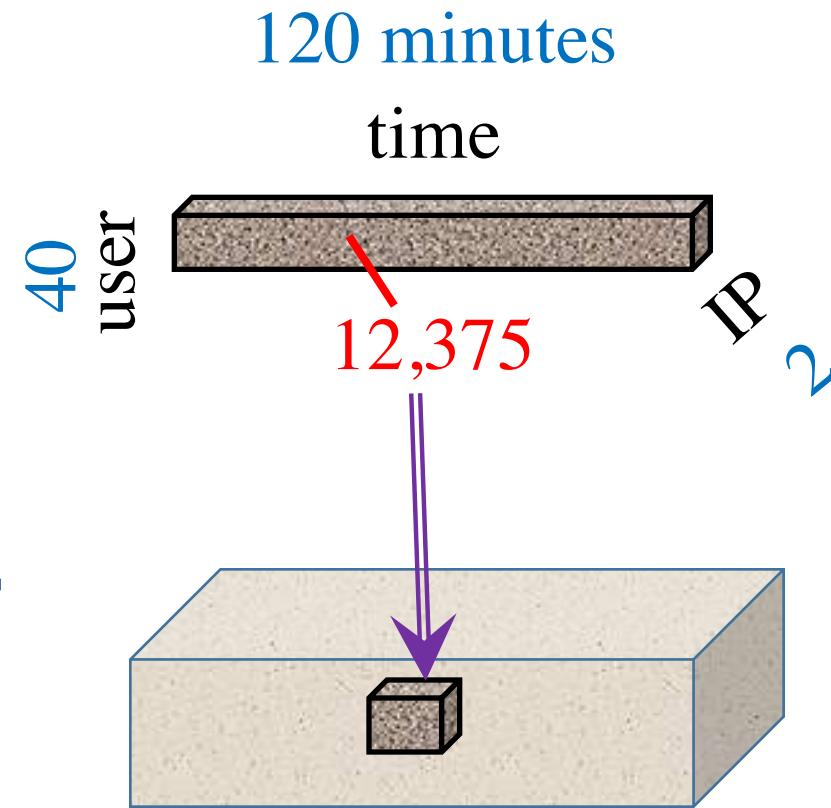
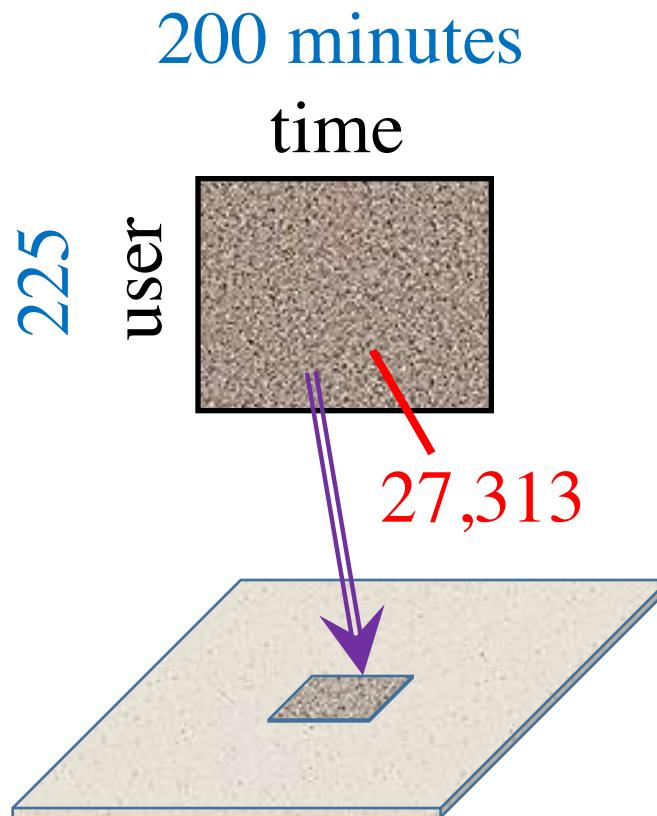
$$f_{K-1} \left([n_k]_{k=1}^{K-1}, c, [N_k]_{k=1}^{K-1}, C \right) = f_K \left(([n_k]_{k=1}^{K-1}, N_K), c, [N_k]_{k=1}^K, C \right)$$

Not including a mode is the same as including all values for that mode.



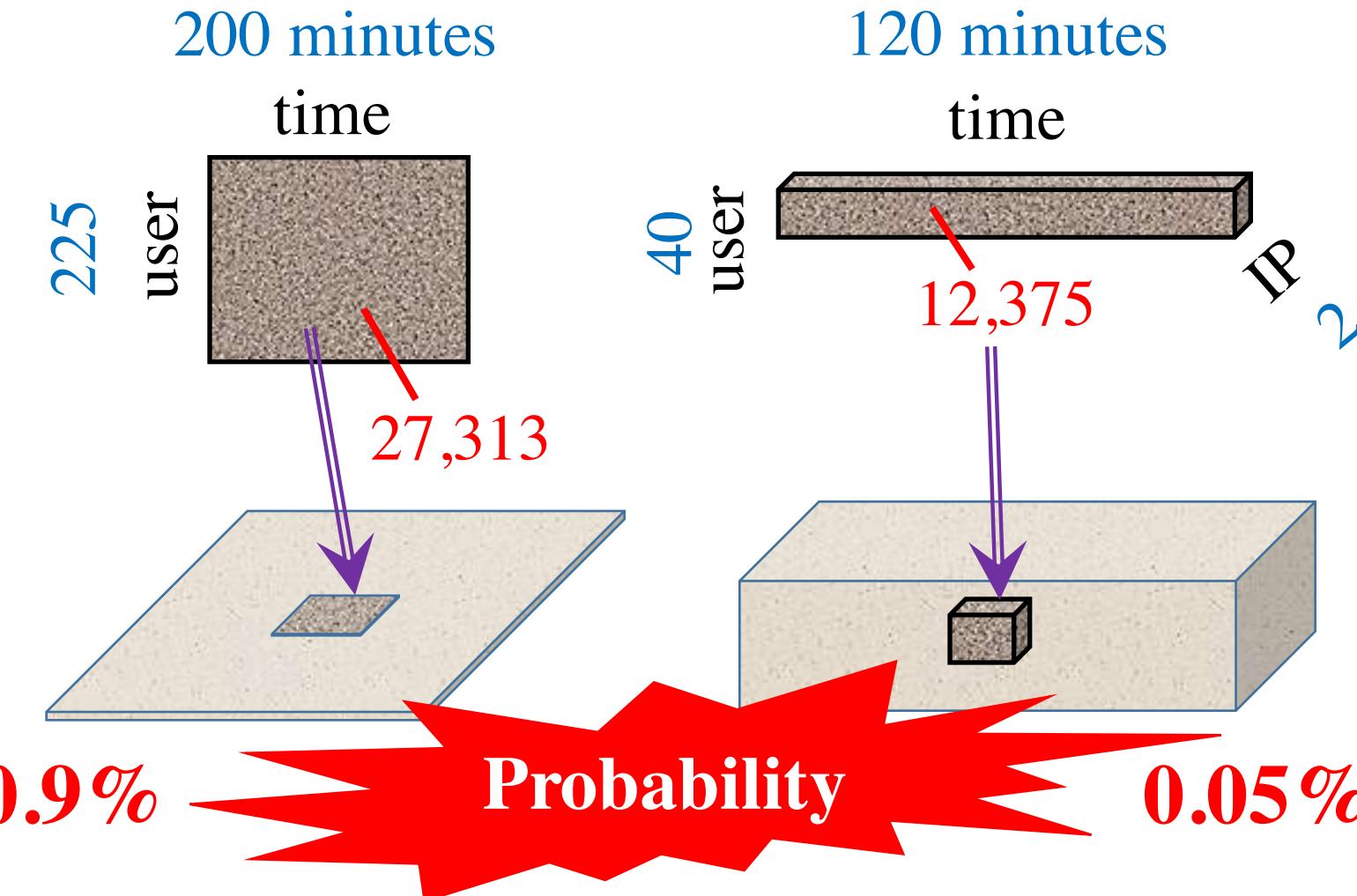
- New information (more modes) can only make our blocks more suspicious

Scoring the Suspiciousness



Q: Which is more suspicious?

Scoring the Suspiciousness





A General Suspiciousness Metric

- ❑ Negative log likelihood of block's probability

$$f(n, c, N, C) = -\log [Pr(Y_n = c)]$$

Lemma Given an $n_1 \times \cdots \times n_K$ block of mass c in $N_1 \times \cdots \times N_K$ data of total mass C , the suspiciousness function is

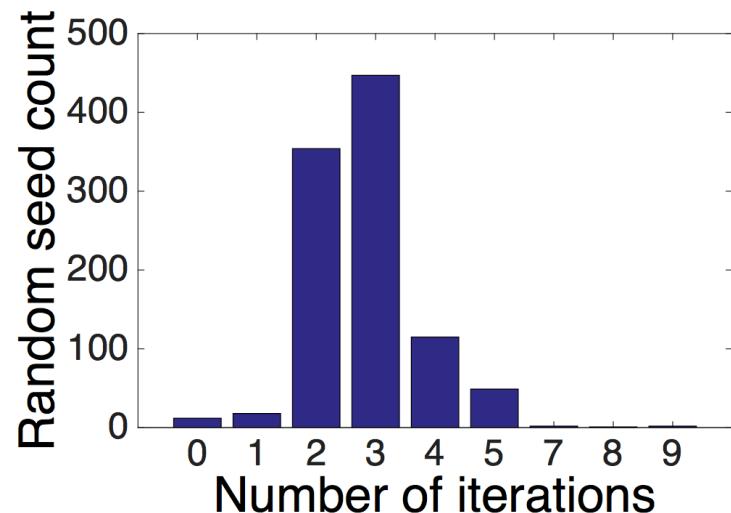
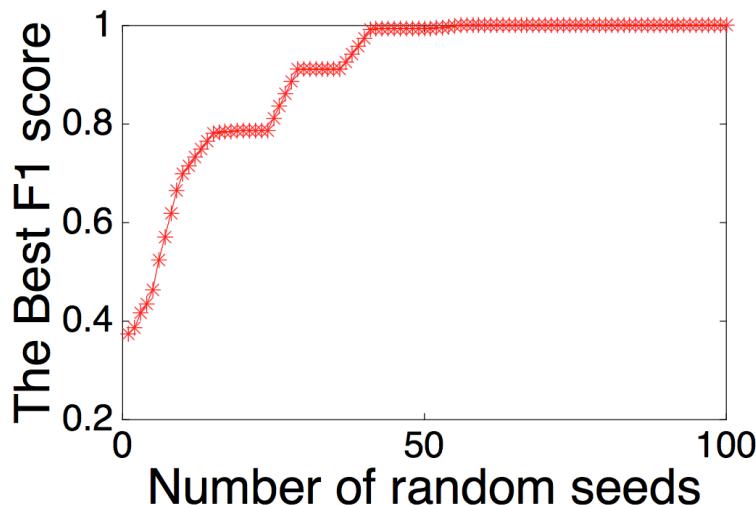
$$f(\mathbf{n}, c, \mathbf{N}, C) = c(\log \frac{c}{C} - 1) + C \prod_{i=1}^K \frac{n_i}{N_i} - c \sum_{i=1}^K \log \frac{n_i}{N_i}$$

Using ρ as the block's density and p is the data's density, we have the simpler formulation

$$\hat{f}(\mathbf{n}, \rho, \mathbf{N}, p) = \left(\prod_{i=1}^K n_i \right) D_{KL}(\rho || p)$$

CrossSpot Algorithm

- ❑ Local search to maximize the metric
 - ❑ Start with seed blocks
 - ❑ Parameter-free: iteratively update the blocks
 - ❑ Scalable: parallelize to multiple machines





Advantages

		Axioms				
		Density	Size	3 Concentration	Contrast	Multi-modal
Method		Scores				
Metrics	SUSPICIOUSNESS	✓	✓	✓	✓	✓
	Mass	✓	✓	✗	✗	✗
	Density	✓	✓	✗	✓	✗
	Average Degree [9]	✓	✓	✗	✗	N/A
	Singular Value [10]	✓	✓	✓	✓	✗
	CROSSSPOT	✓	✓	✓	✓	✓
Methods	Subgraph [30, 10, 36]	✓	✓	✓	✓	N/A
	CopyCatch [6]	✓	✓	✓	✓	N/A
	EigenSpokes [31]	✗	N/A			
	TrustRank [14, 8]	✗	N/A			
	BP [28, 1]	✗	N/A			

Results: Dense Block Detection

□ Synthetic data

- $1,000 \times 1,000 \times 1,000$ of 10,000 random data
- Block#1: $30 \times 30 \times 30$ of 512 3 modes
- Block#2: $30 \times 30 \times 1,000$ of 512 2 modes
- Block#3: $30 \times 1,000 \times 30$ of 512 2 modes
- Block#4: $1,000 \times 30 \times 30$ of 512 2 modes

	Recall				Overall Evaluation		
	Block #1	Block #2	Block #3	Block #4	Precision	Recall	F1 score
HOSVD ($r=20$)	93.7%	29.5%	23.7%	21.3%	0.983	0.407	0.576
HOSVD ($r=10$)	91.3%	24.4%	18.5%	19.2%	0.972	0.317	0.478
HOSVD ($r=5$)	85.7%	10.0%	9.5%	11.4%	0.952	0.195	0.324
CROSSSPOT	100 %	99.9 %	94.9 %	95.4 %	0.978	0.967	0.972



Results: Tweeting Hashtags

User × hashtag × IP × minute	Mass c	Suspiciousness
$582 \times 3 \times 294 \times \mathbf{56,940}$	5,941,821	111,799,948
$188 \times 1 \times 313 \times \mathbf{56,943}$	2,344,614	47,013,868
$75 \times 1 \times 2 \times 2,061$	689,179	19,378,403

User ID	Time	IP address (city, province)	Tweet text with hashtag
USER-D	11-18 12:12:51	IP-1 (Deyang, Shandong)	#Snow# the Samsung GALAXY SII QQ Service customized version...
USER-E	11-18 12:12:53	IP-1 (Deyang, Shandong)	#Snow# the Samsung GALAXY SII QQ Service customized version...
USER-F	11-18 12:12:54	IP-2 (Zaozhuang, Shandong)	#Snow# the Samsung GALAXY SII QQ Service customized version...
USER-E	11-18 12:17:55	IP-1 (Deyang, Shandong)	#Li Ning - a weapon with a hero# good support activities!
USER-F	11-18 12:17:56	IP-2 (Zaozhuang, Shandong)	#Li Ning - a weapon with a hero# good support activities!
USER-D	11-18 12:18:40	IP-1 (Deyang, Shandong)	#Toshiba Bright Daren# color personality test to find out your sense...
USER-E	11-18 17:00:31	IP-2 (Zaozhuang, Shandong)	#Snow# the Samsung GALAXY SII QQ Service customized version...
USER-D	11-18 17:00:49	IP-2 (Zaozhuang, Shandong)	#Toshiba Bright Daren# color personality test to find out your sense...
USER-F	11-18 17:00:56	IP-2 (Zaozhuang, Shandong)	#Li Ning - a weapon with a hero# good support activities!



Results: Network Attacks

	#	Src-IP \times dst-IP \times port \times second	Mass c	Suspiciousness
CROSSSPOT	1	$411 \times 9 \times 6 \times 3,610$	47,449	552,465
	2	$533 \times 6 \times 1 \times 3,610$	30,476	400,391
	3	$5 \times 5 \times 2 \times 3,610$	18,881	317,529
	4	$11 \times 7 \times 7 \times 3,610$	20,382	295,869
HOSVD	1	$15 \times 1 \times 1 \times 1,336$	4,579	80,585
	2	$1 \times 2 \times 2 \times 1,035$	1,035	18,308
	3	$1 \times 1 \times 1 \times 1,825$	1,825	34,812
	4	$1 \times 13 \times 6 \times 181$	1,722	29,224



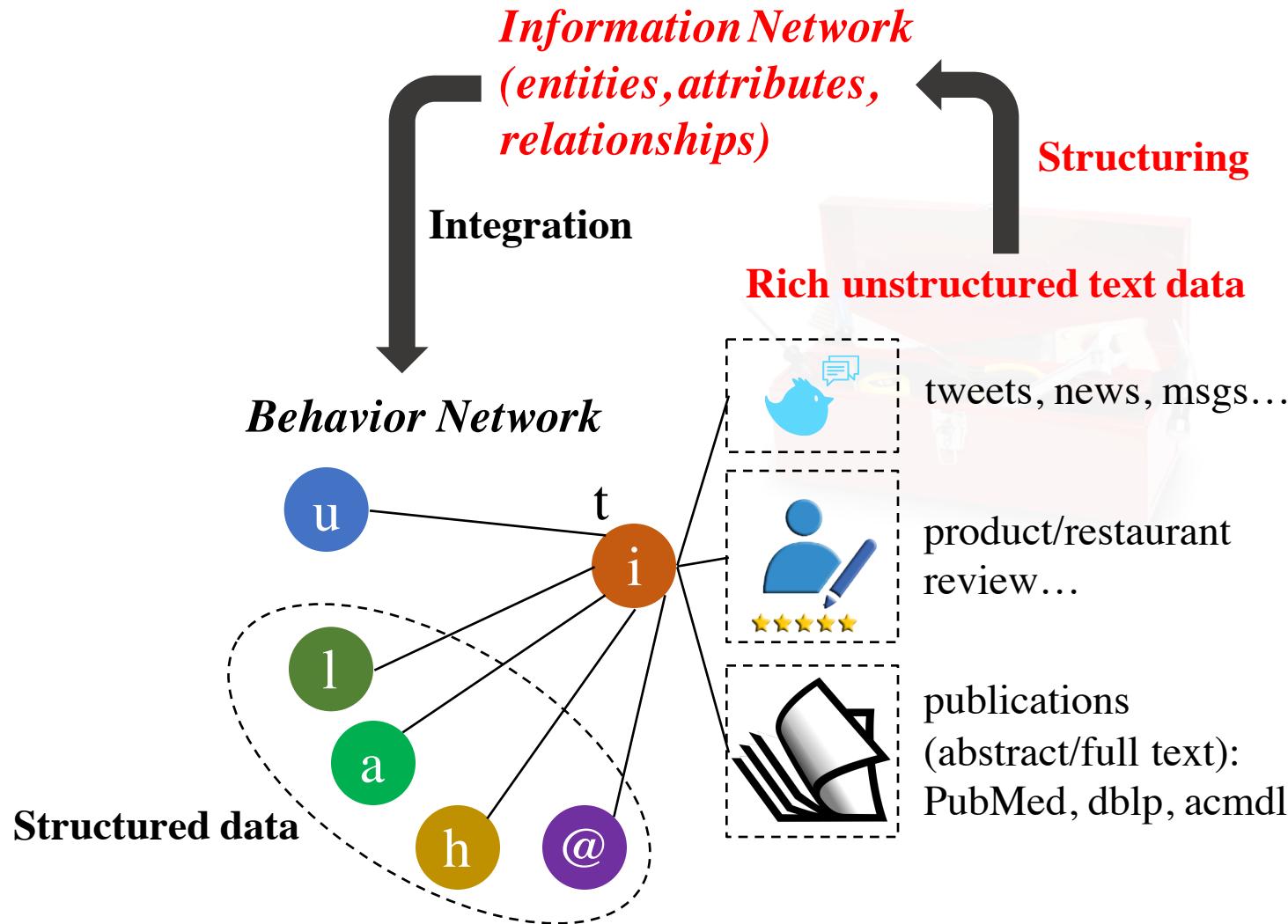
Summary

- ❑ Ill-gotten Facebook Likes, Zombie Followers
- ❑ **Observations, Representations, Models**
 - ❑ **CopyCatch:** Catching ill-gotten Likes by core search
 - ❑ **LockInfer:** Adding seed selection before search
 - ❑ **CatchSync:** Catching smart zombie followers with high recall (recovering power-law distributions)
 - ❑ **CrossSpot:** Defining suspiciousness across dimensions



II. Structuring behavioral content and integrating behavioral analysis with information networks

Data to Network to Knowledge

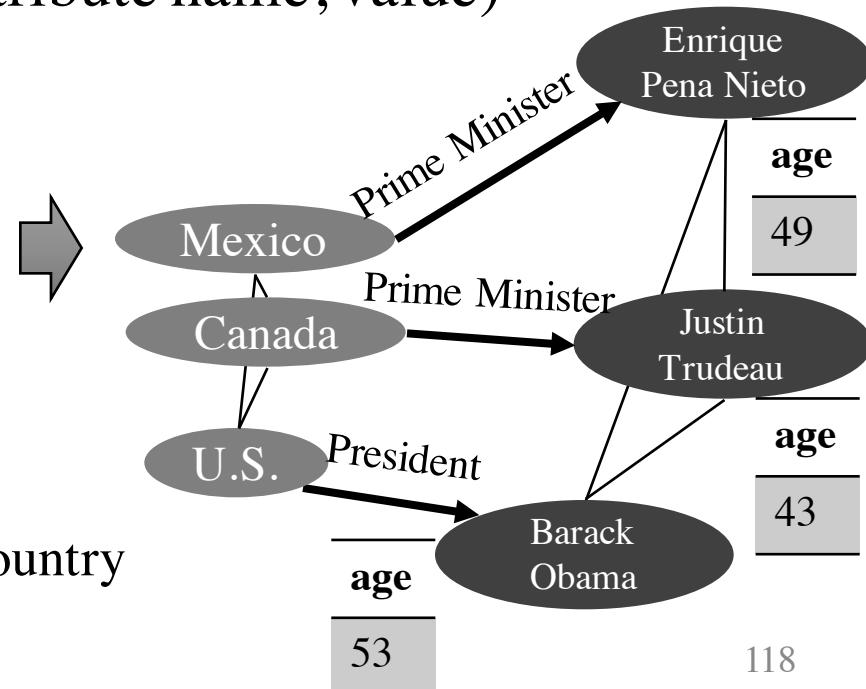


Construction of Heterogeneous Information Networks from Text

💡 Philosophy: Not extensive “labeling” but exploring the power of massive text corpora!

- ❑ Mining phrases (the minimal semantic units)
- ❑ Entity recognition and typing
- ❑ Attribute discovery (entity, attribute name, value)

...here by **Canada Prime Minister Justin Trudeau, 43**, the so-called #APEChottie...of Mexico's **Enrique Pena Nieto, 49**, ... United States President **Barack Obama, 53**, who...

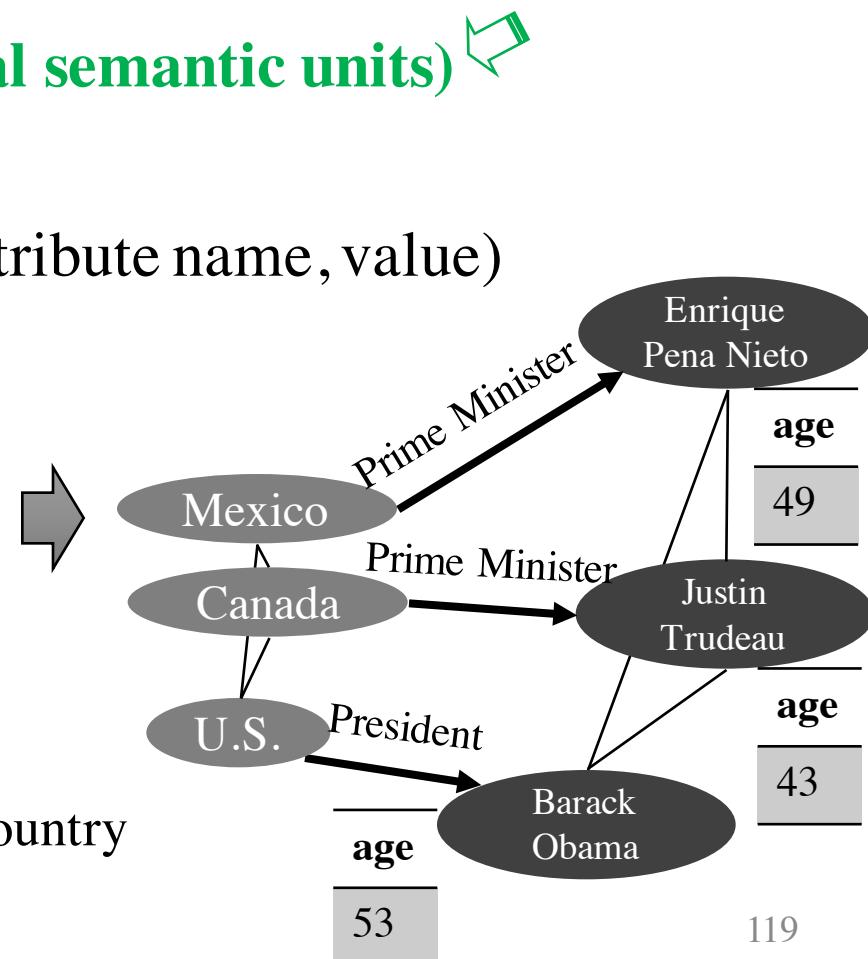


Construction of Heterogeneous Information Networks from Text

💡 Philosophy: Not extensive “labeling” but exploring the power of massive text corpora!

- ❑ Mining phrases (the minimal semantic units)
- ❑ Entity recognition and typing
- ❑ Attribute discovery (entity, attribute name, value)

...here by **Canada Prime Minister Justin Trudeau, 43**, the so-called #APEChottie...of Mexico's **Enrique Pena Nieto, 49**, ... United States President **Barack Obama, 53**, who...



Why Mining Phrases?

- ❑ **Unigrams** are *ambiguous* but **phrases** are natural, *unambiguous* semantic units
 - ❑ Ex.: “United” vs. United States, United Airline, United Parcel Service
- ❑ Mining semantically meaningful phrases
 - ❑ Transform text data from *word granularity* to *phrase granularity*
 - ❑ Enhance the power at manipulating unstructured data using information networks
- ❑ Phrase mining: Most NLP methods may need annotation and training
 - ❑ Annotate hundreds of documents as training data
 - ❑ Train a supervised model based on part-of-speech features
 - ❑ Limitations: High annotation cost
 - ❑ May not be scalable to domain-specific, dynamic, emerging applications
 - ❑ Scientific domains, query logs, or social media, e.g., Yelp, Twitter
- 💡 Minimal/no training but making good use of massing corpora



Strategies for Phrase Mining

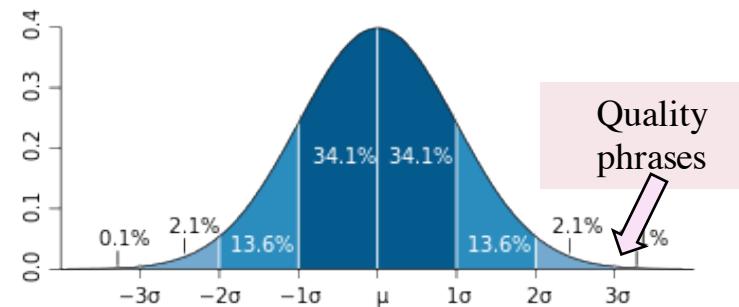
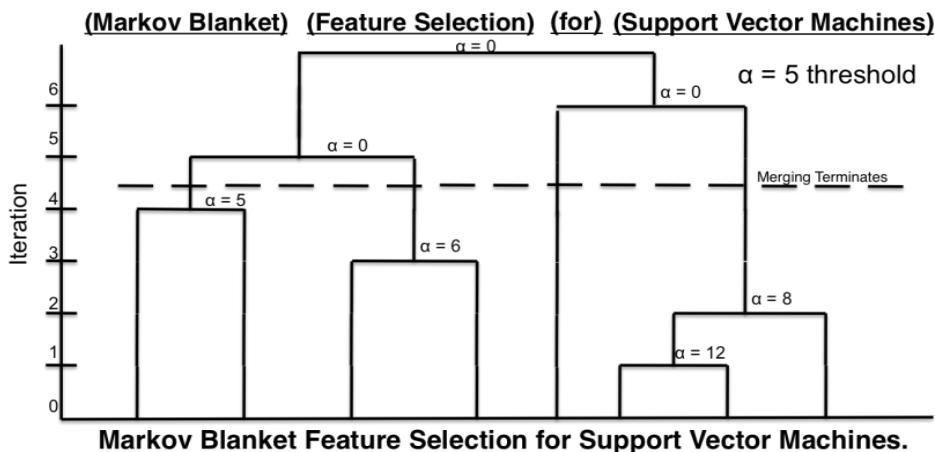
- Strategy 1: Simultaneously inferring phrases and topics
 - Bigram topical model [Wallach'06], topical n-gram model [Wang, et al.'07], phrase discovering topic model [Lindsey, et al.'12]
 - High model complexity: Tends to overfitting; High inference cost: Slow
- Strategy 2: Post topic modeling phrase construction
 - Label topic [Mei et al.'07], TurboTopic [Blei & Lafferty'09], KERT [Danilevsky, et al.'14]
 - Words in the same phrase may be assigned to different topics
 - Ex. knowledge discovery using least squares support vector machine ...
- Our solution 1: ToPMine [El-kishky, et al., VLDB'15]
 - First Phrase Mining then Topic Modeling (No training data at all)
- Our solution 2: SegPhrase+ [Liu, et al., SIGMOD'15]
 - Integrating phrase mining and document segmentation (with minimal training data)



ToPMine: The Overall Phrase Mining Framework

- ❑ ToPMine [El-Kishky et al. VLDB’15]
 - ❑ First phrase construction, then topic mining
 - ❑ Contrast with KERT: First topic modeling, then phrase mining
- ❑ The ToPMine Framework:
 - ❑ Perform **frequent *contiguous pattern*** mining to extract candidate phrases and their counts
 - ❑ Perform agglomerative merging of adjacent unigrams as guided by a significance score — This segments each document into a “***bag-of-phrases***”
 - ❑ The newly formed bag-of-phrases are passed as input to PhraseLDA, an extension of LDA, that constrains all words in a phrase to each sharing the same latent topic

Phrase Mining: Frequent Pattern Mining + Statistical Analysis



Based on significance score [Church et al. '91]:

$$\alpha(P_1, P_2) \approx (f(P_1 \bullet P_2) - \mu_0(P_1, P_2)) / f(P_1 \bullet P_2)^{1/2}$$

- [Markov blanket] [feature selection] for [support vector machines]
- [knowledge discovery] using [least squares] [support vector machine] [classifiers]
- ...[support vector] for [machine learning]...

Phrase	Raw freq.	True freq.
[support vector machine]	90	80
[vector machine]	95	0
[support vector]	100	20



What Kind of Phrases are of “High Quality”?

- ❑ Judging the quality of phrases
 - ❑ Popularity
 - ❑ “information retrieval” vs. “cross-language information retrieval”
 - ❑ Concordance
 - ❑ “powerful tea” vs. “strong tea”
 - ❑ “active learning” vs. “learning classification”
 - ❑ Informativeness
 - ❑ “this paper” (frequent but not discriminative, not informative)
 - ❑ Completeness
 - ❑ “vector machine” vs. “support vector machine”



ToPMine: Experiments on Yelp Reviews

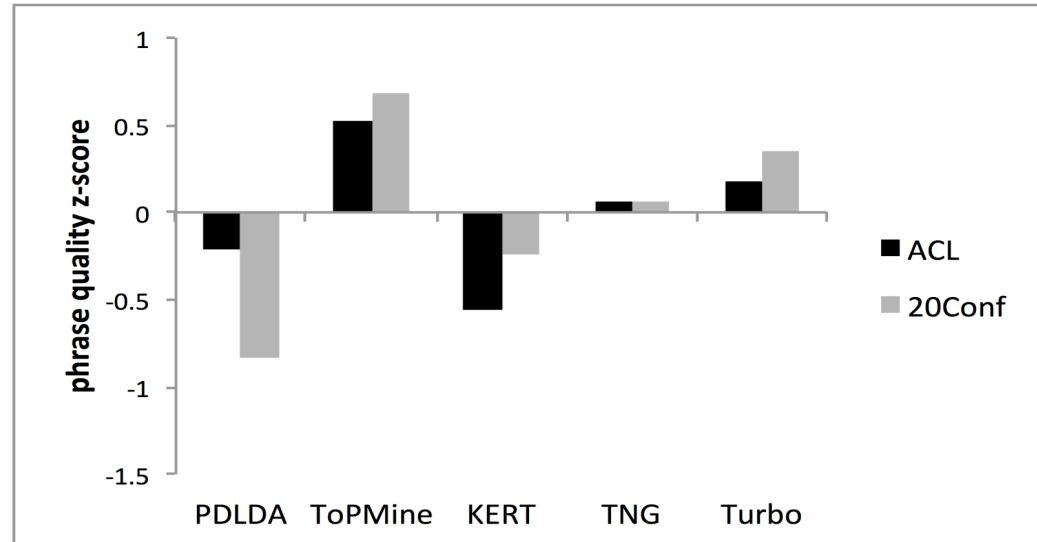
	<i>Topic 1</i>	<i>Topic 2</i>	<i>Topic 3</i>	<i>Topic 4</i>	<i>Topic 5</i>
unigrams	coffee	food	room	store	good
	ice	good	parking	shop	food
	cream	place	hotel	prices	place
	flavor	ordered	stay	find	burger
	egg	chicken	time	place	ordered
	chocolate	roll	nice	buy	fries
	breakfast	sushi	place	selection	chicken
	tea	restaurant	great	items	tacos
	cake	dish	area	love	cheese
	sweet	rice	pool	great	time
n-grams	ice cream	spring rolls	parking lot	grocery store	mexican food
	iced tea	food was good	front desk	great selection	chips and salsa
	french toast	fried rice	spring training	farmer's market	food was good
	hash browns	egg rolls	staying at the hotel	great prices	hot dog
	frozen yogurt	chinese food	dog park	parking lot	rice and beans
	eggs benedict	pad thai	room was clean	wal mart	sweet potato fries
	peanut butter	dim sum	pool area	shopping center	pretty good
	cup of coffee	thai food	great place	great place	carne asada
	iced coffee	pretty good	staff is friendly	prices are reasonable	mac and cheese
	scrambled eggs	lunch specials	free wifi	love this place	fish tacos

ToPMine: Faster and Generating Better Quality Phrases

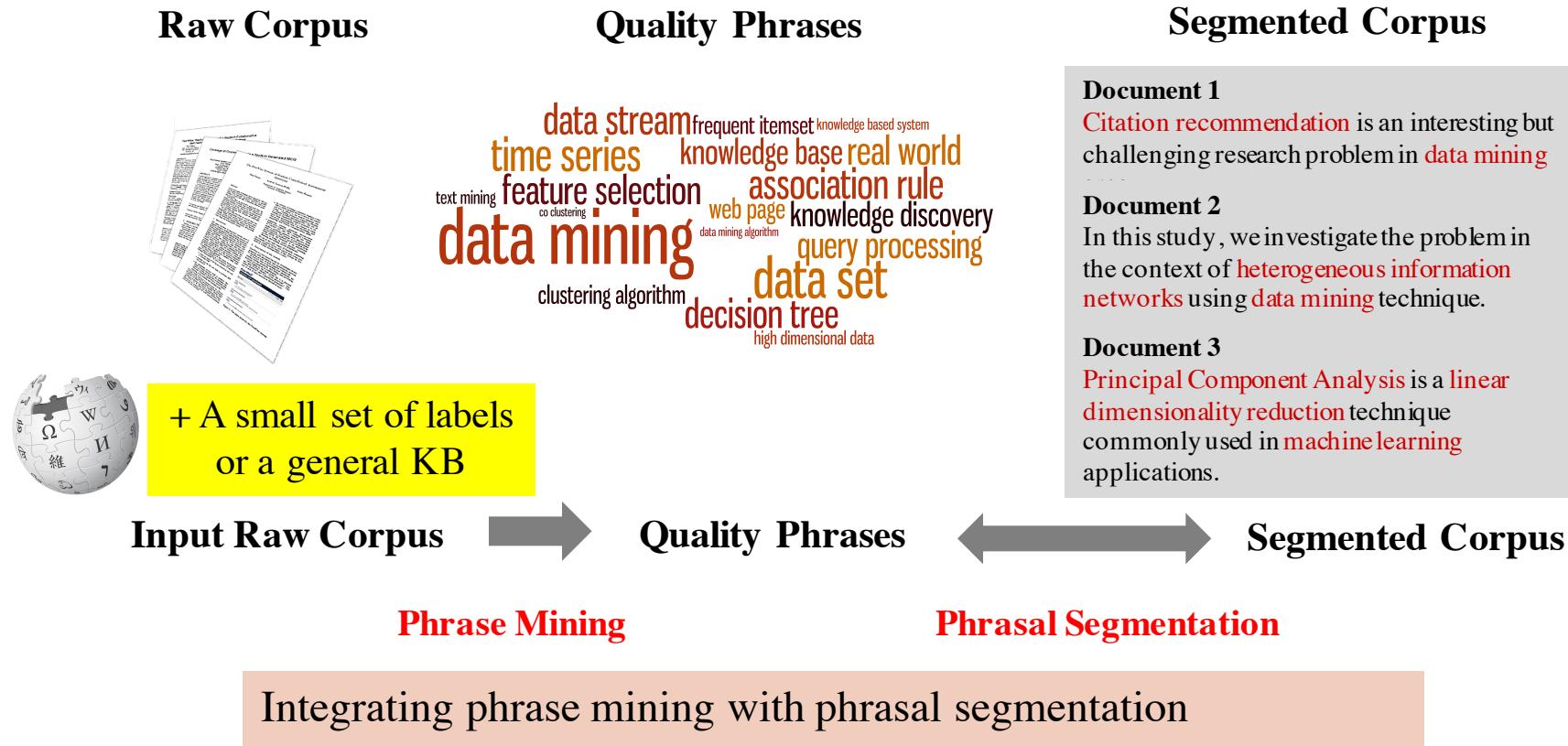
Running time of different algorithms

Method	<i>sam-pled dblp titles (k=5)</i>	<i>dblp titles (k=30)</i>	<i>sampled dblp abstracts</i>	<i>dblp abstracts</i>
PDLDA	3.72(hrs)	~20.44(days)	1.12(days)	~95.9(days)
Turbo Topics	6.68(hrs)	>30(days)*	>10(days)*	>50(days)*
TNG	146(s)	5.57 (hrs)	853(s)	NA†
LDA	65(s)	3.04 (hrs)	353(s)	13.84(hours)
KERT	68(s)	3.08(hrs)	1215(s)	NA†
ToP-Mine	67(s)	2.45(hrs)	340(s)	10.88(hrs)

Phrase quality measured by z-score



SegPhrase: From Raw Corpus to Quality Phrases and Segmented Corpus





Experiments: Interesting Phrases Generated (From the Titles and Abstracts of SIGMOD)

Query	SIGMOD		
Method	SegPhrase+	Chunking (TF-IDF & C-Value)	
1	data base	data base	
2	database system	database system	
3	relational database	query processing	
4	query optimization	query optimization	
5	query processing	relational database	
...	
51	sql server	database technology	
52	relational data	database server	
53	data structure	large volume	
54	join query	performance study	
55	web service	web service	Only in Chunking
...	Only in SegPhrase+	...	
201	high dimensional data	efficient implementation	
202	location based service	sensor network	
203	xml schema	large collection	
204	two phase locking	important issue	
205	deep web	frequent itemset	
...	



Mining Quality Phrases in Multiple Languages

- ❑ Both ToPMine and SegPhrase+ are extensible to mining quality phrases in multiple languages
- ❑ SegPhrase+ on Chinese (From Chinese Wikipedia)
- ❑ ToPMine on Arabic (From Quran Fus7a Arabic)(no preprocessing)
- ❑ Experimental results of Arabic phrases:
اُوْرَفُك → Those who disbelieve
بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِيْمِ → In the name of God the Gracious and Merciful

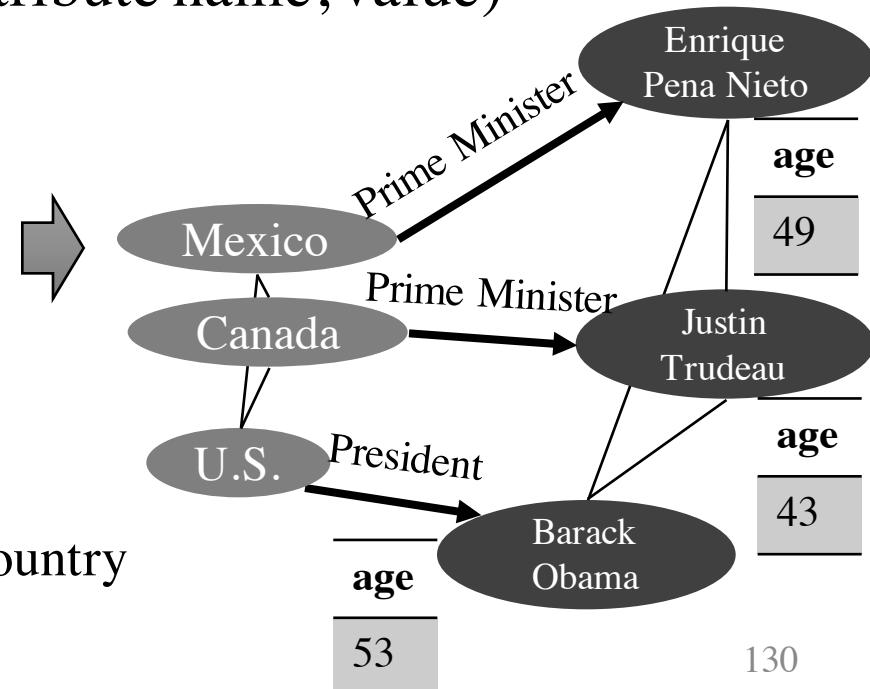
Rank	Phrase	In English
...
62	首席_执行官	CEO
63	中间_偏右	Middle-right
...
84	百度_百科	Baidu Pedia
85	热带_气旋	Tropical cyclone
86	中国科学院_院士	Fellow of Chinese Academy of Sciences
...
1001	十大_中文_金曲	Top-10 Chinese Songs
1002	全球_资讯网	Global Info Website
1003	天一阁_藏_明代_科举_录_选刊	A Chinese book name
...
9934	国家_戏剧_院	National Theater
9935	谢谢_你	Thank you
...

Construction of Heterogeneous Information Networks from Text

💡 Philosophy: Not extensive “labeling” but exploring the power of massive text corpora!

- ❑ Mining phrases (the minimal semantic units)
- ❑ Entity recognition and typing 🔈
- ❑ Attribute discovery (entity, attribute name, value)

...here by **Canada Prime Minister Justin Trudeau, 43**, the so-called #APEChottie...of Mexico's **Enrique Pena Nieto, 49**, ... United States President **Barack Obama, 53**, who...





Why Entity Recognition and Typing from Massive Corpora?

- ❑ Traditional named entity recognition systems are designed for **major types** (e.g., PER, LOC, ORG) and **general domains** (e.g., news)
 - ❑ Require additional steps to adapt to **new domains/types**
 - ❑ Expensive human labor on annotation
 - ❑ 500 documents for entity extraction; 20,000 queries for entity linking
 - ❑ Unsatisfying agreement due to various granularity levels and scopes of types
- ❑ Entities obtained by **entity linking techniques** have *limited coverage* and **freshness**
 - ❑ > 50% unlinkable entity mentions in Web corpus [Lin et al., EMNLP'12]
 - ❑ > 90% in our experiment corpora: tweets, Yelp reviews, ...
- ❑ A new approach: ClusType: Entity Recognition and Typing by Relation Phrase-Based Clustering [Ren, et al., KDD 2015]
 - ❑ Recognizing entity mentions of target types with **minimal/no human supervision** and with **no requirement that entities can be found in a KB** (distant supervision)

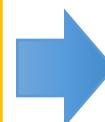
Recognizing Typed Entities

Identifying token span as entity mentions in documents and labeling their types

Target Types

FOOD
LOCATION
JOB_TITLE
EVENT
ORGANIZATION
...

The best BBQ I've tasted in Phoenix! I had the pulled pork sandwich with coleslaw and baked beans for lunch. ... The owner is very nice.

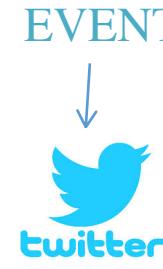
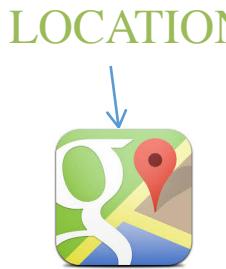
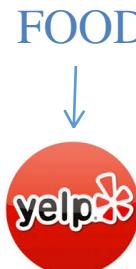
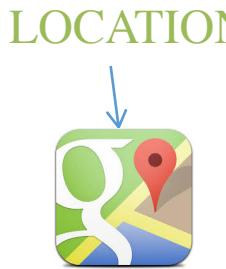
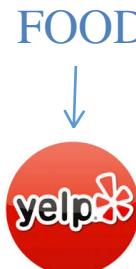


The best **BBQ:Food** I've tasted in **Phoenix:LOC** ! I had the **[pulled pork sandwich]:Food** with **coleslaw:Food** and **[baked beans]:Food** for lunch. ... The **owner:JOB_TITLE** is very nice.

Plain text

Text with typed entities

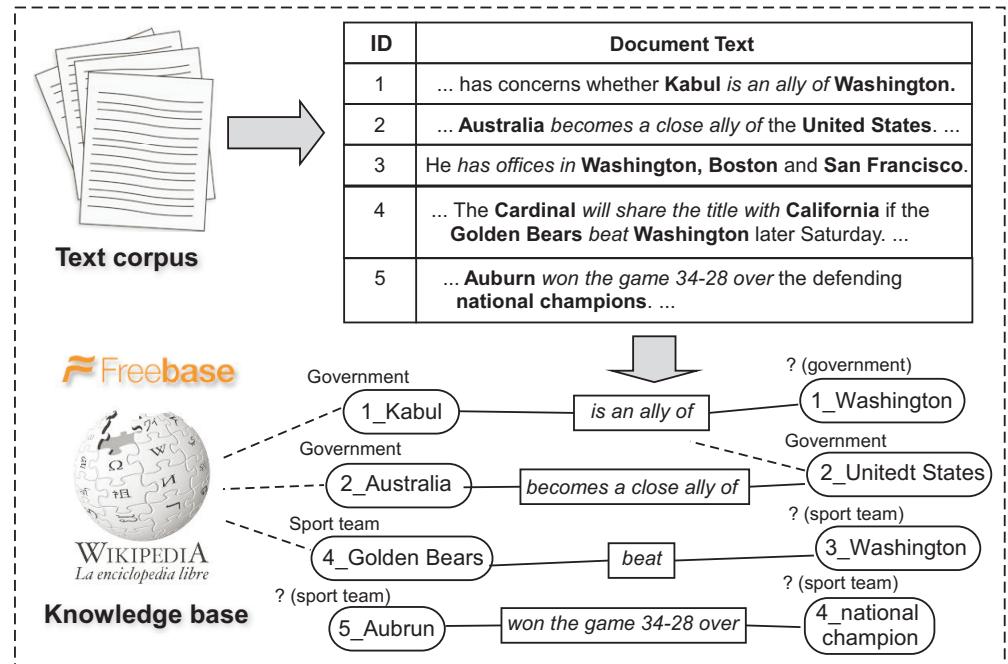
Enabling structured analysis
of unstructured text corpus



ClusType: A Distant Supervision Framework

Problem: *Distantly-supervised entity recognition in a domain-specific corpus*

- ❑ Given: (1) a domain-specific corpus D , (2) a knowledge base (e.g., Freebase), (3) a set of target types (T) from a KB
- ❑ Detect candidate entity mentions in D , and categorize each candidate mention by target types or Not-Of-Interest (NOI)

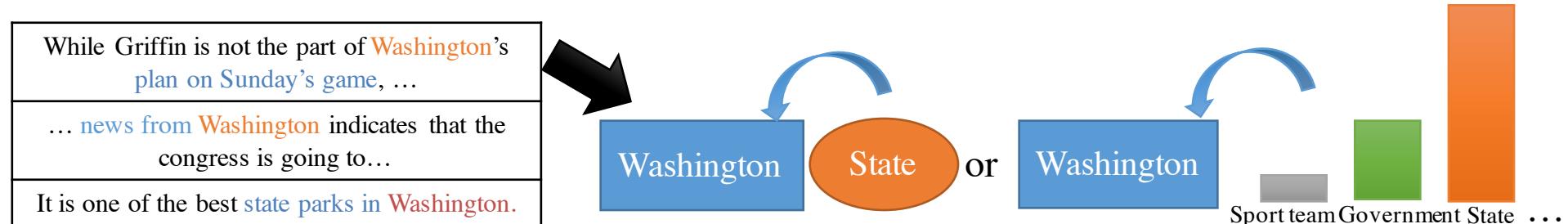


Solution:

- ❑ Detect entity mentions from text
- ❑ Map candidate mentions to KB entities of target types
- ❑ Use confidently mapped {mention, type} to infer types of remaining candidate mentions

Entity Recognition and Typing: Challenges and Solutions

- Challenge 1: Domain Restriction: Extensive training, use general-domain corpora, not work well on **specific, dynamic or emerging domains** (e.g., tweets, Yelp reviews)
 - Solution: Domain-agnostic phrase mining: Extracts candidate entity mentions with **minimal linguistic assumption** (e.g., only use POS tagging)
- Challenge 2: Name ambiguity: Multiple entities may share the same surface name
 - Solution: Model **each mention** based on its **surface name** and **context**



- Challenge 3: Context Sparsity: There are many ways to describe the same relation
 - Solution: cluster **relation phrase**, infer synonymous **relation phrases**

Sentence	Freq.
The magnitude 9.0 quake caused widespread devastation in [Kesennuma city]	12
... tsunami that ravaged [northeastern Japan] last Friday	31
The resulting tsunami devastate [Japan]'s northeast	244

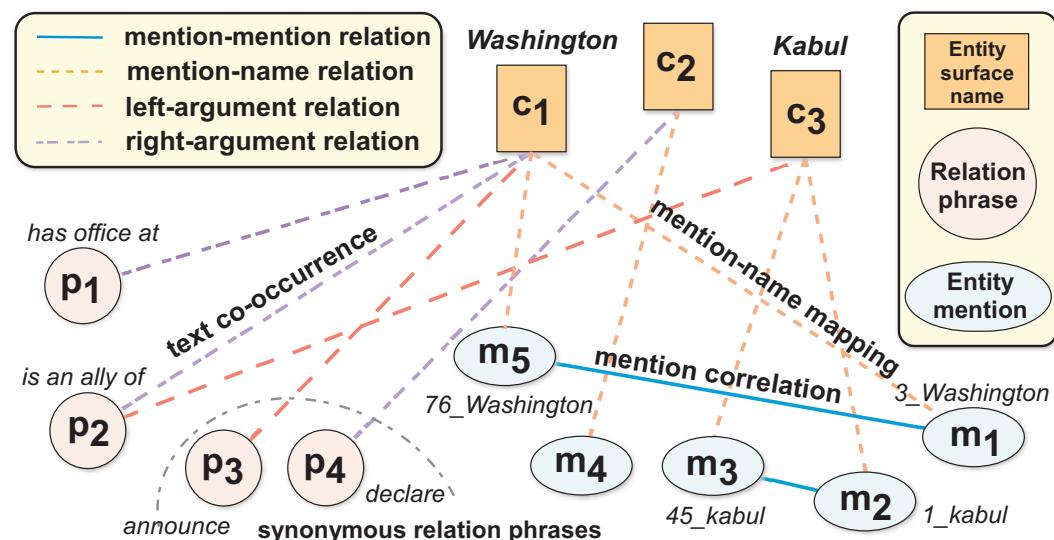
The ClusType Framework: Phrase Segmentation and Heterogeneous Graph Construction

- ❑ POS-constrained phrase segmentation for mining candidate entity mentions and relation phrases, simultaneously
 - ❑ Construct a heterogeneous graph to represent available information in a unified form

Entity mentions are kept as individual objects **to be disambiguated**

Linked to entity surface
names & relation phrases

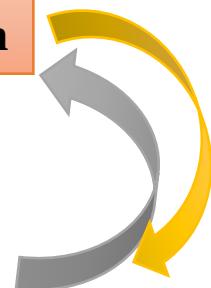
Weight assignment: The more two objects are likely to share the same label, the larger the weight will be associated with their connecting edge



The ClusType Framework: Mutual Enhancement of Type Propagation and Relation Phrase Clustering

- With the constructed graph, formulate a **graph-based semi-supervised learning** of two tasks jointly:

Type propagation on heterogeneous graph



Multi-view relation phrase clustering

Derived entity argument types serve as **good feature** for clustering relation phrases

Propagate type information among entities bridges via synonymous relation phrases

Mutually enhancing each other; leads to quality recognition of unlinkable entity mentions



ClusType: A General Framework Overview

❑ Candidate Generation

- ❑ Perform phrase mining on a POS-tagged corpus to extract candidate entity mentions and relation phrases

❑ Construction of Heterogeneous Graphs

- ❑ Construct a heterogeneous graph to encode our insights on modeling the type for each entity mention
- ❑ Collect seed entity mentions as labels by linking extracted mentions to the KB

❑ Relation Phrase Clustering

- ❑ Estimate type indicator for unlinkable candidate mentions with the proposed type propagation integrated with relation phrase clustering on the constructed graph



Candidate Generation

- ❑ Phrase mining incorporating both *corpus-level statistics* and *syntactic constraints*
 - ❑ **Global significance score:** Filter low-quality candidates; **generic POS tag patterns:** remove phrases with improper syntactic structure
 - ❑ Extend ToPMine to partition corpus into segments which meet both significance threshold and POS patterns → candidate entity mentions & relation phrases

Relation phrase: Phrase that denotes a unary or binary relation in a sentence

Pattern	Example
V	disperse; hit; struck; knock;
P	in; at; of; from; to;
V P	locate in; come from; talk to;
VW*(P)	caused major damage on; come lately

V-verb; P-prep; W-{adv | adj | noun | det | pron}

W* denotes multiple W; (P) denotes optional.

Experiment: Entity detection: Performance comparison between our method and an NP chunker

Method	NYT		Yelp		Tweet	
	Prec	Recall	Prec	Recall	Prec	Recall
Our method	0.469	0.956	0.306	0.849	0.226	0.751
NP chunker	0.220	0.609	0.296	0.247	0.287	0.181

Recall is most critical for this step, since later we cannot detect the misses (i.e., false negatives)

Type Inference: A Joint Optimization Problem

$$\begin{aligned} \mathcal{O}_{\alpha, \gamma, \mu} = & \mathcal{F}(\mathbf{C}, \mathbf{P}_L, \mathbf{P}_R) + \mathcal{L}_\alpha(\mathbf{P}_L, \mathbf{P}_R, \{\mathbf{U}^{(v)}, \mathbf{V}^{(v)}\}, \mathbf{U}^*) \\ & + \Omega_{\gamma, \mu}(\mathbf{Y}, \mathbf{C}, \mathbf{P}_L, \mathbf{P}_R). \end{aligned} \quad (2)$$

$$\begin{aligned} \mathcal{F}(\mathbf{C}, \mathbf{P}_L, \mathbf{P}_R) = & \sum_{i=1}^n \sum_{j=1}^l W_{L,ij} \left\| \frac{\mathbf{C}_i}{\sqrt{D_{L,ii}^{(\mathcal{C})}}} - \frac{\mathbf{P}_{L,j}}{\sqrt{D_{L,jj}^{(\mathcal{P})}}} \right\|_2^2 \\ & + \sum_{i=1}^n \sum_{j=1}^l W_{R,ij} \left\| \frac{\mathbf{C}_i}{\sqrt{D_{R,ii}^{(\mathcal{C})}}} - \frac{\mathbf{P}_{R,j}}{\sqrt{D_{R,jj}^{(\mathcal{P})}}} \right\|_2^2 \end{aligned}$$

Mention modeling & mention correlation

$$\begin{aligned} \Omega_{\gamma, \mu}(\mathbf{Y}, \mathbf{C}, \mathbf{P}_L, \mathbf{P}_R) = & \|\mathbf{Y} - f(\Pi_C \mathbf{C}, \Pi_L \mathbf{P}_L, \Pi_R \mathbf{P}_R)\|_F^2 \\ & + \frac{\gamma}{2} \sum_{c \in \mathcal{C}} \sum_{i,j=1}^{M_c} W_{ij}^{(c)} \left\| \frac{\mathbf{Y}_i}{\sqrt{D_{ii}^{(c)}}} - \frac{\mathbf{Y}_j}{\sqrt{D_{jj}^{(c)}}} \right\|_2^2 + \mu \|\mathbf{Y} - \mathbf{Y}_0\|_F^2 \end{aligned}$$

Type propagation between entity surface names and relation phrases

$$\begin{aligned} \mathcal{L}_\alpha(\mathbf{P}_L, \mathbf{P}_R, \{\mathbf{U}^{(v)}, \mathbf{V}^{(v)}\}, \mathbf{U}^*) = & \sum_{v=0}^d \beta^{(v)} (\|\mathbf{F}^{(v)} - \mathbf{U}^{(v)} \mathbf{V}^{(v)T}\|_F^2 + \alpha \|\mathbf{U}^{(v)} \mathbf{Q}^{(v)} - \mathbf{U}^*\|_F^2). \end{aligned} \quad (3)$$

Multi-view relation phrases clustering



ClusType: Experiment Setting

- ❑ Datasets: 2013 New York Times news (~110k docs) [event, PER, LOC, ORG]; Yelp Reviews (~230k) [Food, Job, ...]; 2011 Tweets (~300k) [event, product, PER, LOC, ...]
- ❑ Seed mention sets: < 7% extracted mentions are mapped to Freebase entities
- ❑ Evaluation sets: manually annotate mentions of target types for subsets of the corpora
- ❑ Evaluation metrics: Follows named entity recognition evaluation (Precision, Recall, F1)
- ❑ Compared methods
 - ❑ **Pattern:** Stanford pattern-based learning; **SemTagger:** bootstrapping method which trains contextual classifier based on seed mentions; **FIGER:** distantly-supervised sequence labeling method trained on Wiki corpus; **NNPLB:** label propagation using ReVerb assertion and seed mention; **APOLLO:** mention-level label propagation using Wiki concepts and KB entities;
 - ❑ **ClusType-NoWm:** ignore mention correlation; **ClusType-NoClus:** conducts only type propagation; **ClusType-TwpStep:** first performs hard clustering then type propagation

Comparing ClusType with Other Methods and Its Variants

Performance comparison on three datasets in terms of Precision, Recall and F1 score

Table 5: Performance comparisons on three datasets in terms of Precision, Recall and F1 score.

Data sets	NYT			Yelp			Tweet		
Method	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Pattern [9]	0.4576	0.2247	0.3014	0.3790	0.1354	0.1996	0.2107	0.2368	0.2230
FIGER [16]	0.8668	0.8964	0.8814	0.5010	0.1237	0.1983	0.7354	0.1951	0.3084
SemTagger [12]	0.8667	0.2658	0.4069	0.3769	0.2440	0.2963	0.4225	0.1632	0.2355
APOLLO [29]	0.9257	0.6972	0.7954	0.3534	0.2366	0.2834	0.1471	0.2635	0.1883
NNPLB [15]	0.7487	0.5538	0.6367	0.4248	0.6397	0.5106	0.3327	0.1951	0.2459
ClusType-NoClus	0.9130	0.8685	0.8902	0.7629	0.7581	0.7605	0.3466	0.4920	0.4067
ClusType-NoWm	0.9244	0.9015	0.9128	0.7812	0.7634	0.7722	0.3539	0.5434	0.4286
ClusType-TwoStep	0.9257	0.9033	0.9143	0.8025	0.7629	0.7821	0.3748	0.5230	0.4367
ClusType	0.9550	0.9243	0.9394	0.8333	0.7849	0.8084	0.3956	0.5230	0.4505

- ❑ vs. **FIGER**: Effectiveness of our candidate generation and type propagation
- ❑ vs. **NNPLB** and **APOLLO**: ClusType utilizes not only semantic-rich relation phrase as type cues, but also cluster synonymous relation phrases to tackle context sparsity
- ❑ vs. our **variants**: (i) models mention correlation for name disambiguation; and (ii) integrates clustering in a mutually enhancing way



Comparing on Trained NER System

- Compare with Stanford NER, which is trained on general-domain corpora including ACE corpus and MUC corpus, on three types: PER, LOC, ORG

F1 score comparison with trained NER

Table 6: F1 score comparison with trained NER.

Method	NYT	Yelp	Tweet
Stanford NER [6]	0.6819	0.2403	0.4383
ClusType-NoClus	0.9031	0.4522	0.4167
ClusType	0.9419	0.5943	0.4717

[6] J. R. Finkel, T. Grenager and C. Manning. Incorporating non-local information into information extraction systems by Gibbs sampling. In ACL'05.

- ClusType and its variants outperform Stanford NER on both dynamic corpus (NYT) and domain-specific corpus (Yelp)
- ClusType has lower precision but higher Recall and F1 score on Tweet → Superior recall of ClusType mainly come from domain-independent candidate generation

Example Output and Relation Phrase Clusters

Example output of ClusType and the compared methods on the Yelp dataset

ClusType	SemTagger	NNPLB
The best BBQ:Food I've tasted in Phoenix:LOC ! I had the [pulled pork sandwich]:Food with coleslaw:Food and [baked beans]:Food for lunch. ...	The best BBQ I've tasted in Phoenix:LOC ! I had the pulled [pork sandwich]:LOC with coleslaw:Food and [baked beans]:LOC for lunch. ...	The best BBQ:Loc I've tasted in Phoenix:LOC ! I had the pulled pork sandwich:Food with coleslaw and baked beans:Food for lunch:Food
I only go to ihop:LOC for pancakes:Food because I don't really like anything else on the menu. Ordered [chocolate chip pancakes]:Food and a [hot chocolate]:Food .	I only go to ihop for pancakes because I don't really like anything else on the menu. Ordered [chocolate chip pancakes]:LOC and a [hot chocolate]:LOC .	I only go to ihop for pancakes because I don't really like anything else on the menu. Ordered chocolate chip pancakes and a hot chocolate .

❑ Extracts more mentions and predicts types with higher

Example relation phrase clusters and corpus-wide frequency from the NYT dataset

ID	Relation phrase
1	recruited by (5.1k); employed by (3.4k); want hire by (264)
2	go against (2.4k); struggling so much against (54); run for re-election against (112); campaigned against (1.3k)
3	looking at ways around (105); pitched around (1.9k); echo around (844); present at (5.5k);

- ❑ Not only synonymous relation phrases, but also both sparse and frequent relation phrase can be clustered together
- ❑ → boosts sparse relation phrases with type information of frequent relation phrases

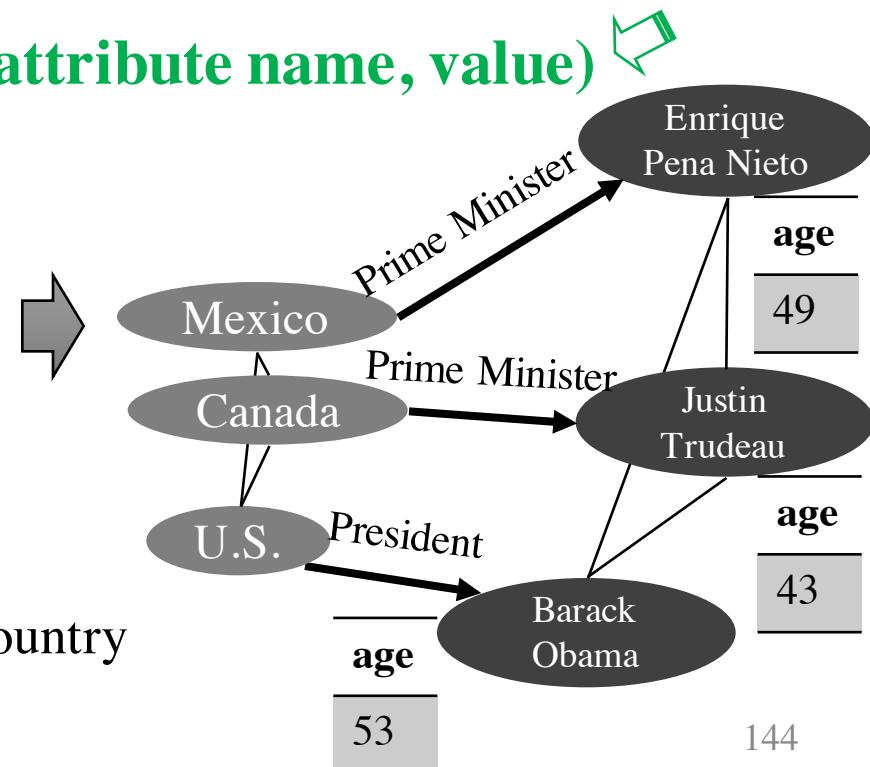
Construction of Heterogeneous Information Networks from Text

💡 Philosophy: Not extensive “labeling” but exploring the power of massive text corpora!

- ❑ Mining phrases (the minimal semantic units)
- ❑ Entity recognition and typing

- ❑ **Attribute discovery (entity, attribute name, value)**

...here by **Canada Prime Minister Justin Trudeau, 43**, the so-called #APEChottie...of Mexico's **Enrique Pena Nieto, 49**, ... United States President **Barack Obama, 53**, who...



Attributed Network Construction

- ❑ Automatic Attribute discovery: Given a class (*e.g.*, \$Country)
 - ❑ Feature as a characteristic (*e.g.*, “population”)
 - ❑ Value: the feature value (*e.g.*, \$Digit or NULL)
 - ❑ Relationship with another class (*e.g.*, “prime minister”)
 - ❑ Value: the other class (*e.g.*, \$Person.Politician.PrimeMinister)
- ❑ Google’s [VLDB’14, WWW’16] based on **fact-seeking** queries
 - ❑ Challenge 1: (Class, Attribute name, **Attribute value**)
 - ❑ Challenge 2: Just text documents (news, tweets, etc.). **NO query**.

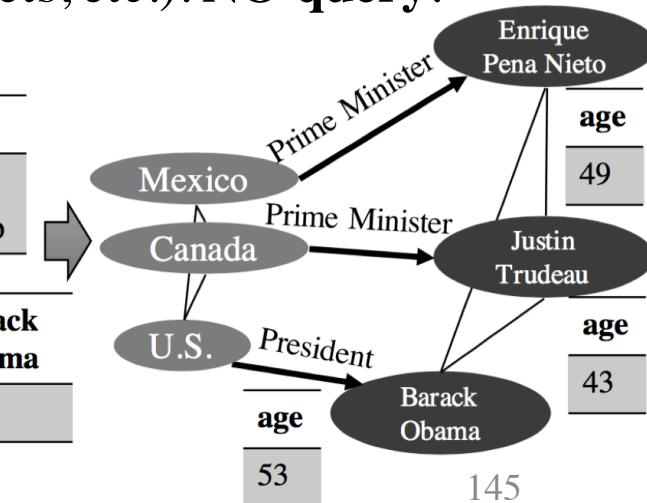
“canada prime minister”, “trudeau age”,
 “united states president”, “obama age”,
 “mexico prime minister” ...

Unfortunately, we don’t have the query data.

...here by Canada Prime Minister Justin Trudeau, 43, the so-called #APEChottie...of Mexico’s Enrique Pena Nieto, 49, ... United States President Barack Obama, 53, who...

Fortunately, we have large text corpus.

		Canada	Mexico
Prime Minister	Justin Trudeau	Enrique Pena Nieto	
	Justin Trudeau	Enrique Pena Nieto	Barack Obama
age	43	49	53





Related Work

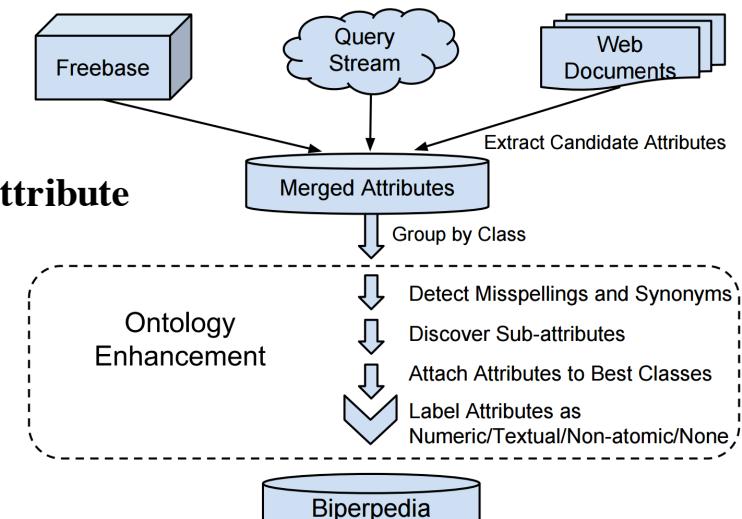
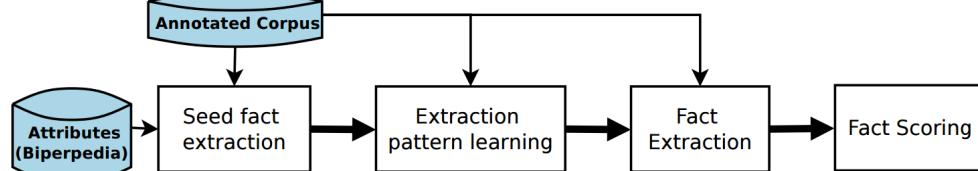
	Data source (✓: available and unlimited with web crawler)	Supervision level (✓: distant or unsupervised for massive data)	Attribute discovery		
			Fine-grained classes	Attribute names	Attribute values
co-EM [6]	✗, product description	?, semi-supervised	✗	✓	✓
SE+R [2]	✗, query log	?, weakly	✗	✓	✗
DvsQ [†] [23]	✗, query log	✓, unsupervised	✓	✓	✗
CAE [†] [22]	✗, query log	✓, unsupervised	✓	✓	✗
WSIEQ [†] [20]	✗, query log	?, weakly	✓	✓	✗
HCAE [†] [21]	✗, query log	?, weakly	✓	✓	✗
WSEST [†] [24]	✗, query log+HTML table	?, weakly	✓	✓	✗
KNEXT [29]	✓, text corpus	✗, supervised	✓	✓	✗
FAR [12]	✗, HTML table	✗, supervised	✓	✓	✓
TYPICALITY [13]	✗, query log+HTML table	✗, supervised	✓	✓	✗
BIPERPEDIA [‡] [7]	✗, query log+HTML table	✓, distant	✓	✓	✗
RENOUN [‡] [30]	✓, text corpus	✗, supervised	✓	✗	✓
ARI [‡] [8]	✗, query log	✓, distant	✓	✓	✗
UPSF [33]	✓, text corpus	✓, unsupervised	✗	✗	✓
MetaPAD	✓, text corpus	✓, distant	✓	✓	✓

[†]These related papers were published by Dr. Marius Pașca et al., Google Inc. from 2007 to 2008.

[‡]These related papers were published by Dr. Alon Halevy et al., Google Inc. from 2014 to 2016.

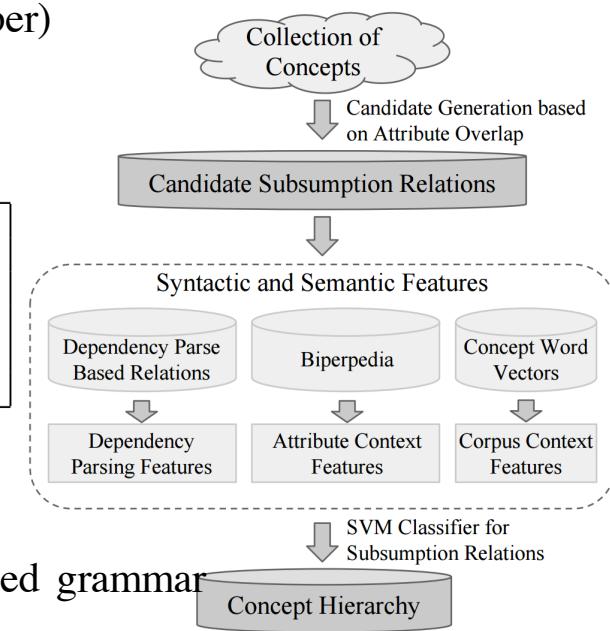
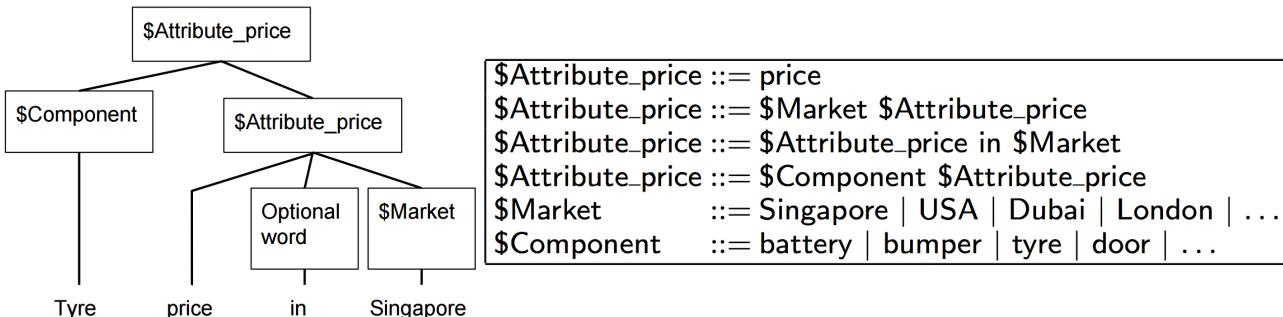
Google's Approaches on Attribute Extraction

- Given Google's **query log**, web text and knowledge bases
 - "Obama wife name" ... "Japan asian population", "Brazil female latino population", "Princeton economist" ...
 - "Obama's wife, Michelle Obama, is a lawyer...", "Princeton economist Paul Krugman was awarded..." ...
 - Obama: \$Person, \$President; Japan, Brazil: \$Location, \$Country; Princeton: \$Organization, \$University...
- Biperpedia (VLDB'14): **Attribute Name Extraction** from query log
 - \$Person: wife name, daughter name
 - \$Country: asian population, female latino population
 - \$University: economist
- ReNoun (EMNLP'14): **Fact Extraction for Noun Phrase Attribute**
 - (Obama, wife, Michelle Obama)
 - (Princeton, economist, Paul Krugman)



Google's Approaches on Attribute Extraction

- Latte (WebDB'15 Best Paper): **Concept (Type) Hierarchy Extraction** with attribute features
 - {country, address, zip code}: \$University (sub) - \$Location (super)
 - {online payment, non profit, tax return}: \$University (sub) - \$Organization (super)
 - {daughter name, wife name, age}: \$President (sub) - \$Person (super)



- ARI (WWW'16): **Attribute Name Structure Extraction** with rule-based grammar
 - Long-tail distribution of attribute names
 - \$Person: \$FamilyMember (name) - daughter, wife, mother, daughter name, wife name
 - \$Country: (\$Gender) (\$Ethnicity) population - asian population, female latino population

Data-Driven: Meta Pattern Mining

- **Meta Pattern:** a sequence of class symbols, words, phrases and punctuation marks that appear contiguously in the text, and serves as a whole semantic unit.

News:

...he's gotten older and grayer, and he's been eclipsed at an Asian economic forum here by **Canada Prime Minister Justin Trudeau, 43**, the so-called #APEChottie... He's also the youngest leader at the Asia Pacific Economic Cooperation forum, six years the junior of **Mexico's Enrique Pena Nieto, 49**, ... **Obama, 53**, who becomes the elder statesman...

- 1. \$Person, \$Digit,
- 2. \$Location.Country Prime_Minister
\$Person.Politician.PrimeMinister

Tweets:

...Protestors march to **Gordon Square** for **12** -year-old **Tamir Rice**...

- 1. protestors march to \$Location.Square
- 2. \$Digit -year-old \$Person.Victim

PubMed abstract:

... Endocarditis caused by **Streptococcus pneumoniae**...
Pericarditis due to **Neisseria meningitidis** ...

- \$Cardiovasular_Diseases caused by \$Bacteria
- \$Cardiovasular_Diseases due to \$Bacteria

MetaPAD Framework

Integrated Data-Driven
Text Mining



Meta Pattern Mining



Attribute Extraction
from Meta Patterns

... Canada Prime Minister Justin Trudeau ...
... Barack Obama , 53, ...

Quality phrase mining (SegPhrase, SIGMOD'15)

... Canada **Prime_Minister Justin_Trudeau** ...
... **Barack_Obama** , 53, ...

Entity recognition and typing with distant
supervision (ClusType, KDD'15)

... **\$Location** Prime_Minister **\$Person** ...
... **\$Person** , **\$Digit** , ...

Fine-grained typing (PLE, KDD' 16)

... **\$Country** Prime_Minister **\$PrimeMinister** ...
... **\$President**, **\$Digit** , ...

MetaPAD Framework

Integrated Data-Driven
Text Mining



Meta Pattern Mining



Attribute Extraction
from Meta Patterns

Quality Meta-Pattern Classifier

Frequency

“prime_minister \$PrimeMinister” vs “young \$PrimeMinister”

Completeness

*“\$Country prime_minister \$PrimeMinister” vs
“\$Country prime_minister”*

Informativeness

*“\$Person ’s brother , \$Person ,” vs “\$Person and
\$Person”*

Coverage

*“\$Person ’s signature healthcare law”: only
“Barack Obama”*

Classifier: Random forest

MetaPAD Framework

Integrated Data-Driven
Text Mining



Meta Pattern Mining



Attribute Extraction
from Meta Patterns

...xxx \$Country Prime_Minister \$PrimeMinister xxx...
...xxx \$President , \$Digit , xxx...

Quality Meta-Pattern Classifier

\$Location Prime_Minister \$Person
\$Person, \$Digit , \$Country Prime_Minister \$PrimeMinister
\$President , \$Digit ,

Synonym Meta-Pattern Detection

- (1) Shared instances
- (2) J.W. similar words

\$Location Prime_Minister \$Person
\$Location PM \$Person
Prime_Minister \$Person of \$Location

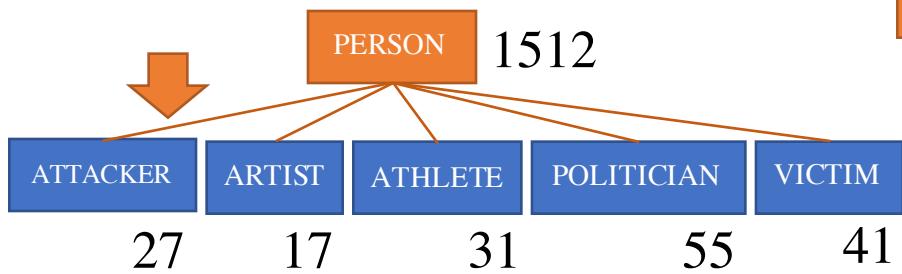
\$Person , \$Digit ,
\$Person , a \$Digit -year-old
\$Person , age \$Digit

Re-typing for Appropriate Granularity

\$Country Prime_Minister \$PrimeMinister
\$Person , \$Digit ,

Top-Down Re-Typing for Granularity

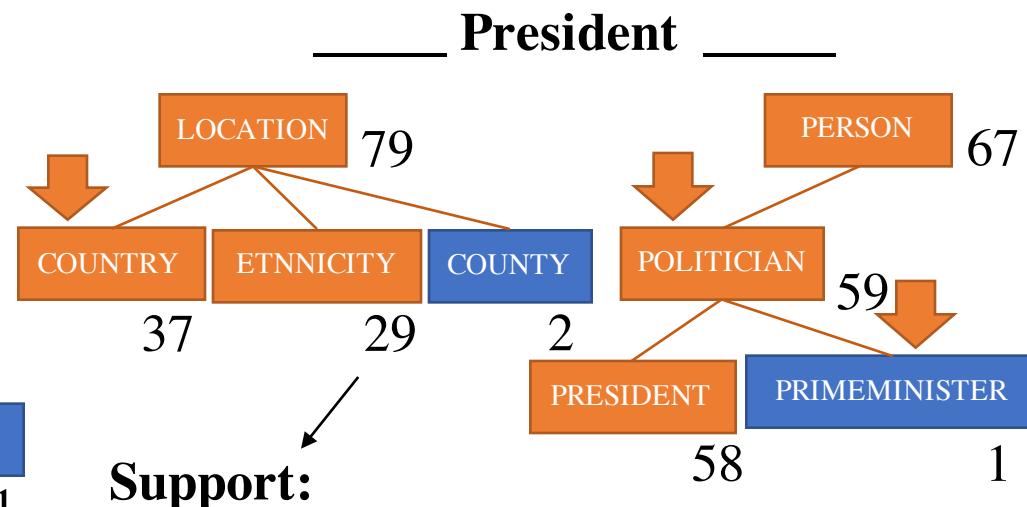
____, a \$Digit -year-old



Graininess:

$$\alpha = (27 + 17 + \dots + 41) / 1512$$

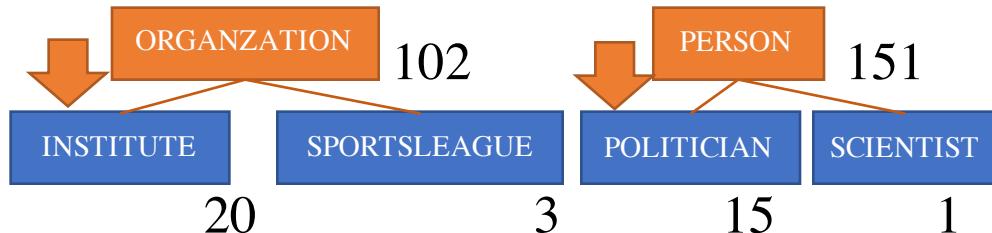
small (< 0.8), stop going down



Support:

$$\beta = 29 / \max(37, 29, 2)$$

big (> 0.1), keep \$Ethnicity



Similar for Bottom-Up...



Experimental Results

Class=\$PERSON (METAPAD: 10,361 names, 4,839 pairs)			Class=\$COUNTRY (METAPAD: 1,132 names, 3,930 pairs)		
Name:BIPERPEDIA	(Name, -)	(Name, Value Type)	Name:BIPERPEDIA	(Name, -)	(Name, Value Type)
Mr.	Dr.	(-year-old,\$DIGIT)	president	president	(ambassador,\$COUNTRY)
Dr.	Mr.	(president,\$ORGANIZATION)	people	government	(president,\$PRESIDENT)
president	president	(spokesman,\$ORGANIZATION)	government	war	(visit,\$PERSON)
wife	director	(director,\$ORGANIZATION)	capital	border	(dead,\$DIGIT)
-year-old	spokesman	(wife,\$PERSON)	visit	volcano	(prime minister,\$PRIMEMINISTER)
death	chief	(chairman,\$ORGANIZATION)	economy	sanctions	(senator,\$SENATOR)
coach	professor	(governor,\$USSSTATE)	prime minister	ambassador	(embassy,\$COUNTRY)
love	head	(spokeswoman,\$ORGANIZATION)	part	earthquake	(condemn,\$ORGANIZATION)
son	coach	(leader,\$ORGANIZATION)	leaders	capital	(district judge,\$PERSON)
...
code case homicide	staff sergeant	(told reporters,\$WEEKDAY)	nuclear dossier	volcano eruption	(protests,\$NEWSAGENCY)
snow pants	army chief	(board member,\$ORGANIZATION)	similar box	security	(-magnitude earthquake,\$DIGIT)
fellow director	basketball coach	(hack,\$COMPANY)	episcopal oversight	parliament	(second biggest,\$ORGANIZATION)
Class=\$INSTITUTE (METAPAD: 402 names, 198 pairs)			Class=\$BASKETBALLPLAYER (METAPAD: 58 names, 40 pairs)		
Name:BIPERPEDIA	(Name, -)	(Name, Value Type)	Name:BIPERPEDIA	(Name, -)	(Name, Value Type)
professor	professor	(professor,\$PERSON)	guard	forward	(points,\$DIGIT)
students	students	(law professor,\$PERSON)	star	points guard	(center,\$TEAMNAME)
president	graduate	(political science professor,\$PERSON)	game	game	(freshman,\$SPORTSLEAGUE)
campus	law professor	(student,\$PERSON)	forward	freshman	(forward,\$TEAMNAME)
law professor	campus	(grad,\$PERSON)	career	center	(point guard,\$TEAMNAME)
graduate	degree	(signee,\$PERSON)	teammate	get better	(all-star,\$SPORTSLEAGUE)
director	dean	(economics professor,\$PERSON)	point guard	basketball player	(games,\$DIGIT)
study	faculty	(basketball coach,\$PERSON)	points	full highlights	(rebounds,\$DIGIT)
researchers	expert	(finance professor,\$PERSON)	season	jumper	(ast,\$DIGIT)
...
foul	commitment	(class,\$YEAR)	understudy	retirement	(PG,\$TEAMNAME)
socialism speech	dorm	(superintendent,\$PERSON)	birthday boy	shoes	(career earnings,\$DIGIT \$DIGITUNIT)
good summary	program	(-year-old student,\$DIGIT \$PERSON)	injury meme	suspended without pay	(sue,\$PERSON)



Experimental Results

Class=\$LOCATION; Value Type=\$MONTH,\$DAY,\$YEAR			Class=\$ORGANIZATION; Name="ceo"		
#	Meta Patterns	#	Meta Patterns		
1	\$LOCATION \$MONTH \$DAY, \$YEAR	1	\$ORGANIZATION CEO \$PERSON		
2	\$COUNTRY, \$WEEKDAY, \$MONTH \$DAY, \$YEAR	2	\$COMPANY CEO \$BUSINESSPERSON		
3	\$LOCATION on \$MONTH \$DAY, \$YEAR	3	\$ORGANIZATION's \$PERSON		
#	Entity	Attribute Value	#	Entity	Attribute Value
1	Pearl Harbor	December 7, 1941	1	Apple	Tim Cook
2	Green Bay	Sunday, Jan 11, 2015	2	Facebook	Mark Zuckerberg
3	Malta ¹	Friday, Nov 27, 2015	3	Hewlett-Packard	Carly Fiorina
...
5862	Beijing ²	October 11, 2013	765	Boston Medical Center	Kate Walsh
5863	Finland ³	April 8, 2015	766	Association of Private Sector Colleges and Universities	Steve Gunderson
Class=\$PERSON; Name="-year-old" ⁷			Class=\$PERSON; Name="president"; Value Type=\$ETHNICITY		
#	Meta Patterns	#	Meta Patterns		
1	\$DIGIT-year-old \$PERSON	1	\$ETHNICITY President \$PRESIDENT		
2	\$PERSON, \$DIGIT,	2	\$ETHNICITY leader \$PRESIDENT		
3	\$PERSON, a \$DIGIT-year-old	3	\$ETHNICITY government of President \$PRESIDENT		
#	Entity	Attribute Value	#	Entity	Attribute Value
1	Tamir Rice	12	1	Vladimir Putin	Russian
2	Bobbi Kristina Brown	21	2	Francois Hollande	French
3	Michael Brown	18	3	Raul Castro	Cuban
...
4993	Jay Nixon	58	254	Mohammed Morsi	Egyptian
4994	Xanana Gusmao	68	255	Klaus Iohannis	Romanian

¹Commonwealth Heads of Government Meeting. ²UCI World Tour of Beijing. ³Finnish parliamentary election begins.



Experimental Results

F1 score	WPB ('10, 100M)	CNA ('97-'10, 200M)	APR ('15, 200M)	TWT ('15, 1GB)
Total (vs Biperpedia -q)	↑67.7%	↑48.3%	↑189.5%	↑208.0%
w/ Meta pattern classifier	↑30.1%	↑27.0%	↑127.1%	↑195.6%
w/ Granularity	↑20.8%	↑15.6%	↑17.3%	↑3.1%
w/ Integrated text mining techs	↑13.8%	↑9.3%	↑13.0%	↑0.8%

\$Cardiovasular_Diseases due to \$Bacteria

\$Cardiovasular_Diseases caused by \$Bacteria

\$Bacteria	\$Cardiovascular_Diseases
Streptococcus pneumoniae	Endocarditis
Neisseria meningitidis	Pericarditis
Haemophilus paraphrophilus	Endocarditis
Proteus	Endocarditis
Listeria monocytogenes	Pericarditis
Corynebacterium	Endocarditis
Actinomyces	Endocarditis
Coxiella	Endocarditis
Pasteurella pneumotropica	Endocarditis
Cardiobacterium	Endocarditis

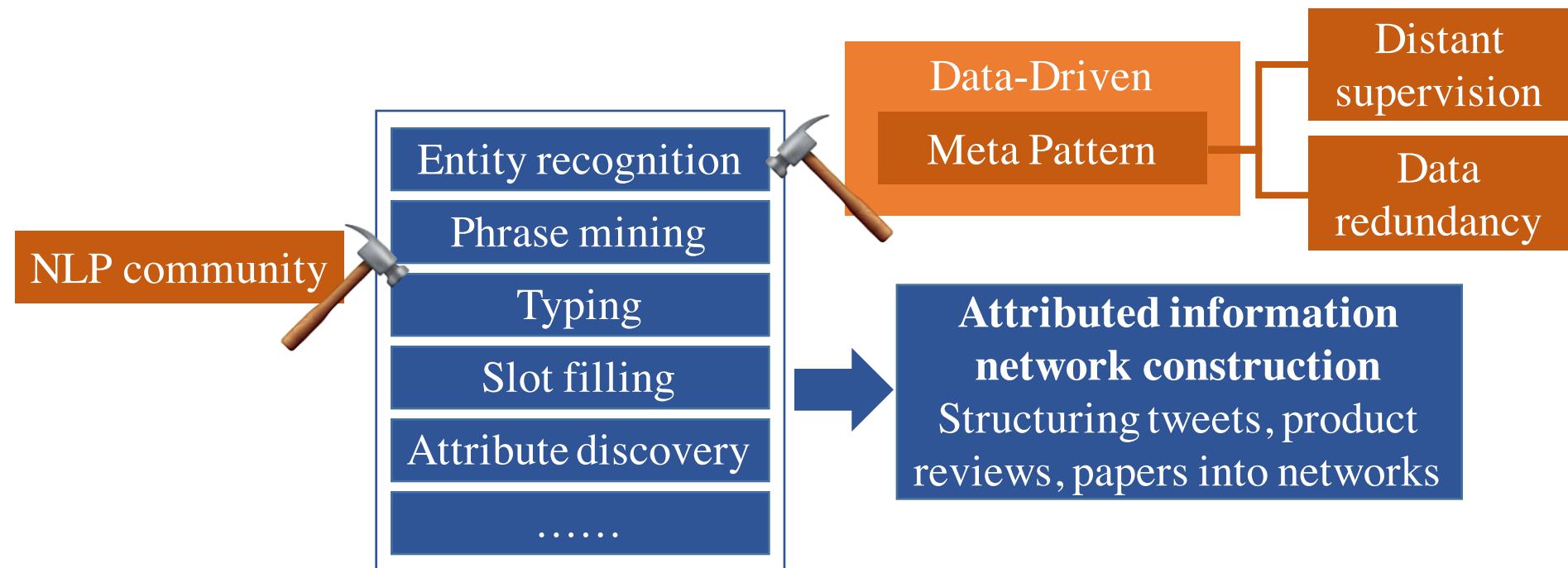
\$Enzymes_and_Coenzymes inhibitor \$Chemical

\$Chemical	\$Enzymes_and_Coenzymes
chelerythrine	protein kinase C
fondaparinux	Factor Xa
calphostin C	protein kinase C
bisindolylmaleimide	protein kinase C

\$Diagnosis : \$Digit +/- \$Digit kg/m (\$Digit)

\$Diagnosis	\$Digit \$Digit \$Digit
BMI	(31.0 , 6.4 , 2)
BMI	(26 , 4 , 2)
body mass index	(27 , 6 , 2)

Meta Pattern: Data-Driven Approaches for NLP Tasks

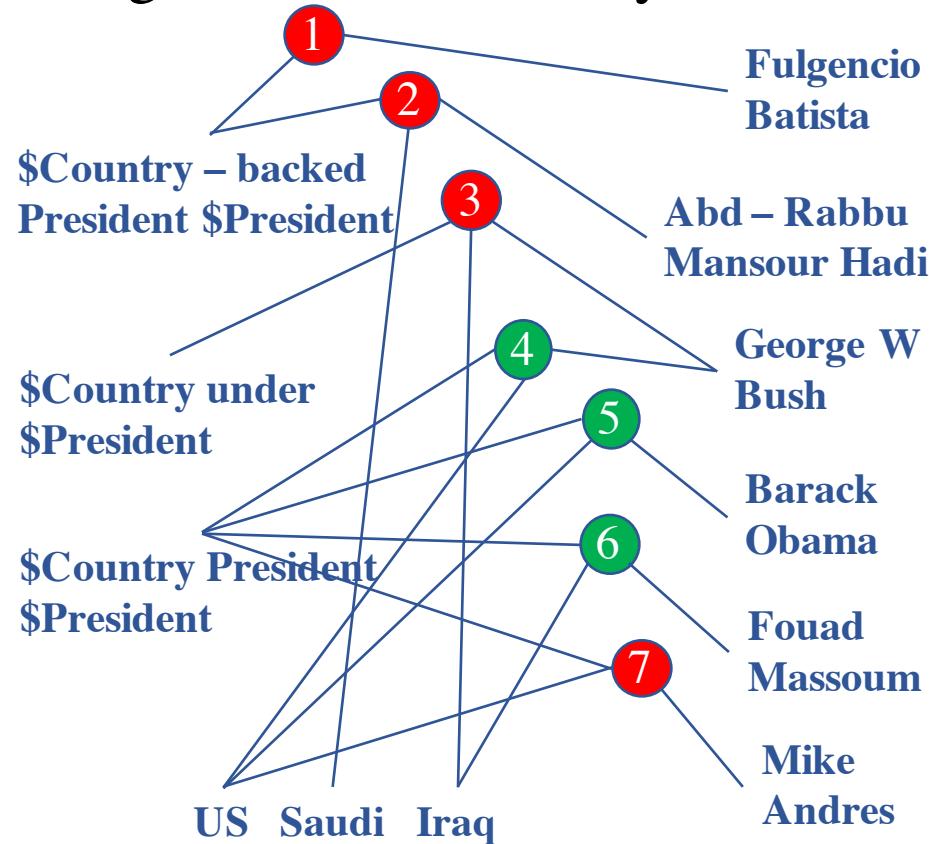


Finding Truth when Structuring

- Find trustworthy facts by modeling “source” reliability

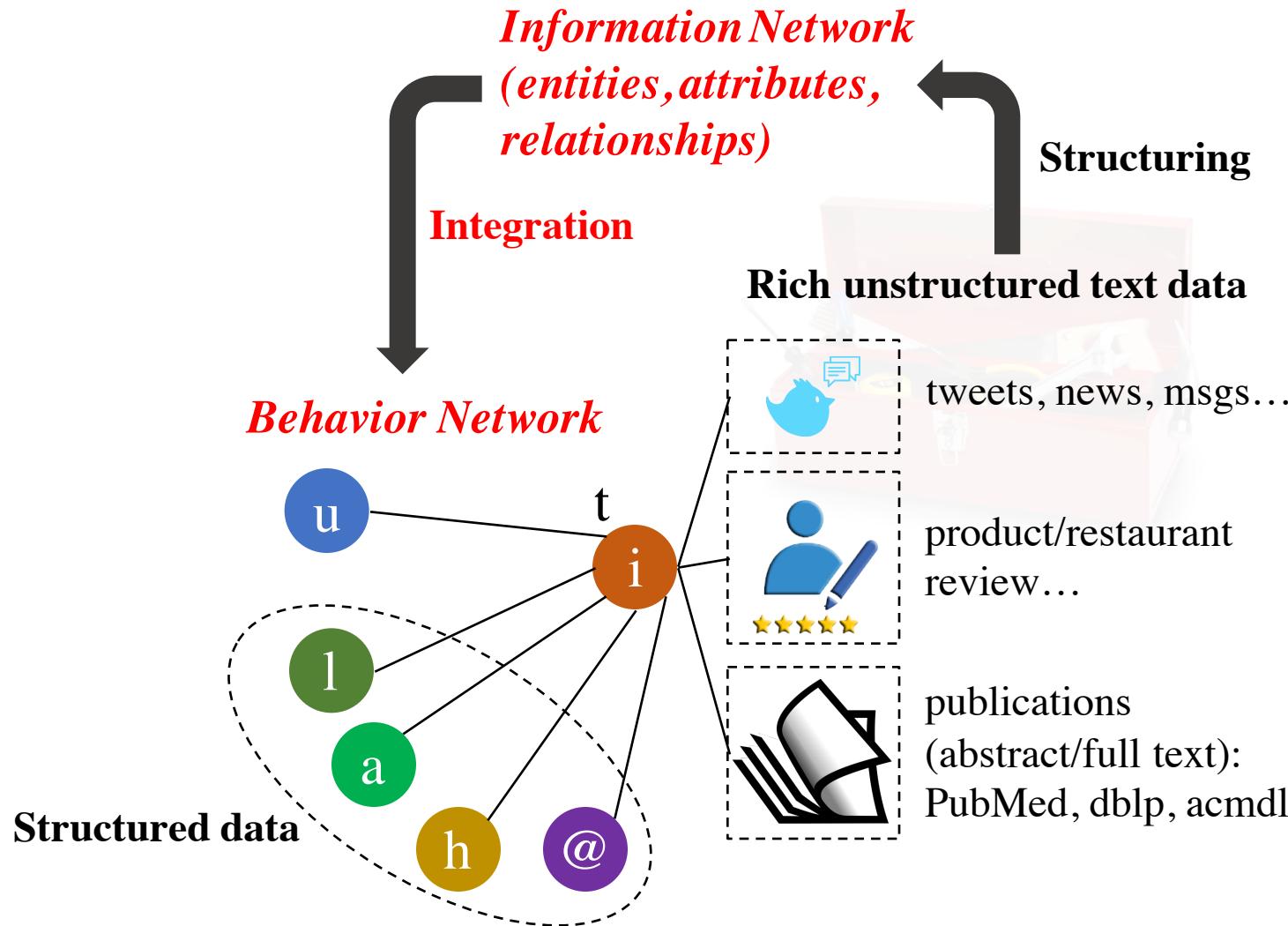
(\$Country, president, \$President) : 0.829

Meta Pattern	Acc. (FP/P)
\$Country 's President \$President	0.984 (1/61)
President \$President of \$Country	1.000 (0/24)
\$Country 's President \$President ,	1.000 (0/16)
” \$Country President \$President	1.000 (0/7)
...	...
President \$President said \$Country	0.833 (1/6)
\$Country President \$President	0.807 (16/83)
\$Country , President \$President	0.650 (7/20)
\$Country - backed President \$President	0.500 (3/6)
\$Country under \$President	0.500 (1/2)



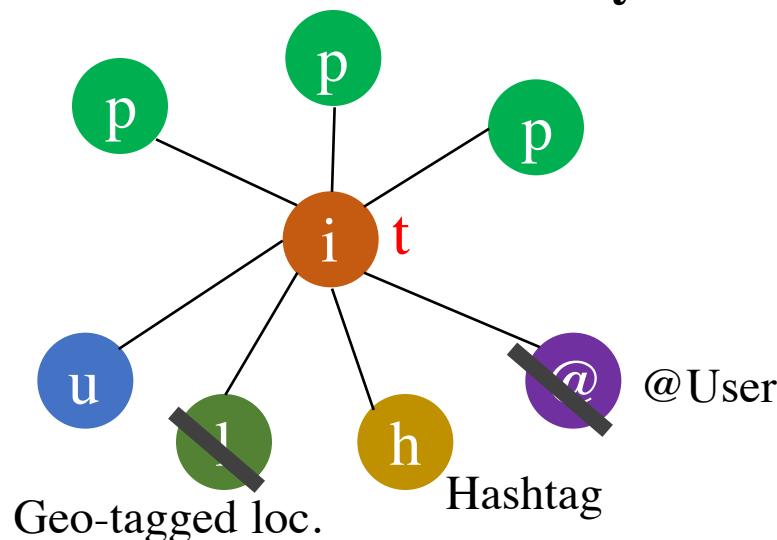
- ... Fidel Castro and his brother Raul led winning a revolution toppling **US - backed President Fulgencio Batista** .
- ... control of the country and at reinstating **Saudi - backed President Abd - Rabbu Mansour Hadi** .
- ... was profoundly forward - leaning and outspoken about the importance of invading **Iraq under George W Bush** .
- ... better delivering on those expectations , " McDonald 's **US President Mike Andres** said in the announcement

Data to Network to Knowledge



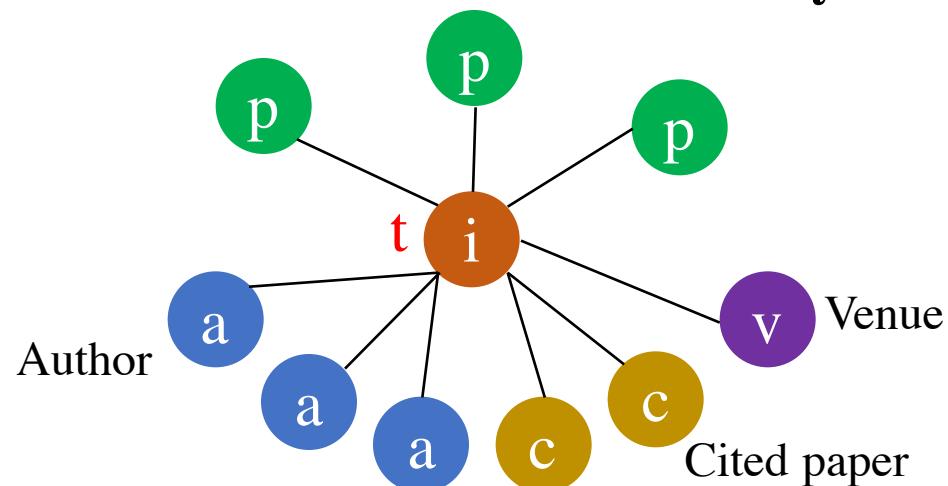
Bring Phrases to Behavior Modeling

- ❑ Tweeting behavior
 - ❑ Event **summary**



20:03:09 @ebekahwsm
this better be the best halftime show ever
in the history of halftimes shows. ever.
#SuperBowl

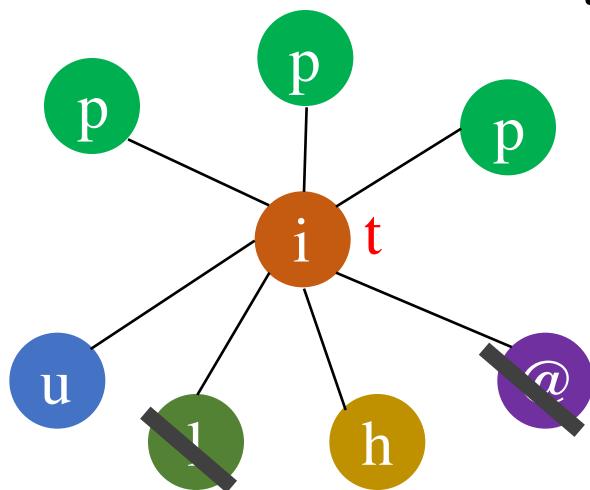
- ❑ Paper-publishing behavior
 - ❑ Research trend **summary**



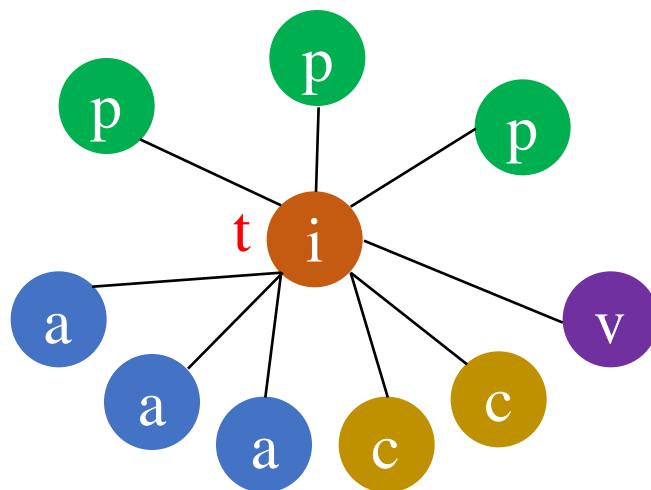
2009 P. Melville, W. Gryc, R. Lawrence,
“Sentiment analysis of blogs by combining
lexical knowledge with text classification”,
KDD’09. Refs: p81623, p84395...

Tensor Fails

- ❑ Tweeting behavior
 - ❑ Event **summary**

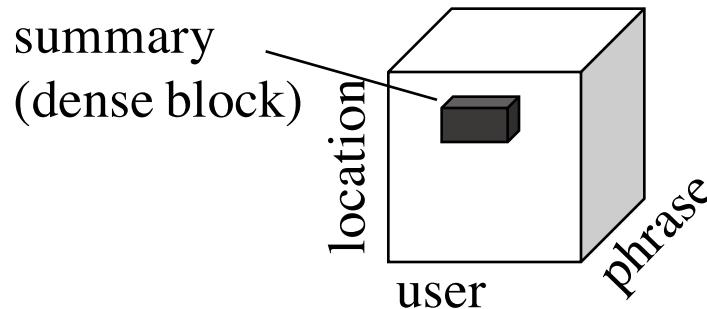


- ❑ Paper-publishing behavior
 - ❑ Research trend **summary**



Q: How to represent and summarize **dynamic multi-contextual** behaviors?

A set of values in dimensions (*one-guaranteed value, empty value, multi-values*)



Two-Level Matrix and “Tartan”

Time slice t

User-Phrase-URL” Tartan
(Advertising campaign)

Multicontextual →
(dimensions,
dimensional values)

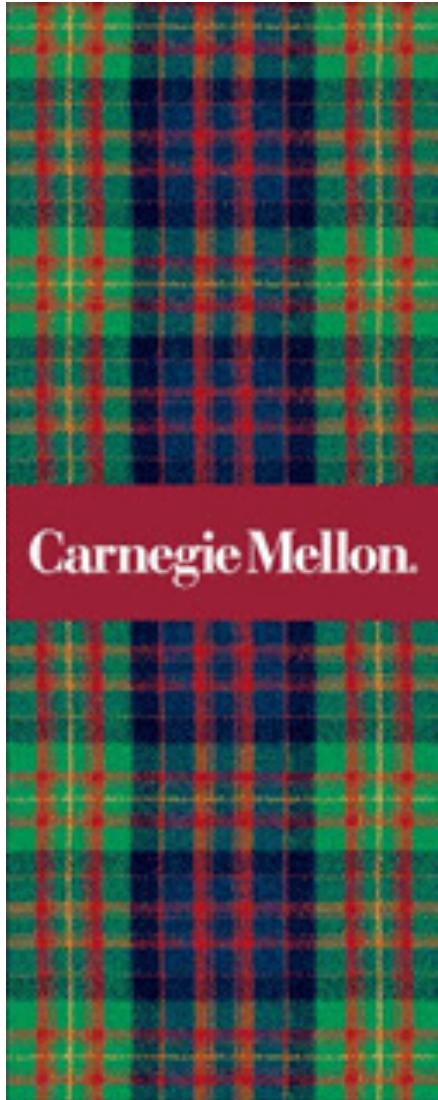
or
(ing)

Dynamic
(consecutive
time slices)

“Phrase-Location-Hashtag” Tartan
(Local event)

162

CMU Tartans



Optimize with MDL Principle

- Maximize the number of bits by encoding the Tartan

$$f(\mathcal{A}, \mathcal{X}) = L(\mathcal{X}^{\mathcal{A}}) - L(\mathcal{A}) - L(\mathcal{X}^{\mathcal{A}} \setminus \mathcal{A}).$$

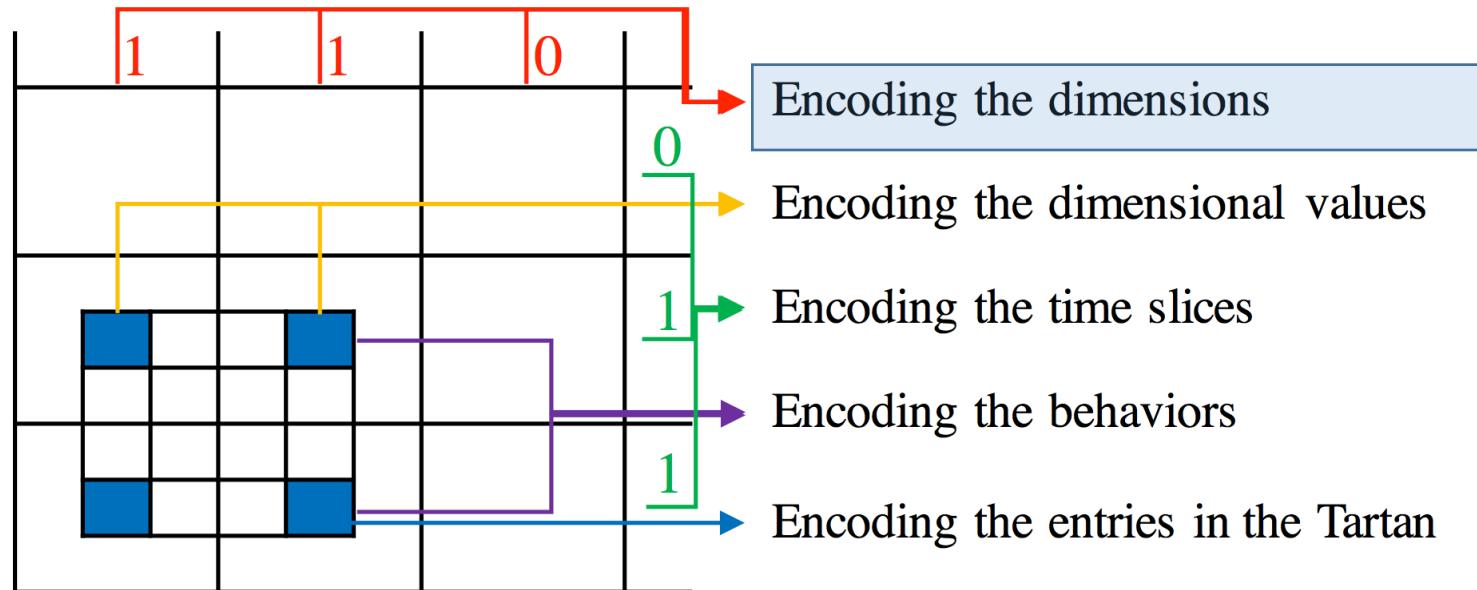
User	Phrase	URL	Loc.	Hashtag	
...	
1 1	1 1 1 2	1 1	1 1	1 1	
...	
Time slice t	“User-Phrase-URL” Tartan (Adver)				
...	1 ... 1 1 ... 1	1 1	1 1	1 1	
...	
Behavior (tweeting)	2 0 1 1	1 1	1 1	1 1	
...	
t+1	1 ... 1 1 ... 1	1 1	1 1	1 1	
t+2	2 2 1 1	1 1	1 1	1 1	‘Phrase-Location-Hashtag’ Tartan (Local event)

$L(\mathcal{X}^{\mathcal{A}}) = g(V + C, C) + L_{\mathcal{D}}(\mathcal{A}) + L_{\mathcal{T}}(\mathcal{A}) + \sum_{d \in \mathcal{D}} \log^* N_d + \sum_{t \in \mathcal{T}} \log^* E^{(t)}.$

$L(\mathcal{A}) = L_{\mathcal{D}}(\mathcal{A}) + L_{\mathcal{V}}(\mathcal{A}) + L_{\mathcal{T}}(\mathcal{A}) + L_{\mathcal{B}}(\mathcal{A}) + L_{\mathcal{A}}(\mathcal{A}).$

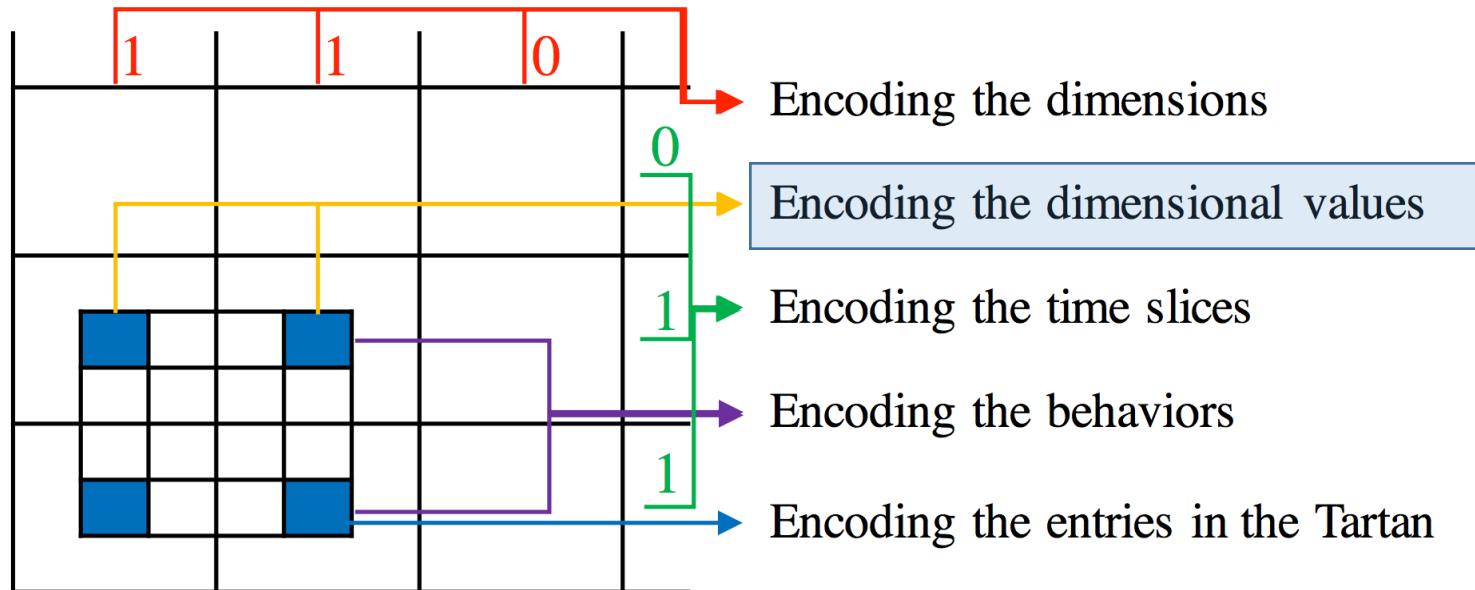
$L(\mathcal{X}^{\mathcal{A}} \setminus \mathcal{A}) = g(V + C - v - c, C - c);$

Encoding Tartan: Dimensions



$$\begin{aligned}
 H_{\mathcal{D}}(X) &= - \sum_{x \in \{0,1\}} P(X = x) \log P(X = x) \\
 &= - \left(\frac{D^{\mathcal{A}}}{D} \log \frac{D^{\mathcal{A}}}{D} + \frac{D - D^{\mathcal{A}}}{D} \log \frac{D - D^{\mathcal{A}}}{D} \right). \\
 L_{\mathcal{D}}(\mathcal{A}) &= \log^* D + \log^* D^{\mathcal{A}} + D \cdot H_{\mathcal{D}}(X) \\
 &= \log^* D + \log^* D^{\mathcal{A}} + g(D, D^{\mathcal{A}}),
 \end{aligned}$$

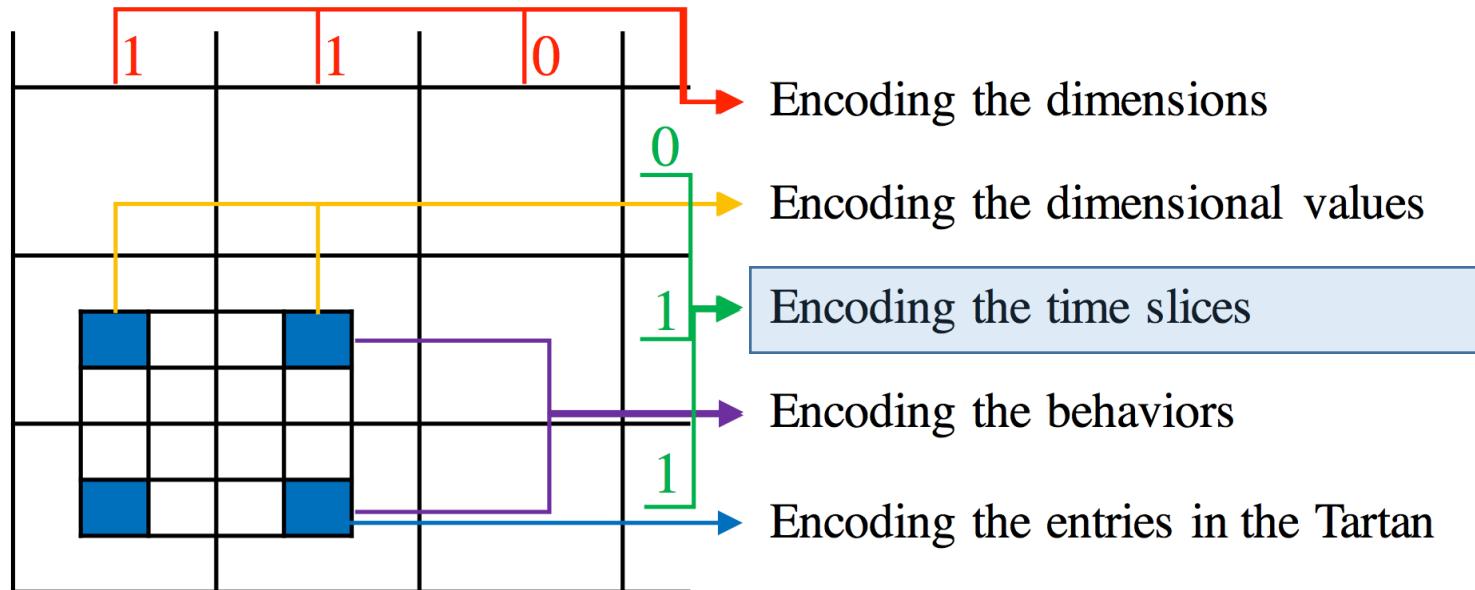
Encoding Tartan: Dimensional Values



$$H_{\mathcal{V}_d}(X) = - \left(\frac{n_d}{N_d} \log \frac{n_d}{N_d} + \frac{N_d - n_d}{N_d} \log \frac{N_d - n_d}{N_d} \right).$$

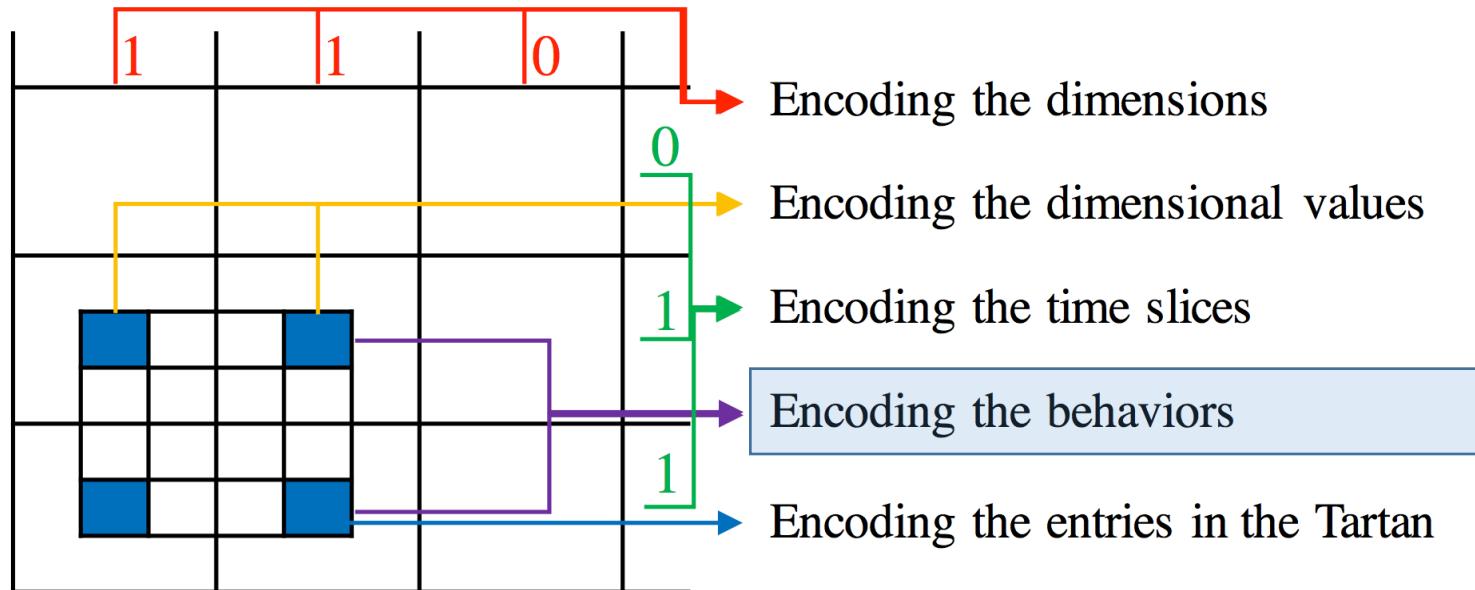
$$L_{\mathcal{V}}(\mathcal{A}) = \sum_{d \in \mathcal{D}} \left(\log^* N_d + \log^* n_d + g(N_d, n_d) \right).$$

Encoding Tartan: Time Slices



$$L_{\mathcal{T}}(\mathcal{A}) = \log^* T + \log^* T^{\mathcal{A}} + \log^* t_{start}$$

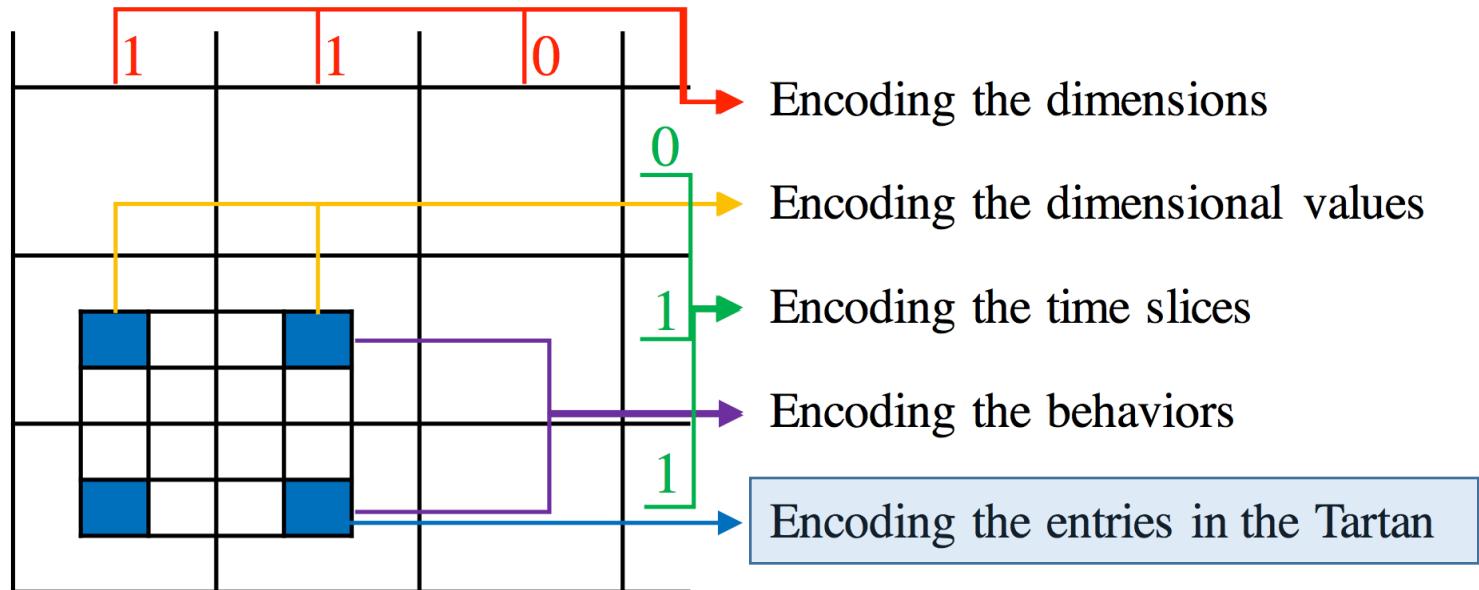
Encoding Tartan: Behaviors



$$H_{\mathcal{B}^{(t)}}(X) = - \left(\frac{e^{(t)}}{E^{(t)}} \log \frac{e^{(t)}}{E^{(t)}} + \frac{E^{(t)} - e^{(t)}}{E^{(t)}} \log \frac{E^{(t)} - e^{(t)}}{E^{(t)}} \right).$$

$$L_{\mathcal{B}}(\mathcal{A}) = \sum_{t \in \mathcal{T}} \left(\log^* E^{(t)} + \log^* e^{(t)} + g(E^{(t)}, e^{(t)}) \right).$$

Encoding Tartan: Entries



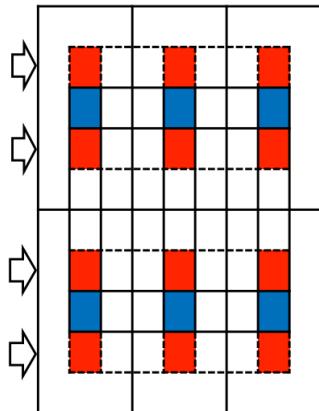
$$v = \left(\sum_{d \in \mathcal{D}} n_d \right) \left(\sum_{t \in \mathcal{T}} e^{(t)} \right).$$

$$c = \sum_{d \in \mathcal{D}, t \in \mathcal{T}} \sum_{b \in \mathcal{B}^{(t)}, i \in \mathcal{V}_d} \chi_d^{(t)}(b, i).$$

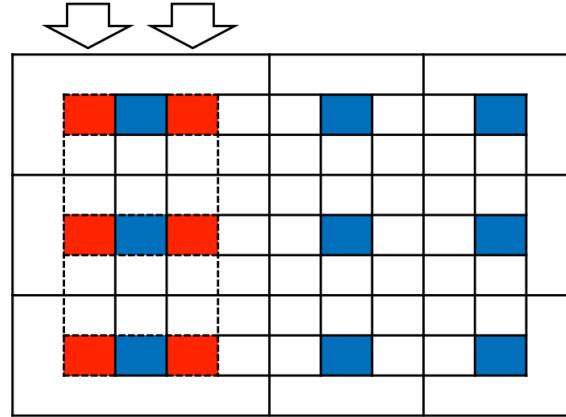
$$H_{\mathcal{A}}(X) = -\left(\frac{c}{v+c} \log \frac{c}{v+c} + \frac{v}{v+c} \log \frac{v}{v+c} \right).$$

$$L_{\mathcal{A}}(\mathcal{A}) = (v + c) H_{\mathcal{A}}(X) = g(v + c, c).$$

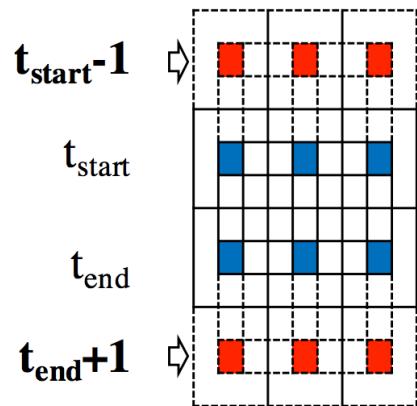
Greedy Search for the Local Optimum



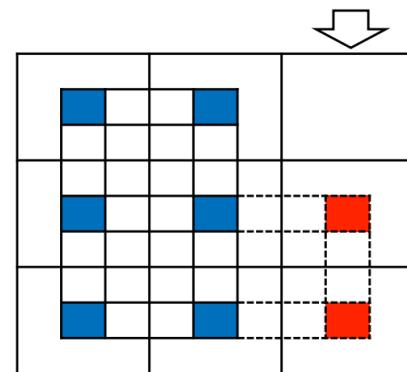
(a) Update the set of behaviors.



(b) Update the set of values.



(c) Update the consecutive time slices.



(d) Update the set of dimensions.

Time complexity:

$$\mathcal{O}(\sum_d N_d \log N_d + \sum_t E^{(t)} \log E^{(t)})$$



Experimental Results

□ DM/ML research trend summaries with DBLP data

Author	Venue	Keyword	Cited	#Paper	Venue	Keyword	#Paper
76 Cheng-xiang Zhai Hui Fang S. Kambhampati	7 SIGIR VLDB TKDE	7 “information retrieval” “data integration” “text classification”	68 p56743 ¹ p62995 p76869	32 2003- 2007	5 ICML NIPS ...	6 “reinforcement learning” “machine learning”	40 1997- 2002

¹ “A language modeling approach to information retrieval”

Author	Venue	Cited	#Paper	Venue	Keyword	#Paper	Author	Venue	Keyword	#Paper
6 Jiawei Han Xifeng Yan	1 SIG- MOD	1 p76095 ²	22 2004- 2010	3 ICDM AAAI TKDE	1 “anomaly detection”	25 2005- 2013	27 C. Faloutsos J. Pei P. S. Yu X. Lin C. Aggarwal...	6 KDD ICDM ICDE TKDE ...	12 “large graphs” “data streams” “evolving data” “evolving graphs” ...	70 2006- 2013

² “Frequent subgraph discovery”

Author	Venue	Keyword	Cited	#Paper	Author	Venue	Keyword	#Paper
12 Ryen White Hang Li Tie-Yan Liu Zhaohui Zheng...	5 SIGIR WWW WSDM CIKM...	3 “web search” “click-through data” “sponsored search”	12 p82630 ³ p116290 p103899 p106191...	32 2006- 2013	8 Qiang Yang Dou Shen Sinno Pan...	3 KDD PAKDD AAAI	6 “transfer learning” “data mining” “localization models”	17 2007- 2010

³ “Optimizing search engines using clickthrough data”



Experimental Results

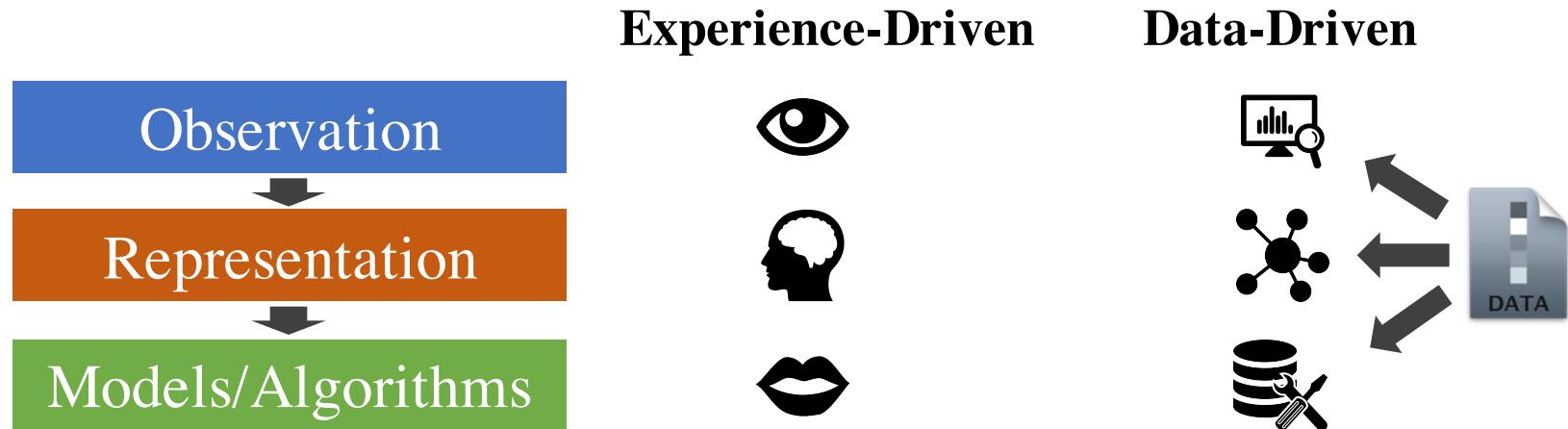
Event summaries with Super Bowl 2013 tweets

							user	phrase	hashtag	URL	3,397 tweets
16:30		16:30:31 <u>My prediction</u> Ravens 34 Niners 31 16:30:57 Ready for the big game :D, <u>my prediction</u> 24-20 SF #SuperBowl	“my prediction”	(3,325)	226	(0)	(0)				Tartan #1: (1 dim) 16:30-17:30
17:00		16:31:14 <u>My prediction for superbowl..</u> 48.. Jets over Bears 17-13 Mark Sanchez MVP 16:32:24 <u>I predict Baltimore Ravens</u> will win 27 to 24 or 25 or 26. Basically it will be a <u>close game</u> .									Tartan #2: (3 dims) 17:00-18:00
17:30		17:30:51 RT @LMAOTWITPICTS: <u>Make Your Prediction. Retweet For 49ers</u> http://t.co/KKksEist 17:31:01 RT @LMAOTWITPICTS: <u>Make Your Prediction. Retweet For 49ers</u> http://t.co/KKksEist 17:31:16 RT @LMAOTWITPICTS: <u>Make Your Prediction. Retweet For 49ers</u> http://t.co/KKksEist 17:31:19 RT @LMAOTWITPICTS: <u>Make Your Prediction. Retweet For 49ers</u> http://t.co/KKksEist	“make your prediction”	(196)	4	1	1				
18:00		18:55:03 RT @49ers: <u>Kaepernick is sacked on 3rd and goal. #49ers K David Akers makes 36-yard FG. Baltimore leads 7-3 with 3:58 left in 1st Qtr. #SB47</u> 18:55:04 RT @49ers: <u>Kaepernick is sacked on 3rd and goal. #49ers K David Akers makes 36-yard FG. Baltimore leads 7-3 with 3:58 left in 1st Qtr. #SB47</u> 18:55:44 RT @Ravens: <u>David Akers is good from 36 yards to make the score 7-3 Ravens. Nice job by the defense to tighten up in the red zone.</u>	“7-3”, “1 st Qtr”	(213)	21	3	(0)				Tartan #3: (2 dims) 18:30-19:30
18:30											
19:00		20:20:01 RT @ExtraGrumpyCat: <u>No Superbowl halftime show will ever surpass this. http://t.co/0VSy7Cv6</u> 20:20:02 RT @WolfpackAlan: <u>No Superbowl halftime show will ever surpass this. http://t.co/6Bll0PXs</u> 20:20:04 RT @ExtraGrumpyCat: <u>No Superbowl halftime show will ever surpass this. http://t.co/0VSy7Cv6</u> 20:20:05 RT @WolfpackAlan: <u>No Superbowl halftime show will ever surpass this. http://t.co/6Bll0PXs</u>	halftime show”	(617)	11	4	4				Tartan #4: (3 dims) 20:00-21:00
19:30											
20:00		20:20:47 (Manhattan, NY)...and every one of those girls took #ballet #Beyonce #superbowl 20:22:01 (New York, NY) I have <u>the biggest lady boner for Beyonce #BeyonceBowl #DestinyBowl #DestinysChild #SuperBowl</u>									Tartan #5: (3 dims) 20:00-21:00
20:30		20:24:32 (Manhattan, NY) No one can ever <u>top that performance by Beyonce EVER. #Beyonce #superbowl #halftimeshow</u>	“beyonce”, #beyonce, #superbowl, #DestinysChild	2	55	17	(0)				
21:00		21:44:42 Ahora si pff #49ers 23-28 #Ravens 21:44:44 Baltimore #Ravens 28-23 San Francisco #49ers 21:44:50 FG Akers #49ers 23-28 #Ravens 3Q 3:10 #SuperBowlXLVII #SuperBowl #NFL	“28-23”, #49ers, #Ravens	(650)	69	11	(0)				Tartan #6: (2 dims) 21:00-22:00
21:30											
22:00		22:42:27 <u>Congratulations Ravens!!!!</u> 22:42:43 <u>Congratulations Ray Lewis and the Ravens.</u> 22:42:43 <u>Game over! Ravens won ray got his retirement ring now all y'all boys and girls go to sleep!</u> 22:42:52 <u>@LetThatBoyTweet: Game over. Ravens win the Super Bowl.</u>	“congratulations”, “game over”	(1942)	248	(0)	(0)				Tartan #7: (1 dim) 22:00-23:30



Summary

- ❑ Structuring text into heterogeneous information networks
- ❑ **Observations, Representations, Models**
 - ❑ **ToPMine/SegPhrase:** Quality phrase mining
 - ❑ **ClusType:** Entity recognition and typing
 - ❑ **MetaPAD:** Data-driven automatic attribute discovery for attributed network construction
 - ❑ Integrating text mining techniques
 - ❑ **Meta Pattern Mining**
- ❑ Integrating phrases into behavioral analysis
- ❑ **Observations, Representations, Models**
 - ❑ **CatchTartan:** Dynamic multicontextual. Tensor fails.

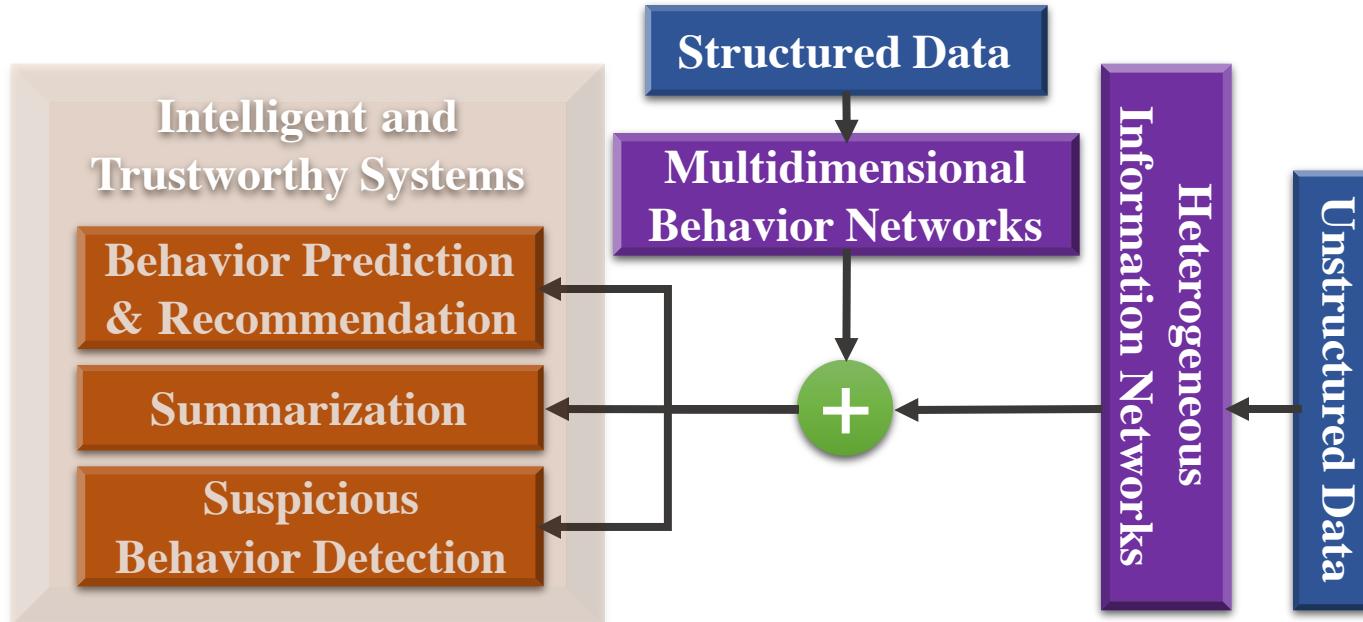


Conclusion

Data-Driven Behavioral Analytics

Data-Driven Behavioral Analytics

- ❑ Mining behavior networks with social and spatiotemporal contexts to support intelligent and trustworthy systems
 - ❑ Mining for behavior prediction and recommendation
 - ❑ Mining for suspicious behavior detection
- ❑ Structuring behavioral content and integrating behavioral analysis with information networks





Acknowledgement



National Natural Science
Foundation of China



Carnegie
Mellon
University



Microsoft®
Research
微软亚洲研究院



176



References

- D. Blei, A. Ng, and M. Jordan. “Latent dirichlet allocation.” JMLR, 2003.
- J. Herlocker, J. Konstan, L. Terveen, J. Riedl. “Evaluating collaborative filtering recommender systems.” ACM TOIS, 2004.
- Y. Koren, R. Bell, C. Volinsky. “Matrix factorization techniques for recommender systems.” Computer, 2009.
- Y. Koren. “Factorization meets the neighborhood: A multifaceted collaborative filtering model.” KDD, 2008.
- Y. Koren. “Collaborative filtering with temporal dynamics.” CACM, 2010.
- M. Balabanovic and Y. Shoham. “FAB: Content-based, collaborative recommendation.” CACM, 1997.
- N. Liu and Q. Yang. “Eigenrank: A ranking-oriented approach to collaborative filtering.” SIGIR, 2008.
- N. Liu, M. Zhao, and Q. Yang. “Probabilistic latent preference analysis for collaborative filtering.” CIKM, 2009.



References

- H. Ma, H. Yang, M. Lyu, and I. King. “Sorec: Social recommendation using probabilistic matrix factorization.” CIKM, 2008.
- H. Ma, T. Zhou, M. Lyu, and I. King. “Improving recommender systems by incorporating social contextual information.” ACM TOIS, 2011.
- H. Ma, D. Zhou, C. Liu, M. Lyu, and I. King. “Recommender systems with social regularization.” WSDM, 2011.
- J. Leskovec, A. Singh, and J. Kleinberg. “Patterns of influence in a recommendation network.” PAKDD, 2006.
- P. Massa and A. Paolo. “Trust-aware recommender systems.” RecSys, 2007.
- M. Jamali and E. Martin. “TrustWalker: A random walk model for combining trust-based and item-based recommendation.” KDD, 2009.
- H. Ma, I. King, and M. Lyu. “Learning to recommend with social trust ensemble.” SIGIR, 2009.
- H. Ma, I. King, and M. Lyu. “Learning to recommend with explicit and implicit social relations.” ACM TIST, 2011.



References

- M. Faloutsos, P. Faloutsos, and C. Faloutsos. “On power-law relationships of the internet topology.” SIGCOMM, 1999.
- A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Weiner. “Graph structure in the web.” Computer Networks, 2000.
- F. Chung and L. Lu. “The average distances in random graphs with given expected degrees.” PNAS, 2002.
- J. Kleinberg. “Authoritative sources in a hyperlinked environment.” JACM, 1999.
- H. Kwak, C. Lee, H. Park, and S. Moon. “What is Twitter, a social network or a news media?” WWW, 2010.
- B. Hooi, H.A. Song, A. Beutel, N. Shah, K. Shin, and C. Faloutsos. “Fraudar: Bounding graph fraud in the face of camouflage.” KDD, 2016.
- C. Aggarwal and J. Han. “Frequent pattern mining.” Springer, 2014.
- J. Han, J. Pei, B. Mortazavi-Asl, Q. Chen, U. Dayal, and M. Hsu. “FreeSpan: Frequent pattern-projected sequential pattern mining.” KDD, 2000.



References

- X. Yan and J. Han. “gspan: Graph-based substructure pattern mining.” ICDM, 2003.
- X. Yan and J. Han. “CloseGraph: Mining closed frequent graph patterns.” KDD, 2003.
- Y. Sun, J. Han, X. Yan, P.S. Yu, and T. Wu. “PathSim: Meta path-based top-k similarity search in heterogeneous information networks.” VLDB, 2011.
- Y. Sun, Y. Yu, and J. Han. “Ranking-based clustering of heterogeneous information networks with star network schema.” KDD, 2009.
- Y. Sun, J. Han, P. Zhao, Z. Yin, H. Cheng, and T. Wu. “RankClus: Integrating clustering with ranking for heterogeneous information network analysis.” EDBT, 2009.
- Y. Sun, R. Barber, M. Gupta, C. Aggarwar, and J. Han. “Co-author relationship prediction in heterogeneous bibliographic networks.” ASONAM, 2011.
- A. El-Kishky, Y. Song, C. Wang, C.R. Voss, and J. Han. “Scalable topical phrase mining from text corpora.” VLDB, 2014.
- J. Liu, J. Shang, C. Wang, X. Ren, and J. Han. “Mining quality phrases from massive text corpora.” SIGMOD, 2015.



References

- X. Ren, A. El-Kishky, C. Wang, F. Tao, C.R. Voss, and J. Han. “Effective entity recognition and typing by relation phrase-based clustering.” KDD, 2015.
- X. Ren, W. He, M. Qu, C.R. Voss, H. Ji, and J. Han. “Label noise reduction in entity typing by heterogeneous partial-label embedding.” KDD, 2016.
- C. Wang, M. Danilevsky, N. Desai, Y. Zhang, P. Nguyen, T. Taula, and J. Han. “A phrase mining framework for recursive construction of a topical hierarchy.” KDD, 2013.
- E.E. Papalexakis, C. Faloutsos, N.D. Sidiropoulos. “ParCube: Sparse parallelizable tensor decompositions.” PKDD, 2012.
- D. Koutra, U. Kang, J. Vreeken, and C. Faloutsos. “VOG: Summarizing and understanding large graphs.” SDM, 2014.
- R. Gupta, A. Halevy, X. Wang, S.E. Whang, and F. Wu. “Biperpedia: An ontology for search applications.” VLDB, 2014.
- M. Yahya, S. Whang, R. Gupta, and A. Halevy. “ReNoun: Fact extraction for nominal attributes.” EMNLP, 2014.
- A. Halevy, N. Noy, S. Sarawagi, S.E. Whang, and X. Yu. “Discovering structure in the universe of attribute names.” WWW, 2016.



References

Q. Li, Y. Li, J. Gao, B. Zhao, W. Fan, and J. Han. “Resolving conflicts in heterogeneous data by truth discovery and source reliability estimation.” SIGMOD, 2014.

Q. Li, Y. Li, J. Gao, L. Su, B. Zhao, M. Demirbas, W. Fan, and J. Han. “A confidence-aware approach for truth discovery on long-tail data.” VLDB, 2014.

F. Ma, Y. Li, Q. Li, M. Qiu, J. Gao, S. Zhi, L. Su, B. Zhao, H. Ji, and J. Han. “Faitcrowd: Fine grained truth discovery for crowdsourced data aggregation.” KDD, 2015.

Y. Li, J. Gao, C. Meng, Q. Li, L. Su, B. Zhao, W. Fan, and J. Han. “A survey on truth discovery.” KDD Explorations Newsletter, 2016.

S. Zhi, B. Zhao, W. Tong, J. Gao, D. Yu, H. Ji, and J. Han. “Modeling truth existence in truth discovery.” KDD, 2015.

S. Kumar, R. West, and J. Leskovec. “Disinformation on the Web: Impact, characteristics, and detection of Wikipedia hoaxes.” WWW, 2016.

S. Kumar, F. Spezzano, and V.S. Subrahmanian. “Identifying malicious actors on social media.” ASONAM, 2016. (tutorial)



Thank you!

**Data-Driven Behavioral Analytics:
Observations, Representations and Models**