



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>



I. Mining behavior networks with social and spatiotemporal contexts

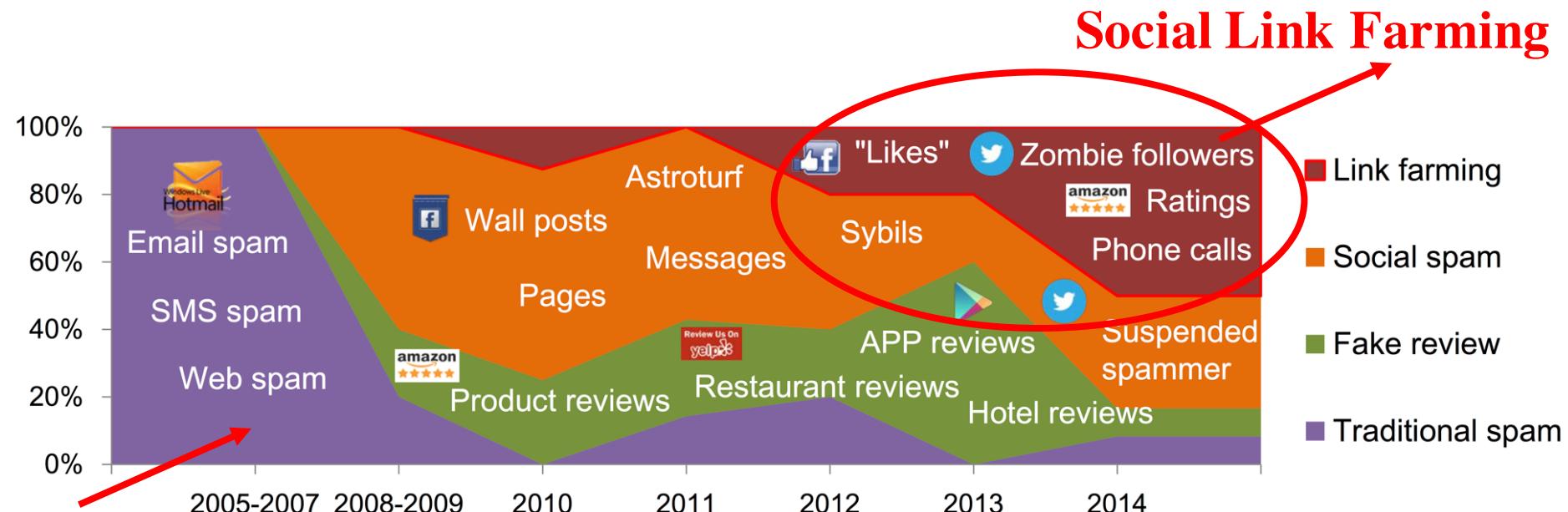
I.2. Suspicious behavior detection



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

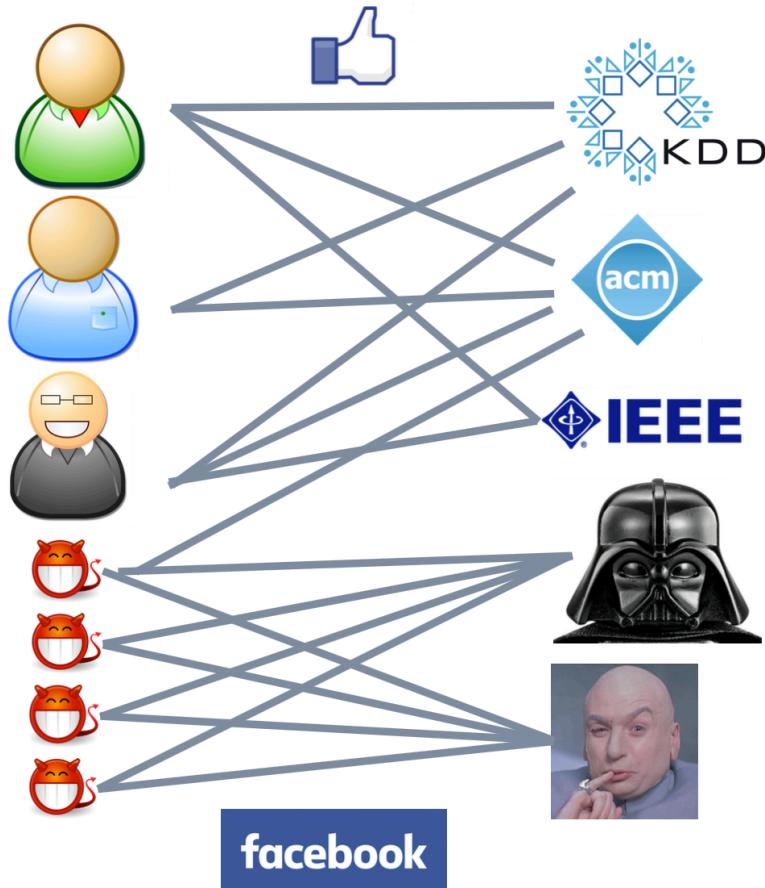
Suspicious Behavior Detection



Meng Jiang, Peng Cui and Christos Faloutsos.

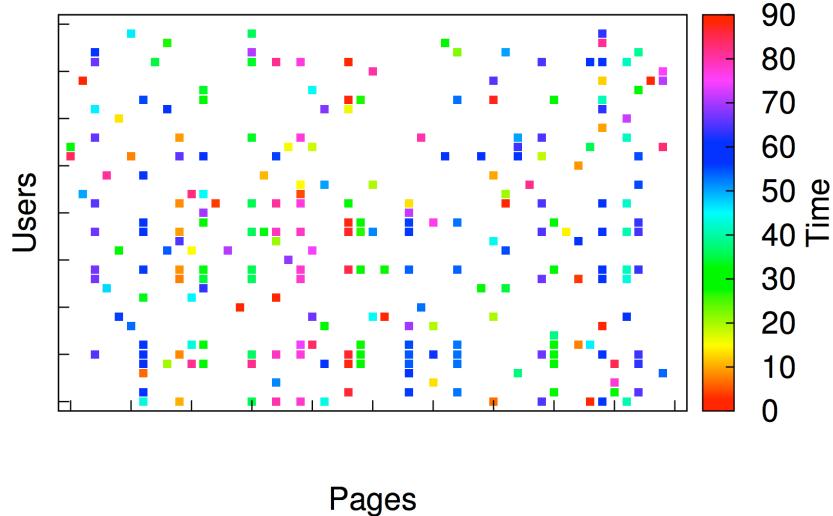
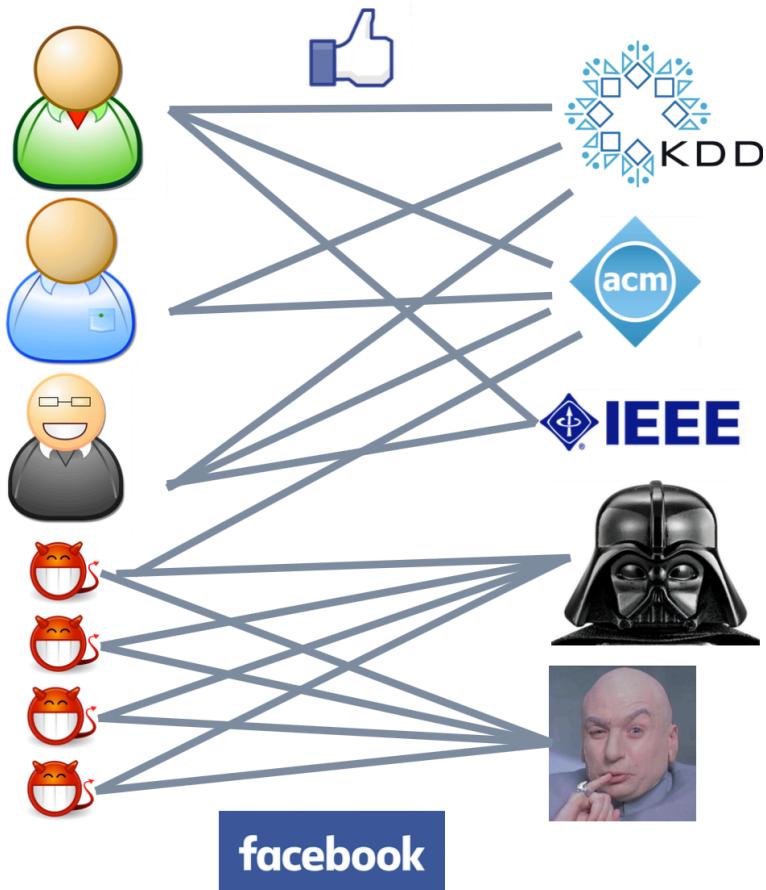
Suspicious Behavior Detection: Current Trends and Future Directions.
IEEE Intelligent Systems (ISSI), 2016.

Ill-gotten Facebook Likes

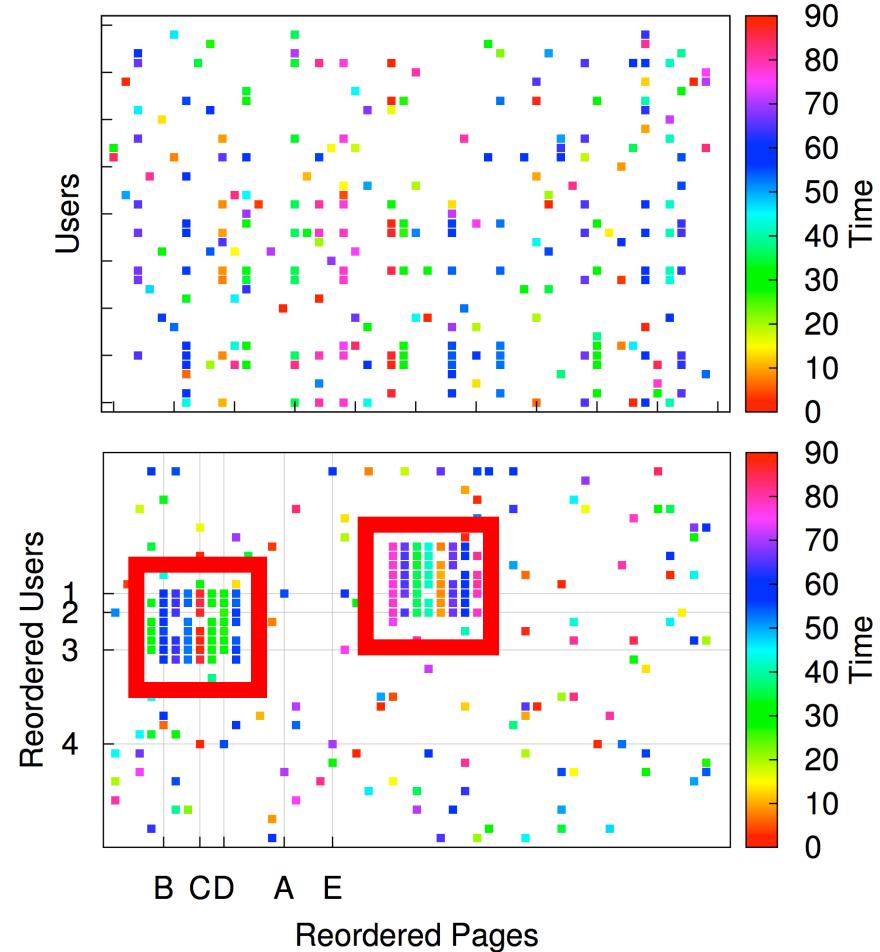
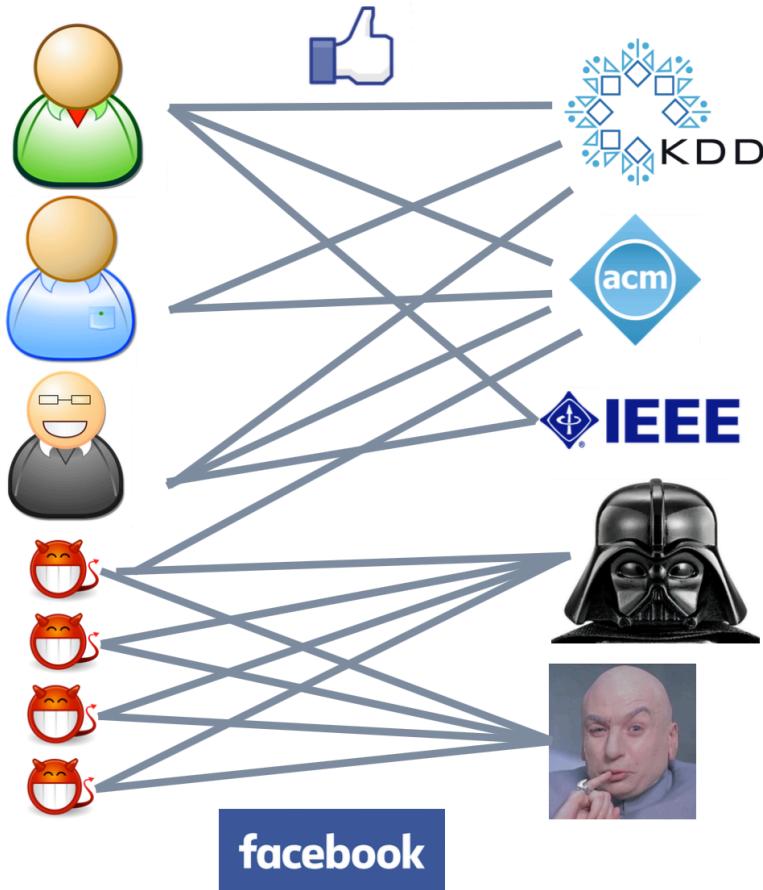


Beutel et al. **CopyCatch: Stopping Group Attacks by Spotting Lockstep Behavior in Social Neworks.** WWW, 2013.

Observation: Graphical View



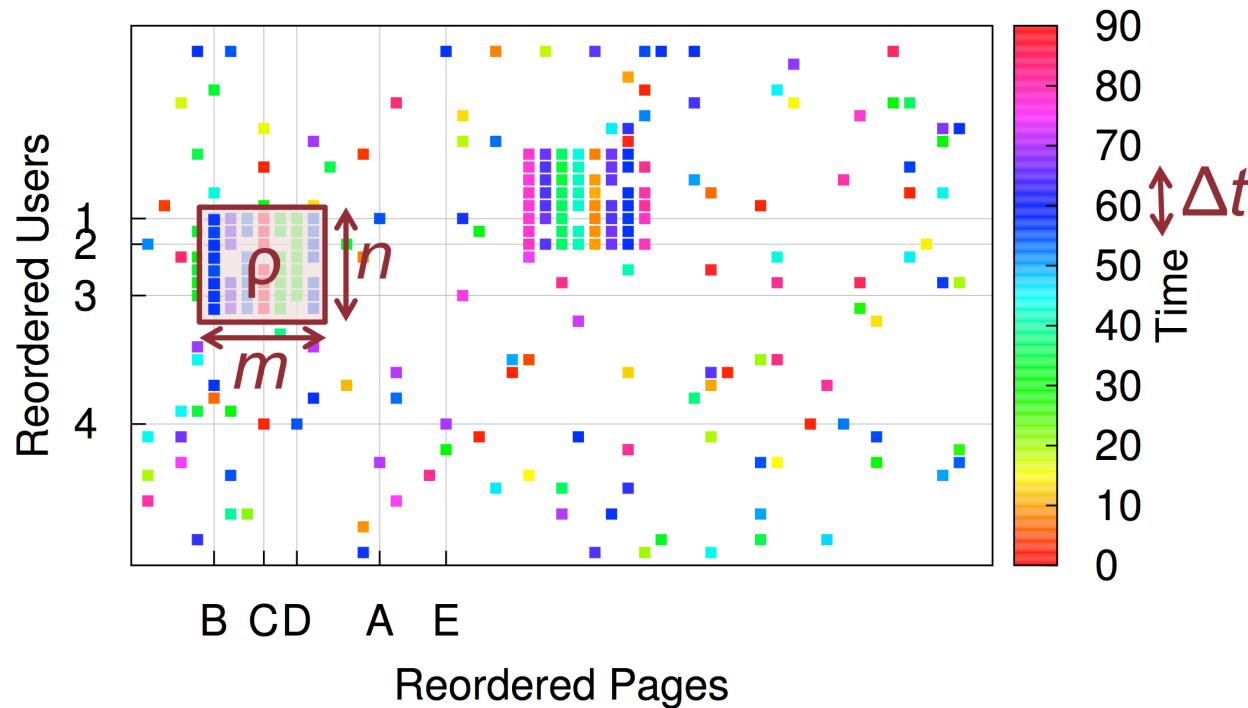
Observation: Reorder Matrix



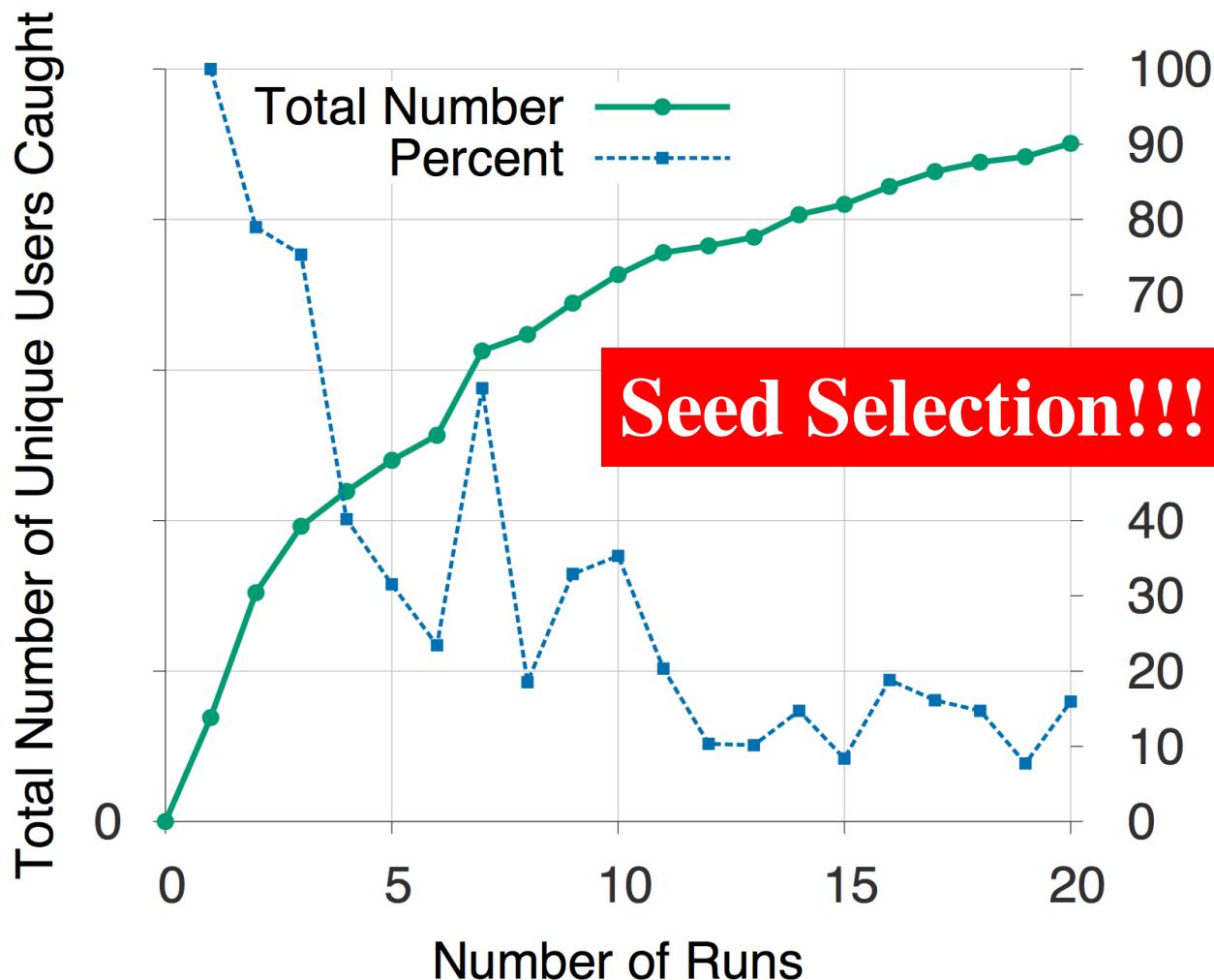
Algorithm: Seed + Search

□ CopyCatch

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



Experimental Result



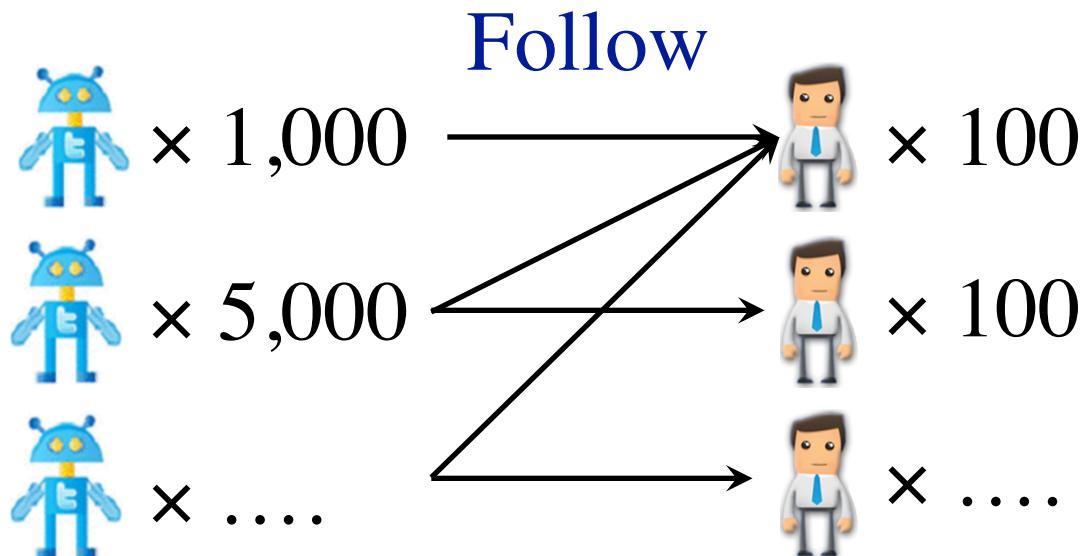
Serious Problem in Weibo



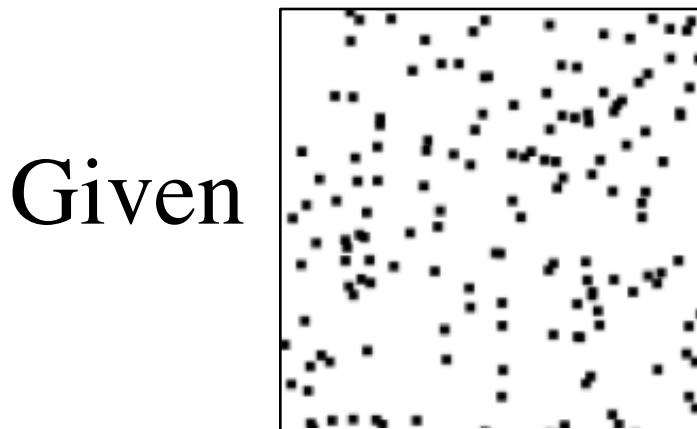
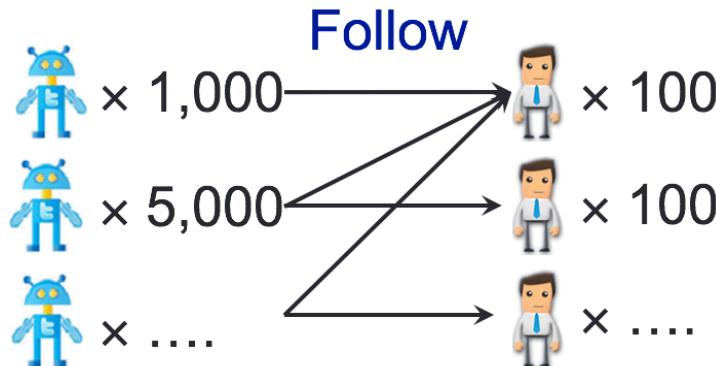
Experience-driven approaches:
features of #followees, #hashtags, #URLs...



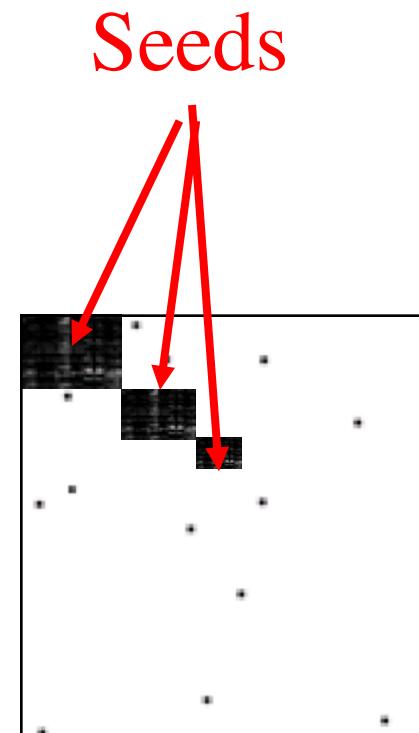
Zombie Followers



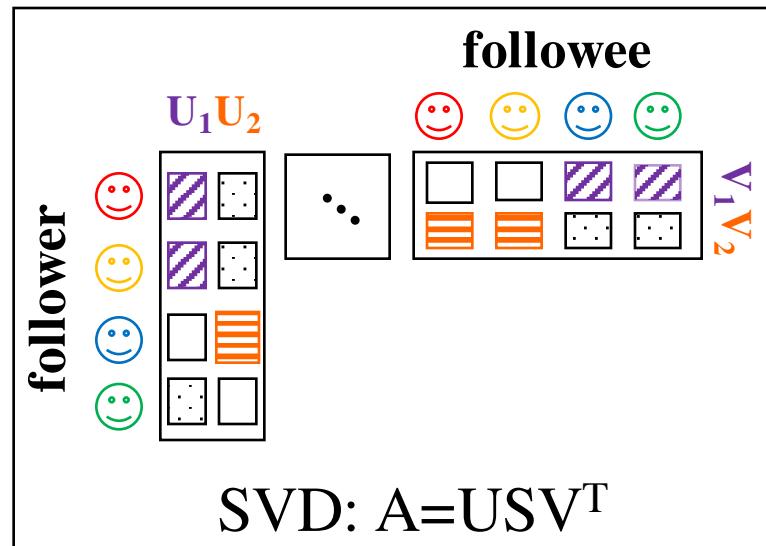
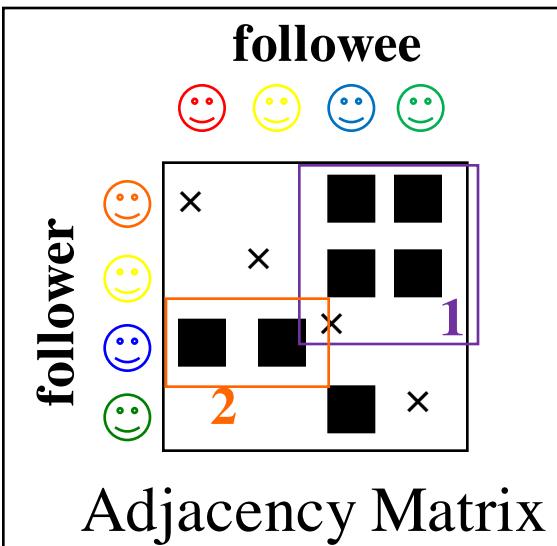
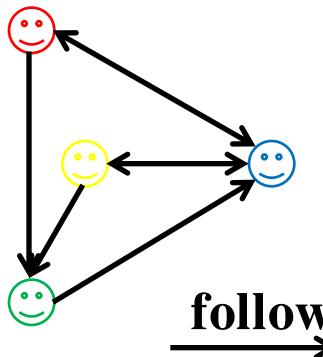
Observation: Reorder Matrix



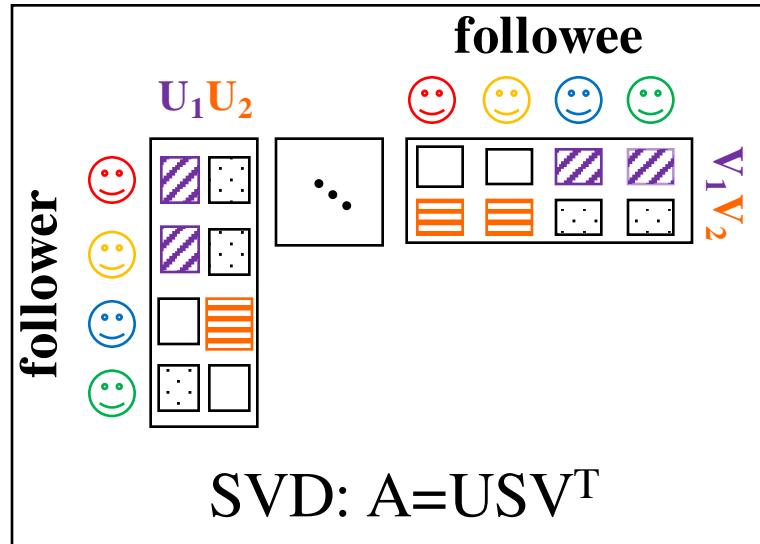
Reorder



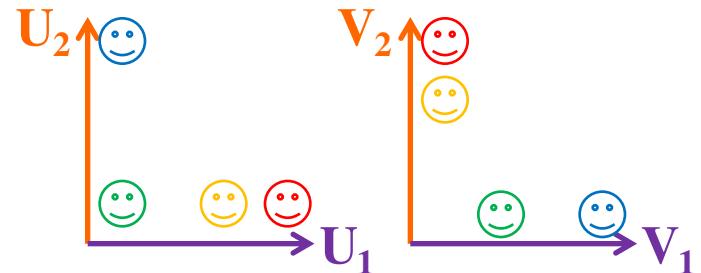
Representation: SVD Reminder



Representation: Spectral Subspace



Pairs of singular vectors:

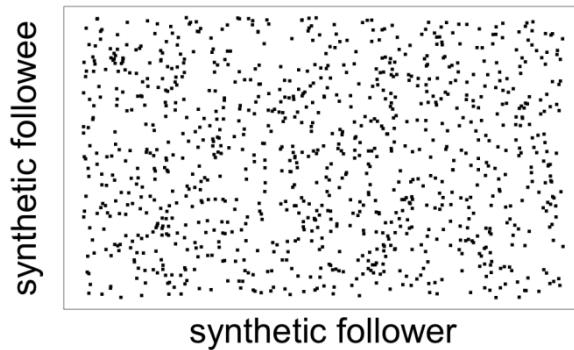


“Spectral Subspace Plot”

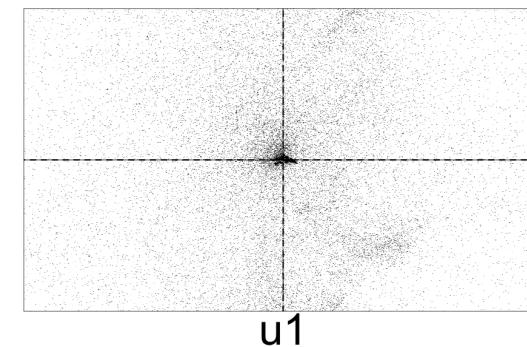
Spectral Subspace Plot: Case #0

- NO lockstep behavior: Scatter

Adjacency Matrix



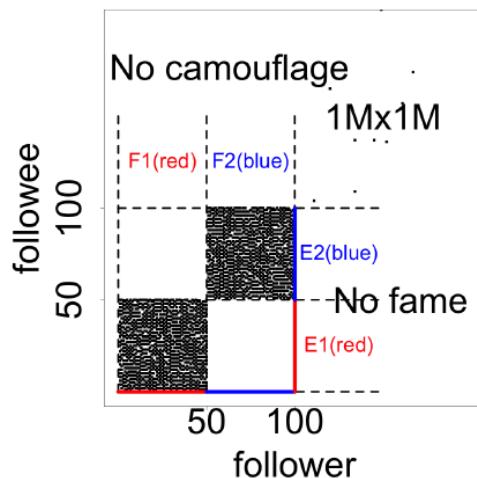
Spectral Subspace Plot



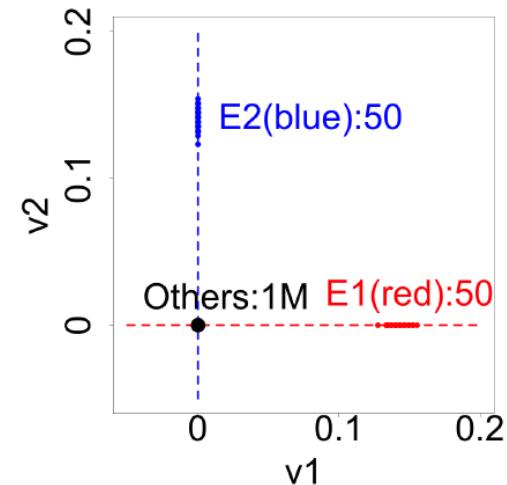
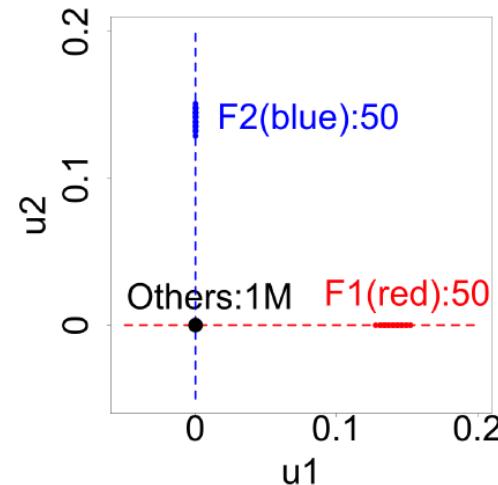
Spectral Subspace Plot: Case #1

- Non-overlapping lockstep: “Rays”

Adjacency Matrix



Spectral Subspace Plot

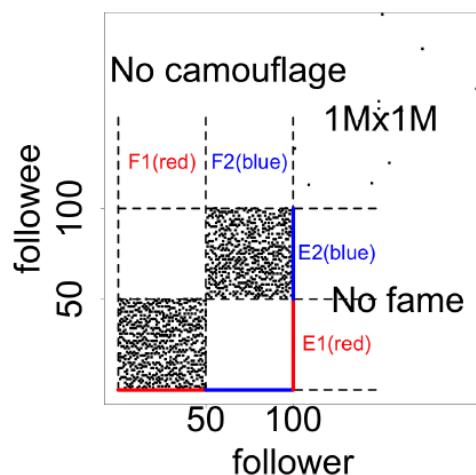


Rule 1 (short “rays”): two blocks, high density (90%), no “camouflage”, no “fame”

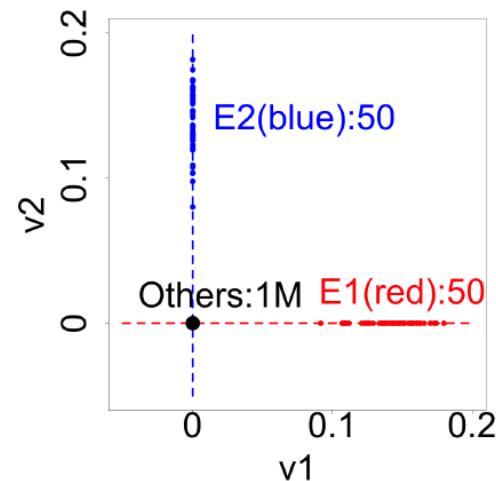
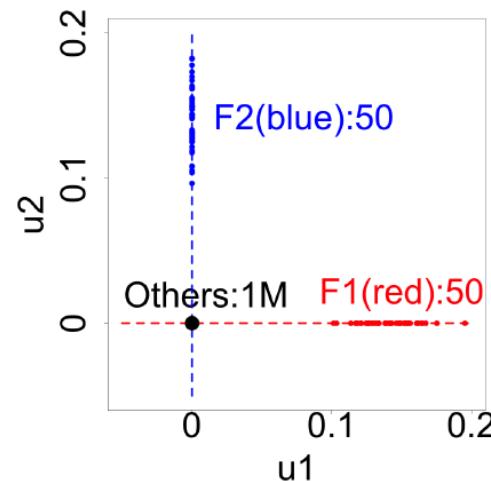
Spectral Subspace Plot: Case #2

- Non-overlapping: Low density, Elongation

Adjacency Matrix



Spectral Subspace Plot

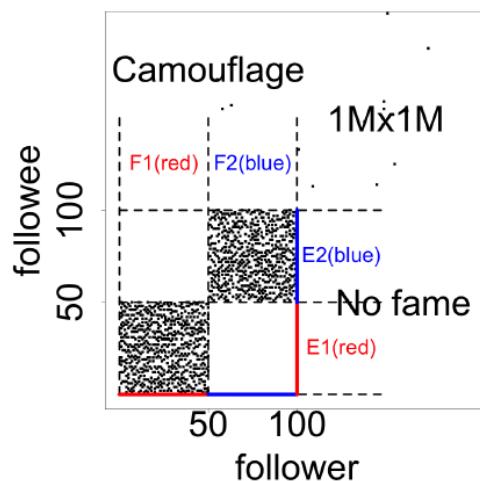


Rule 2 (long “rays”): two blocks, low density (50%), no “camouflage”, no “fame”

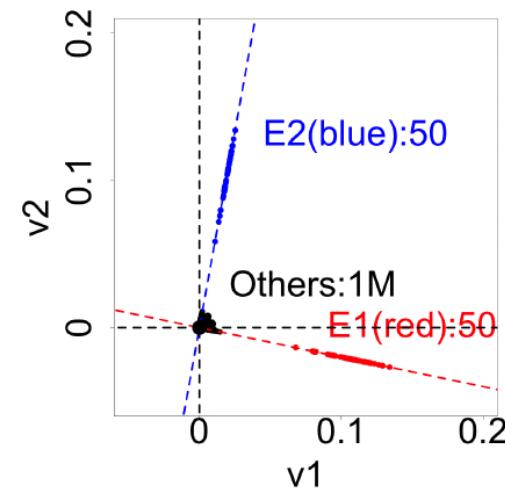
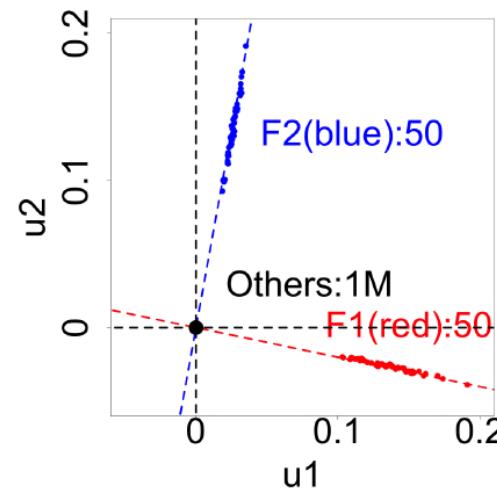
Spectral Subspace Plot: Case #3

- Non-overlapping: Camouflage/Fame, Tilting

Adjacency Matrix



Spectral Subspace Plot

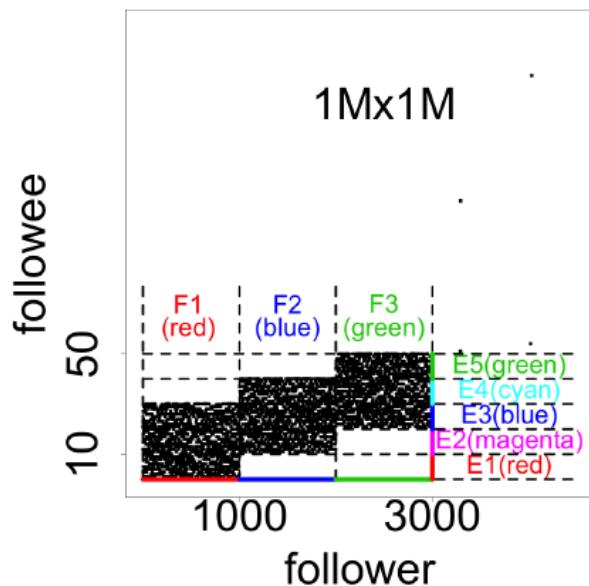


Rule 3 (tilting “rays”): two blocks, with “camouflage”, no “fame”

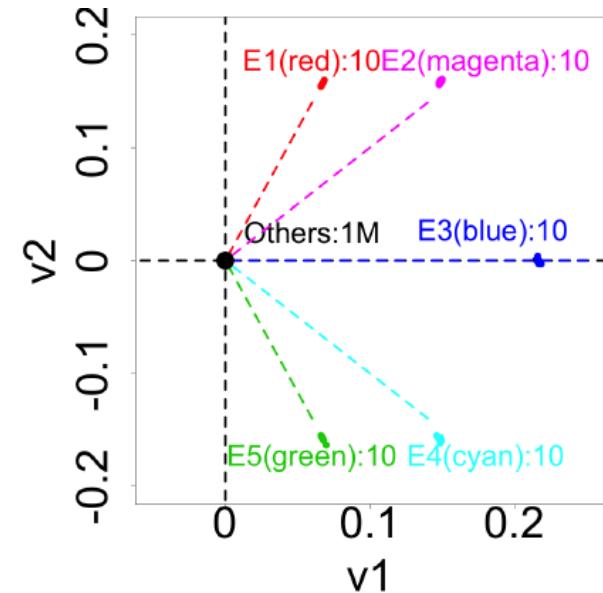
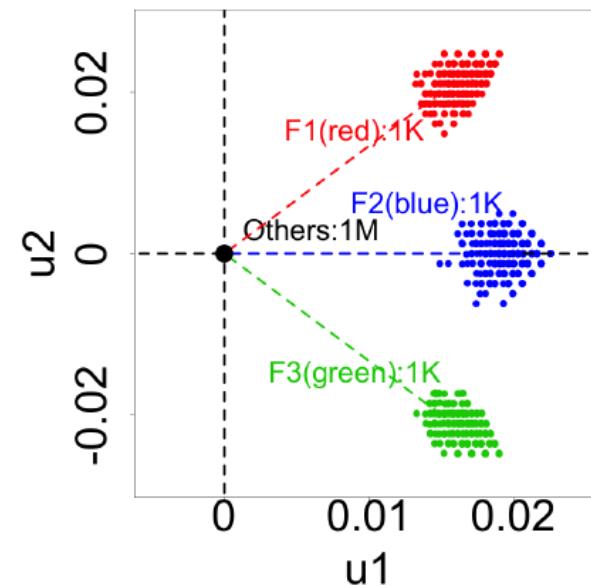
Spectral Subspace Plot: Case #4

- Overlapping: “Staircase”, “Pearls”

Adjacency Matrix



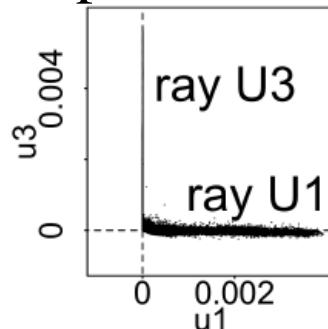
Spectral Subspace Plot



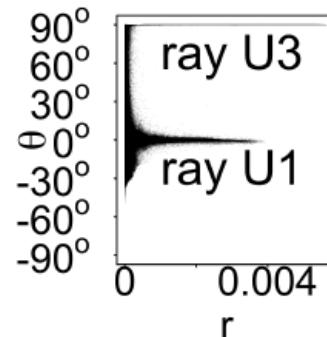
Rule 4 (“pearls”): a “staircase” of three partially overlapping blocks.

Algorithm: Reading & LockInfer

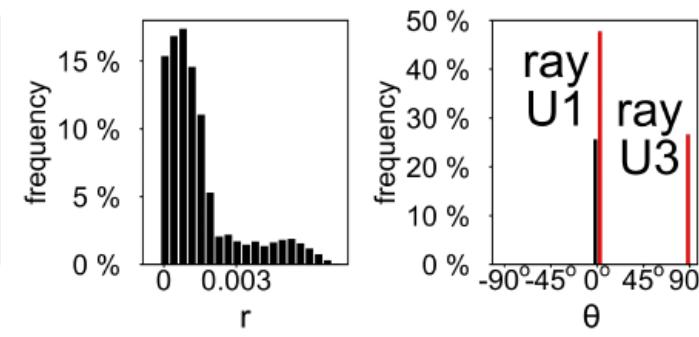
Spectral
Subspace Plot



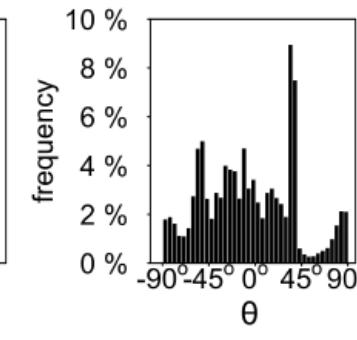
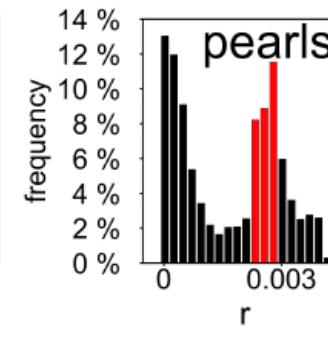
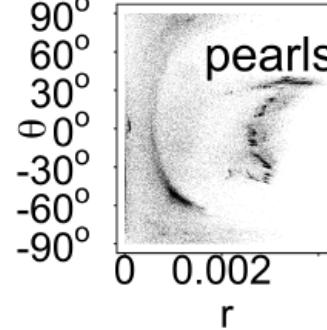
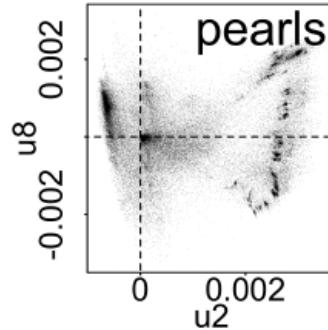
Polar Coordinate
Transform



Histograms



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

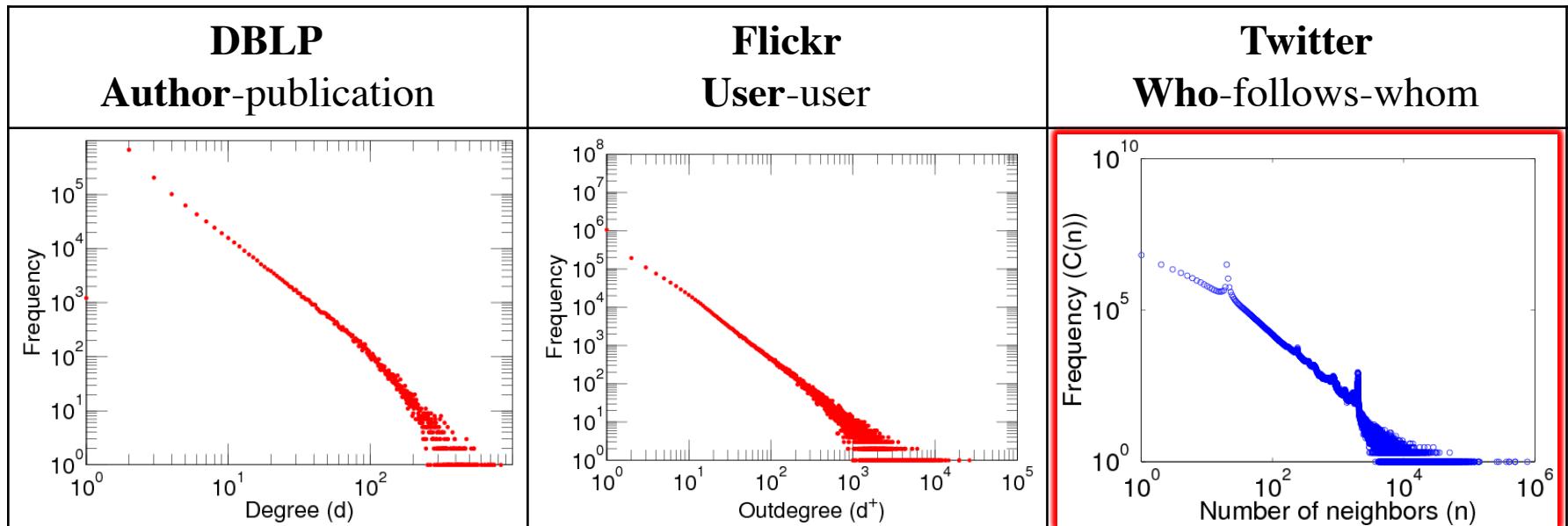


"pearls" show a spike on r frequency at a much-greater-than-zero value

High precision but low recall!!!

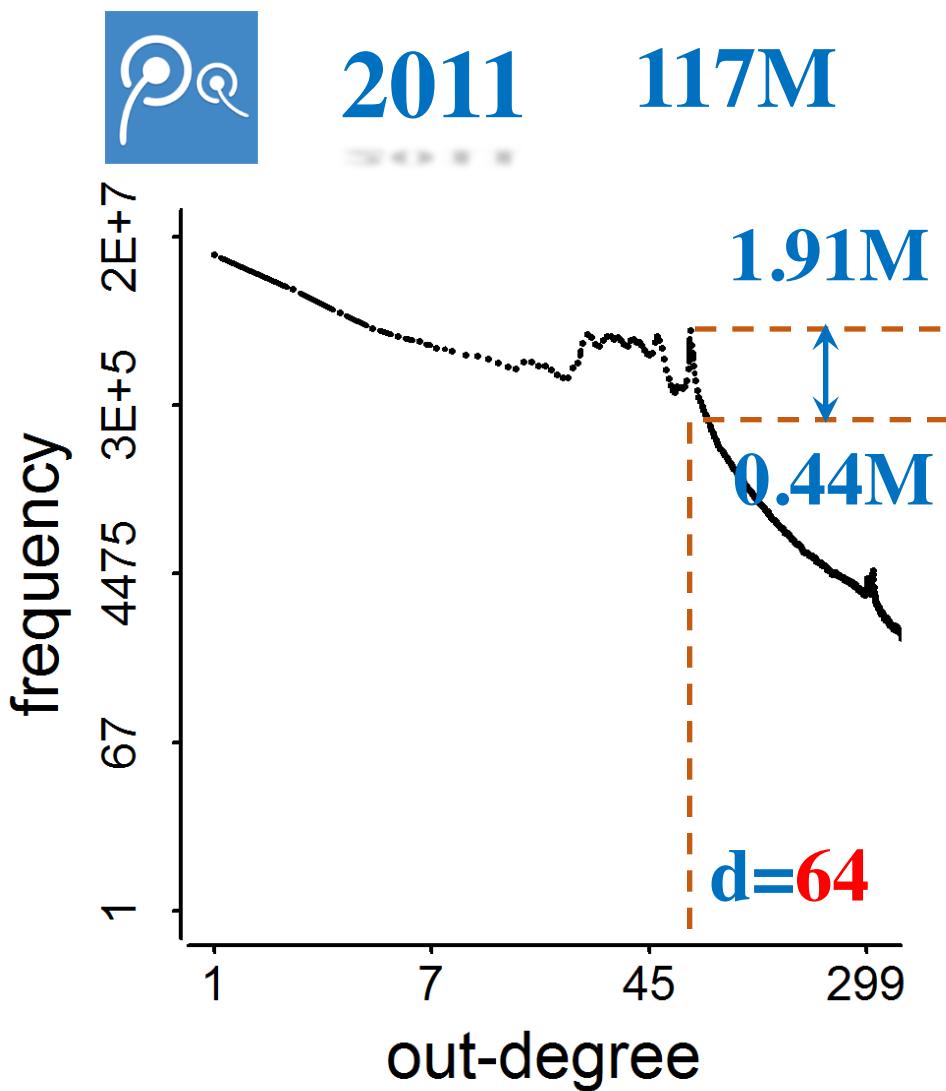
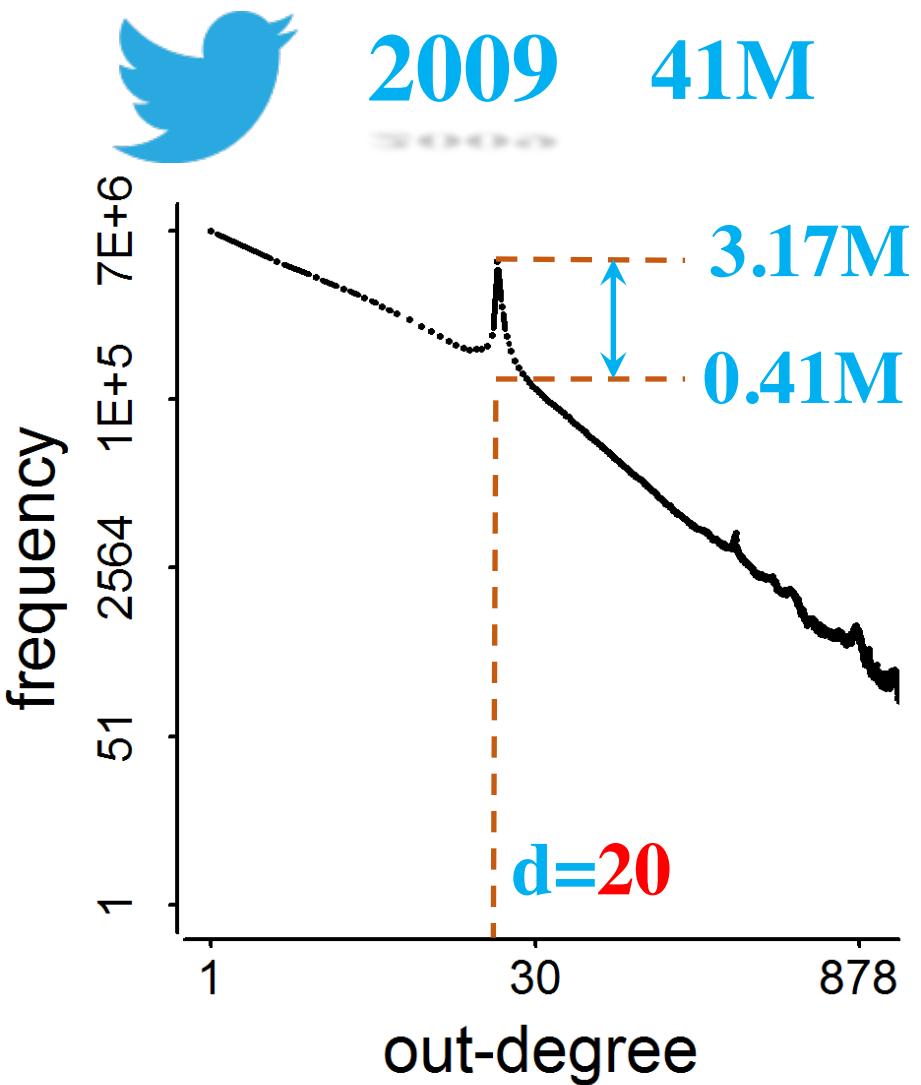
Out-Degree Distributions

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



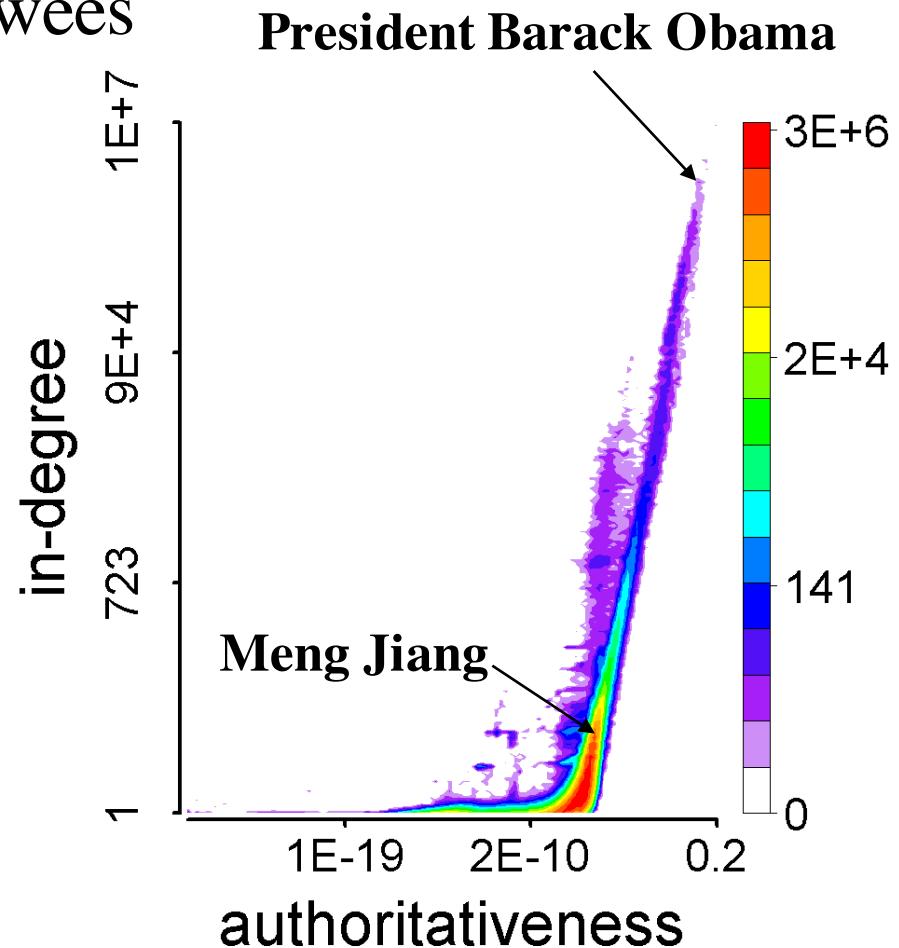
[konect.uni-koblenz.de/networks/]

Spikes!



Observation: How They Behave

- Feature space of followees [Kleinberg. JACM]



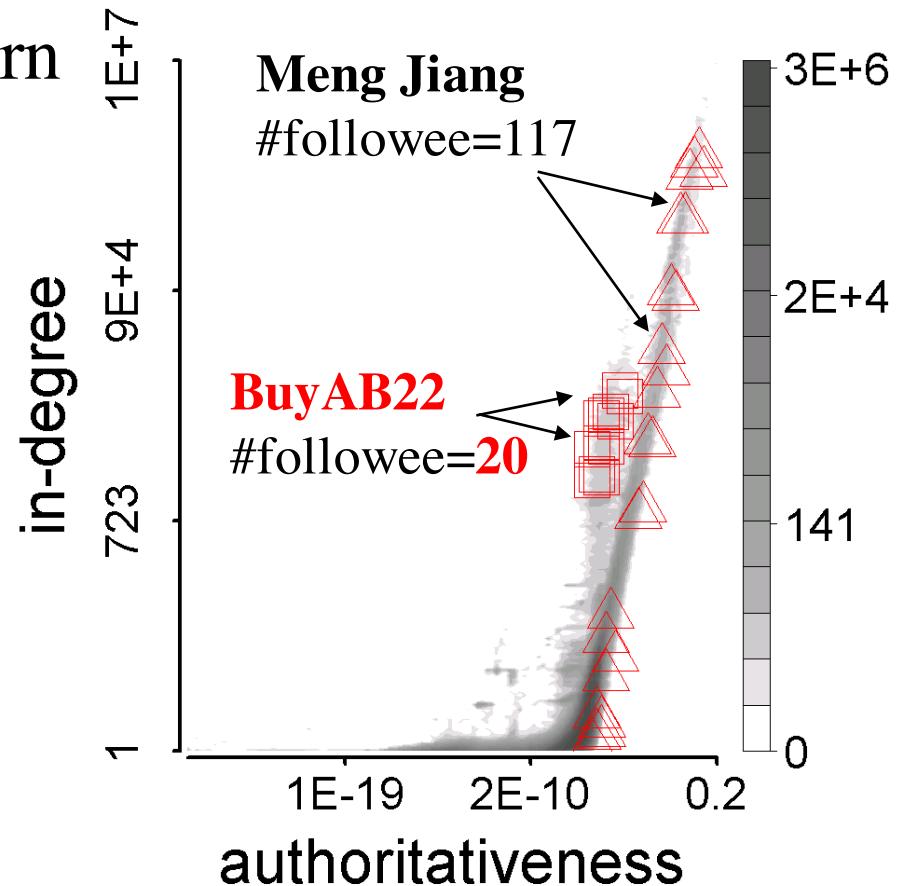
Observation: How They Behave

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

Similar with each other

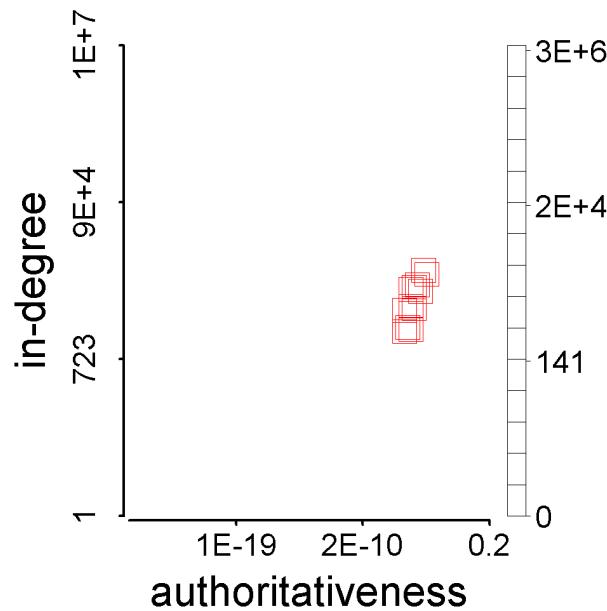
- **Abnormal**

Different from the majority

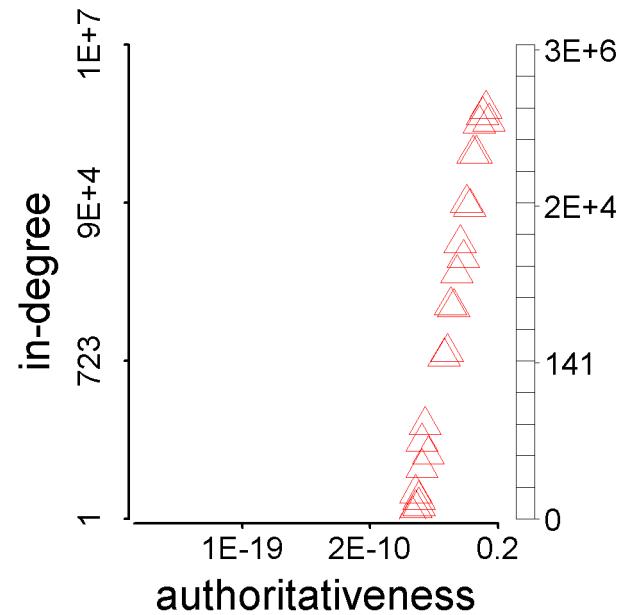


Represent Synchronicity

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

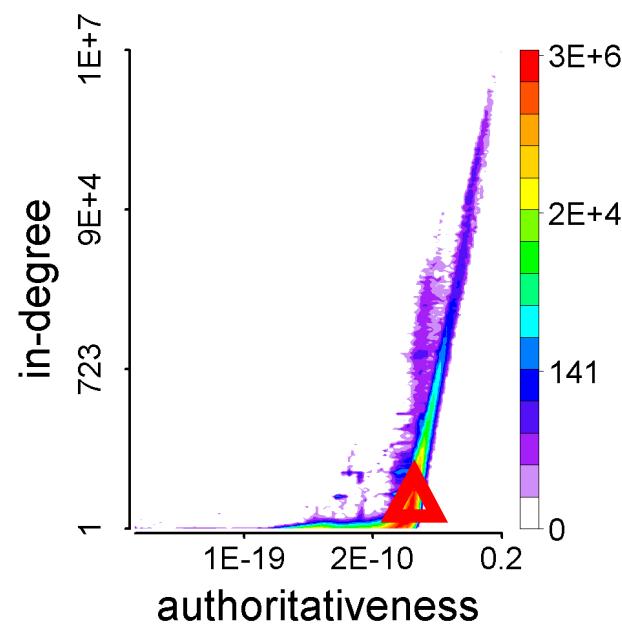
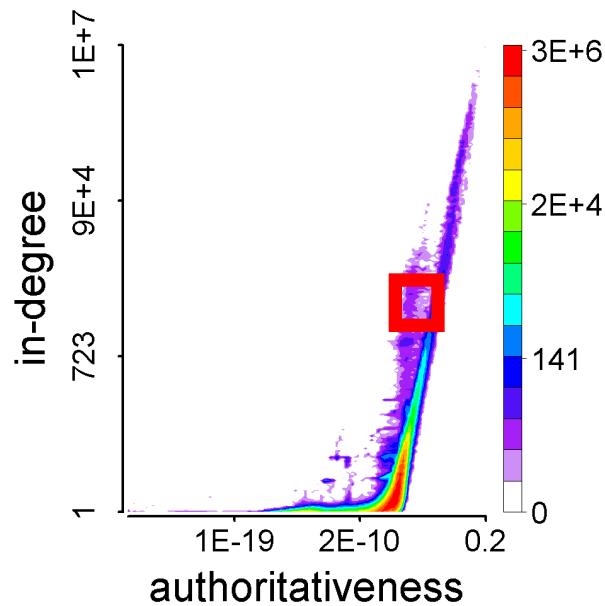


V

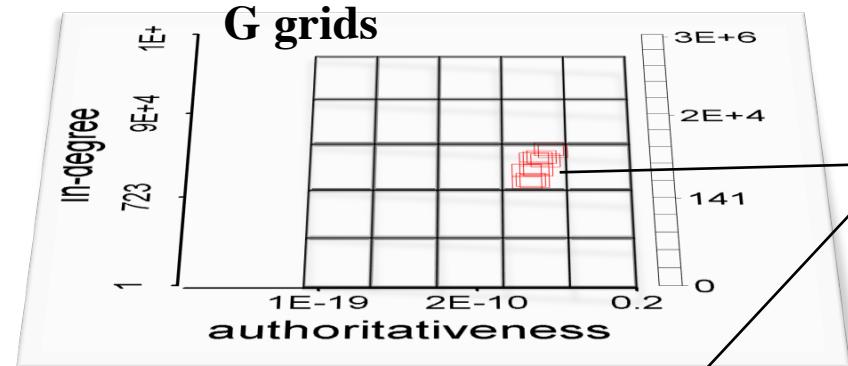


Represent Normality

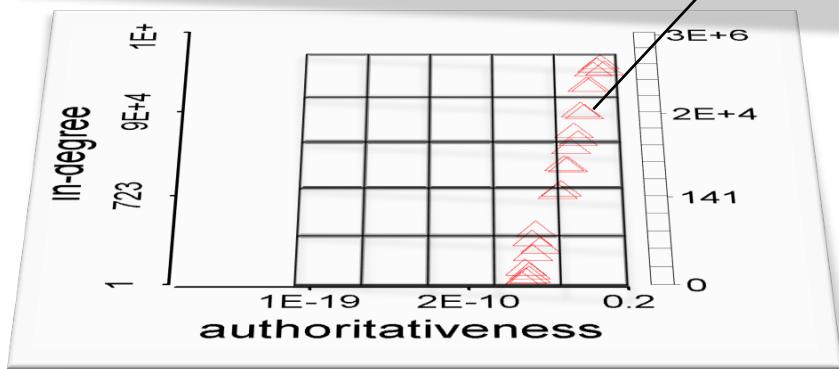
$$\text{norm}(u) = \frac{\sum_{(v,v') \in \mathcal{F}(u) \times \mathcal{U}} \mathbf{p}(v) \cdot \mathbf{p}(v')}{d(u) \times N}$$



Theorem: Synchronicity vs. Normality



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



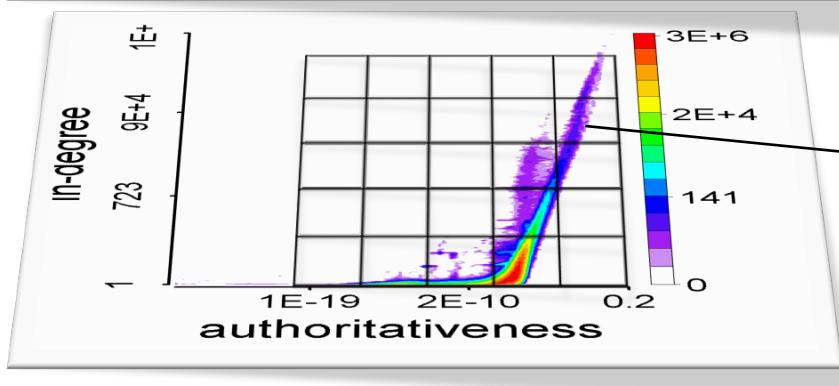
Given normality

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

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

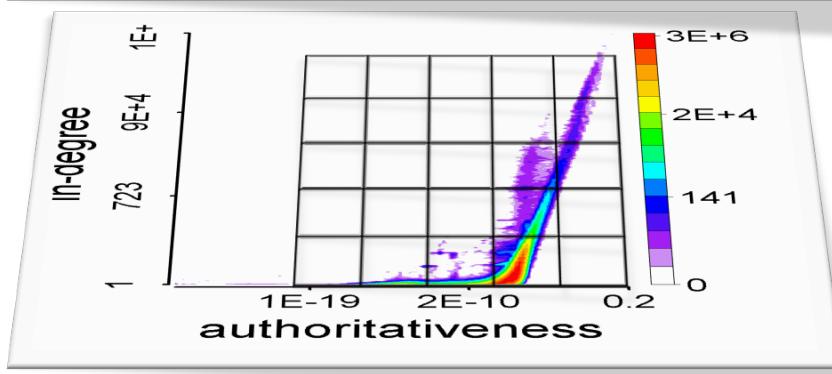
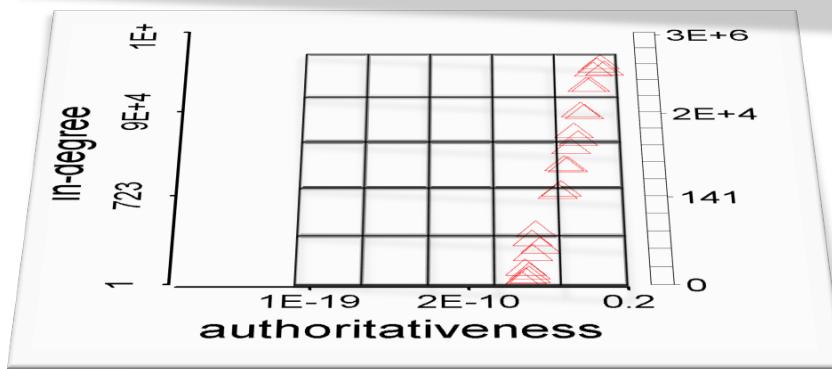
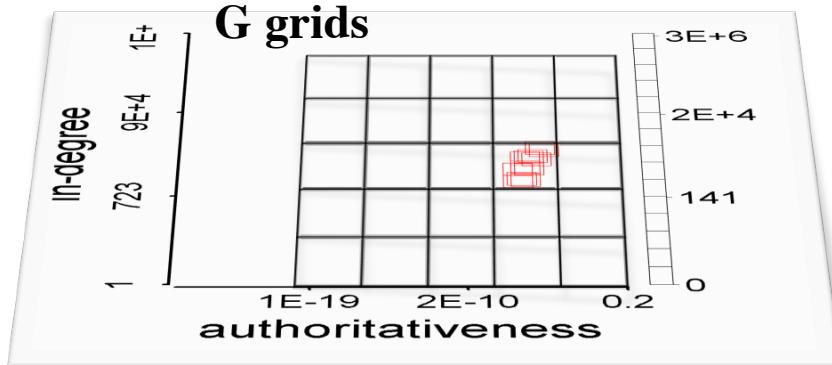
where

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



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 $\times b_g \Sigma$ Σ

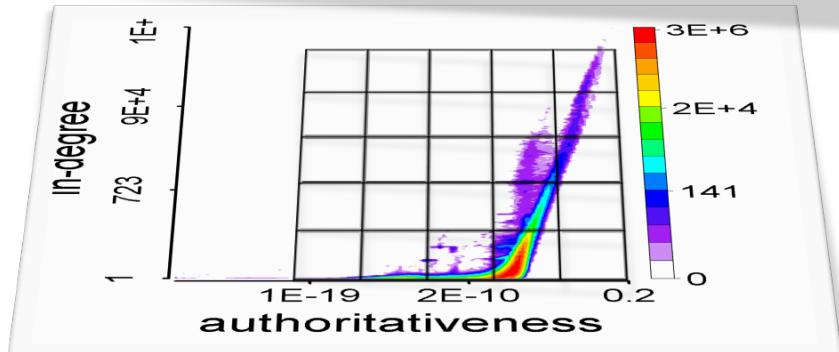
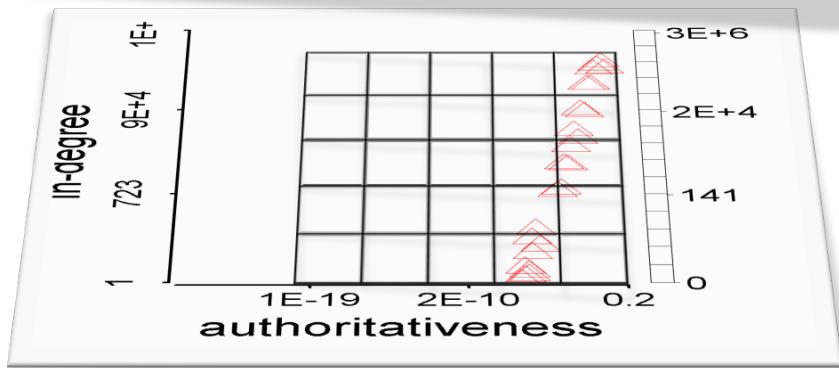
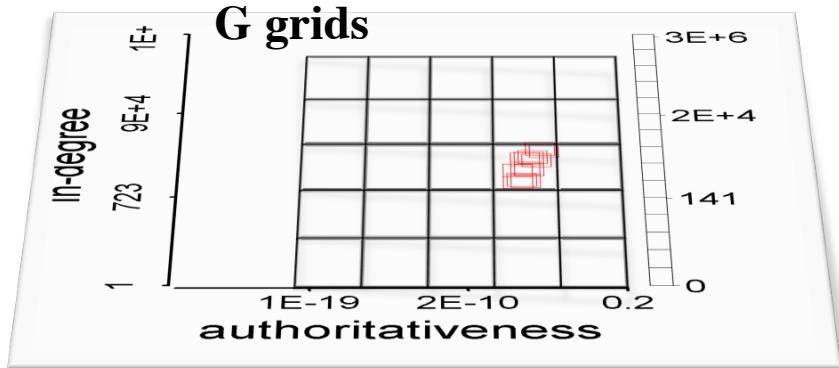
Σ $\times f_g \Sigma$

where $s_b = \sum b_g^2$.

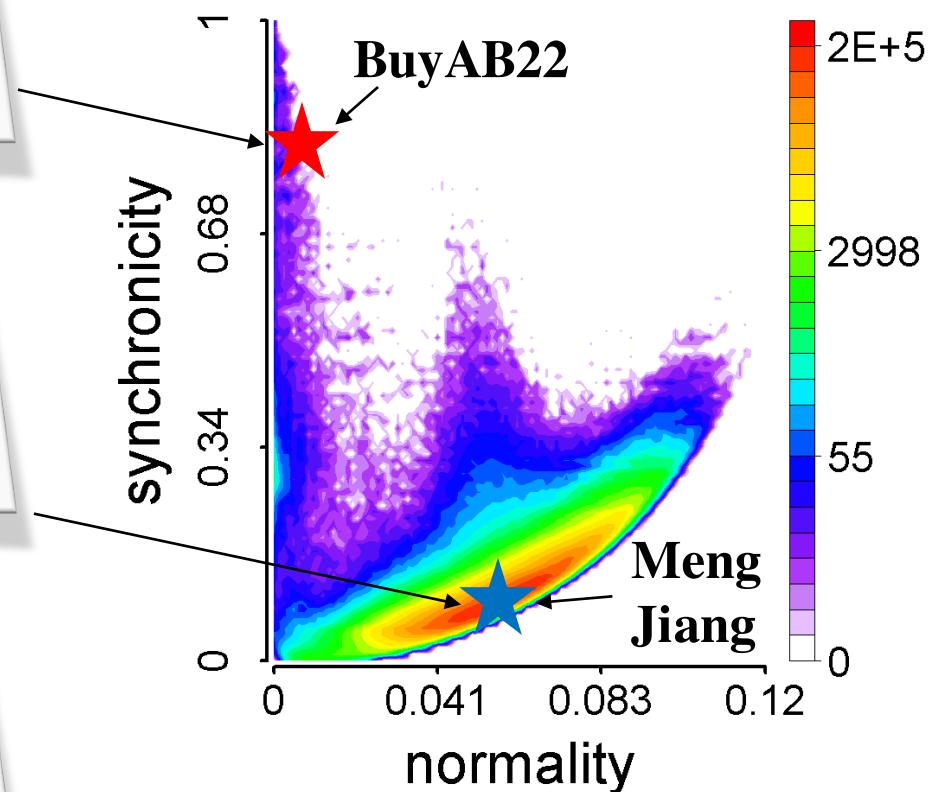
Therefore,

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

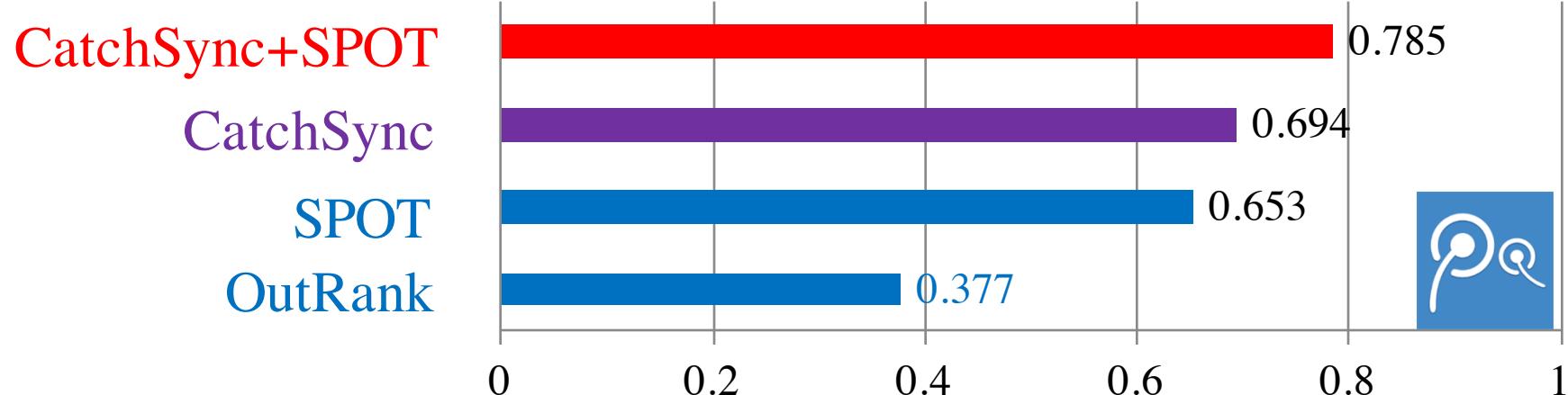
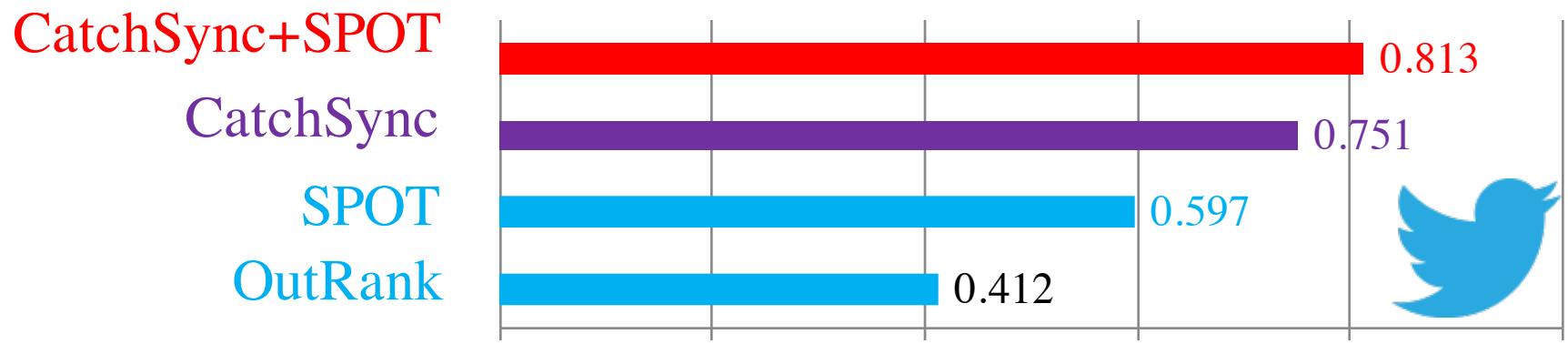
CatchSync Algorithm



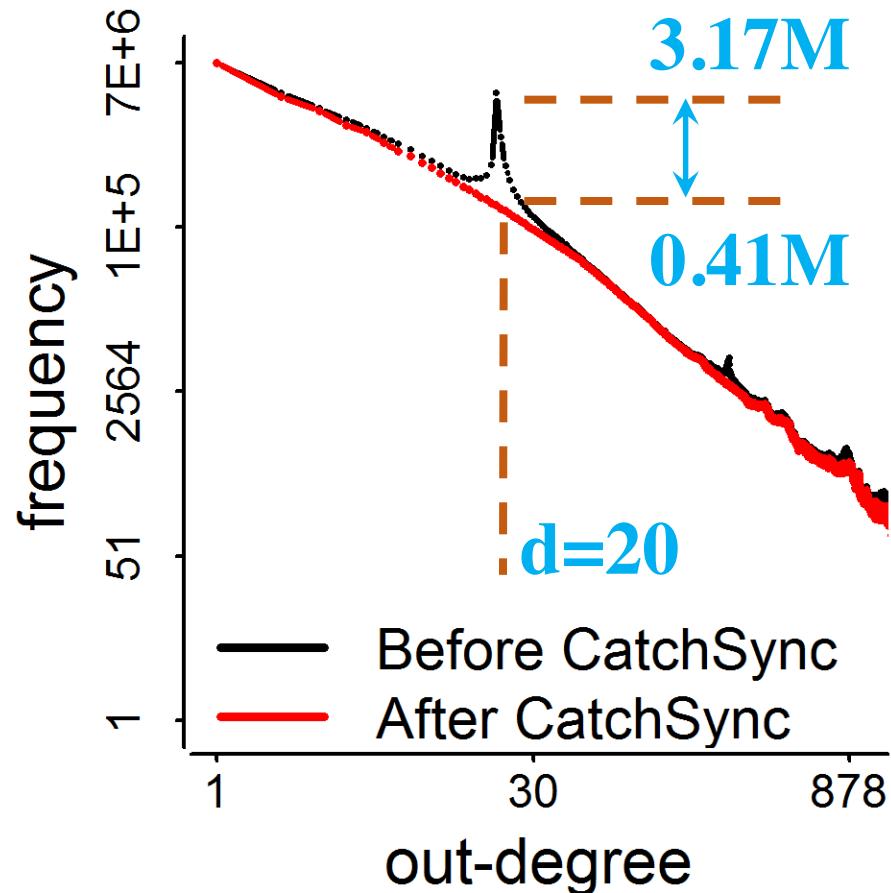
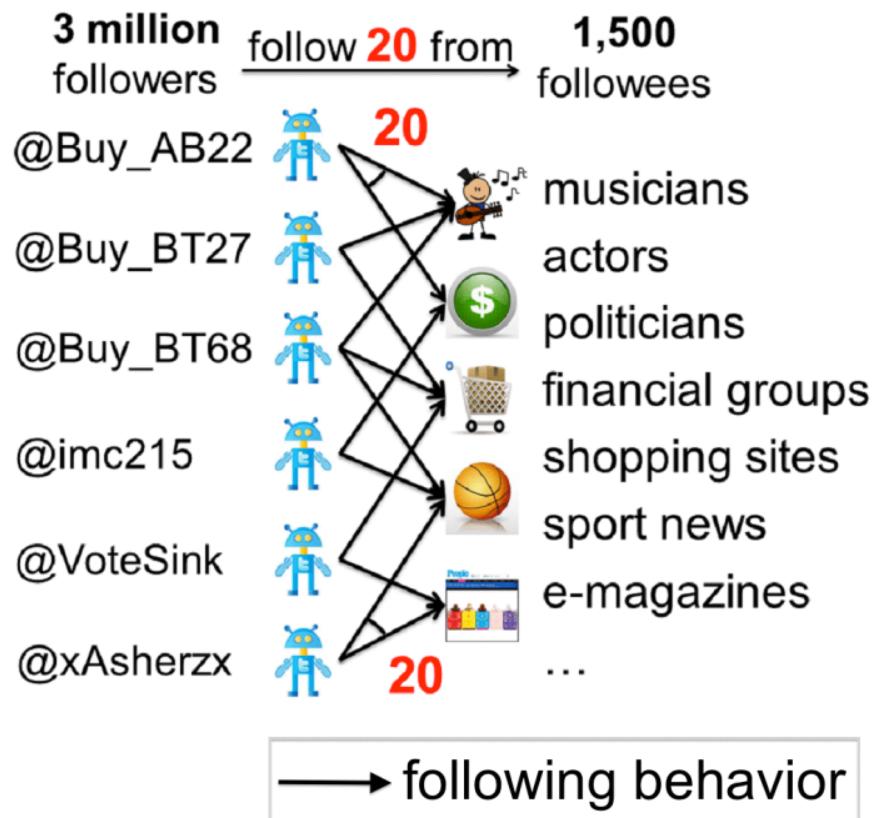
$$s_{\min} = \frac{-G n^2 + 2 n - s_b}{1 - G s_b}$$



Experimental Results



Experimental Results





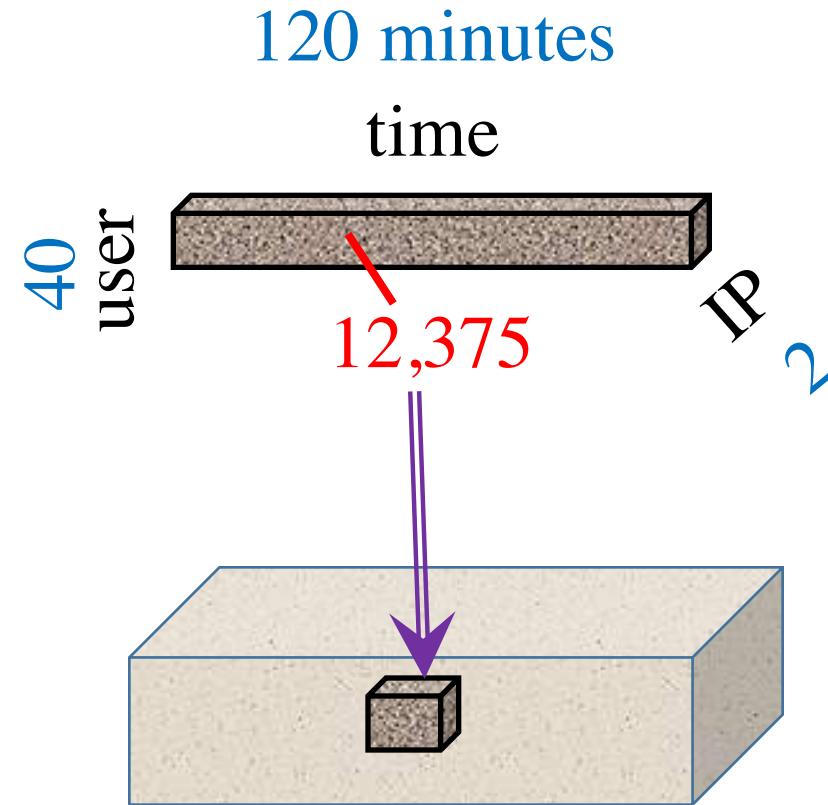
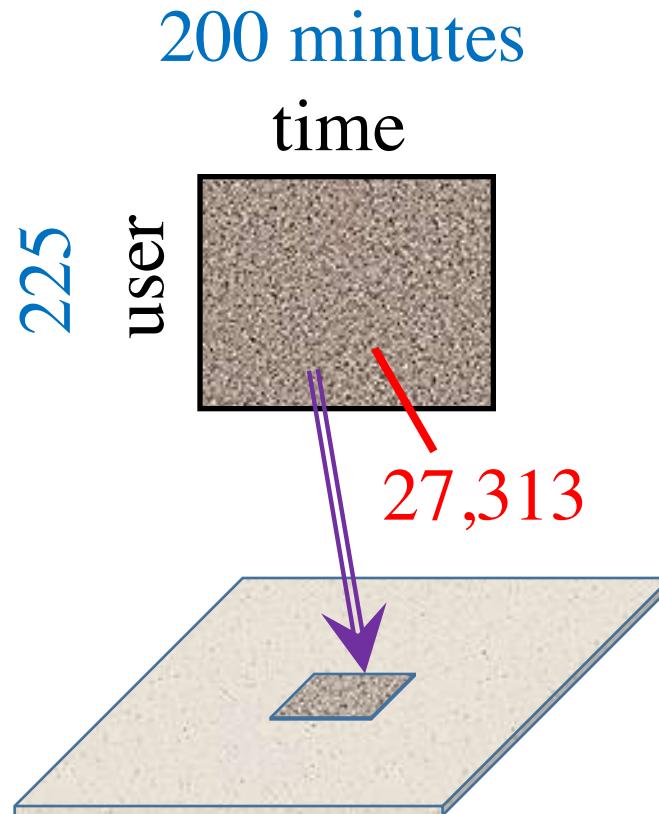
Impact

- ❑ 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

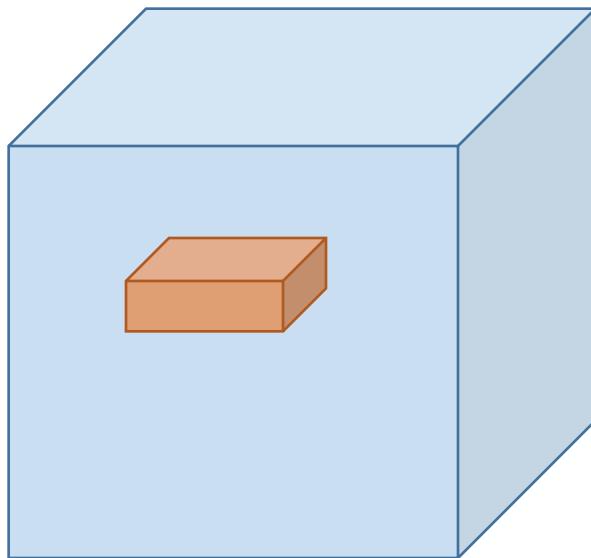


Q: Which is more suspicious?

We need a metric to evaluate the suspiciousness.

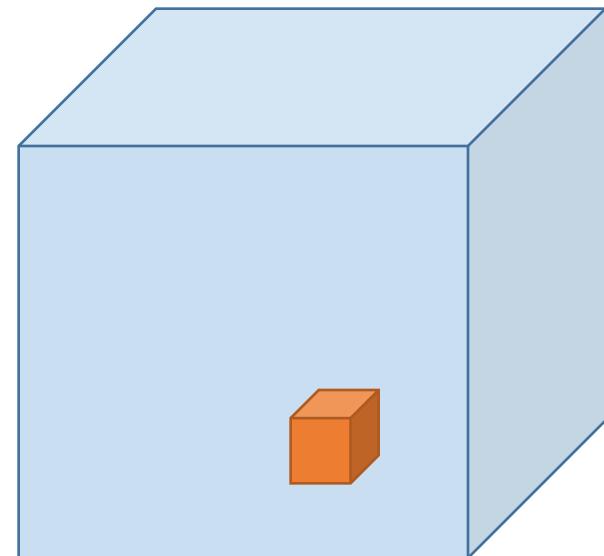
Criteria for Suspiciousness Metric

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

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

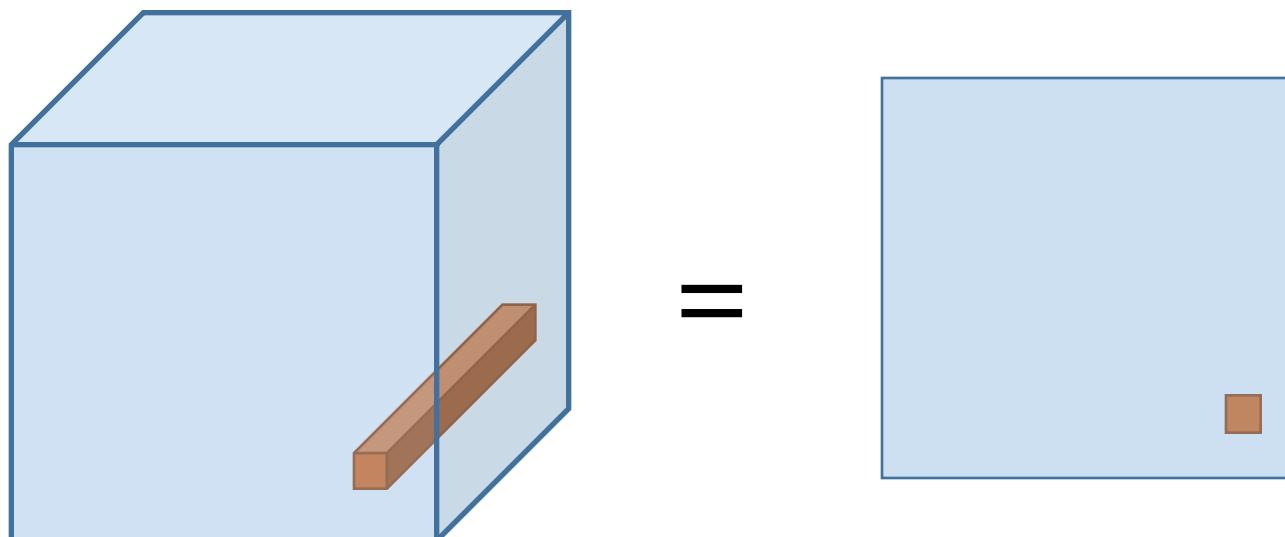
Density Axiom		Contrast Axiom	
	>		
Size Axiom		Concentration Axiom	
	>		

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

Axiom 5: Cross Dimensions

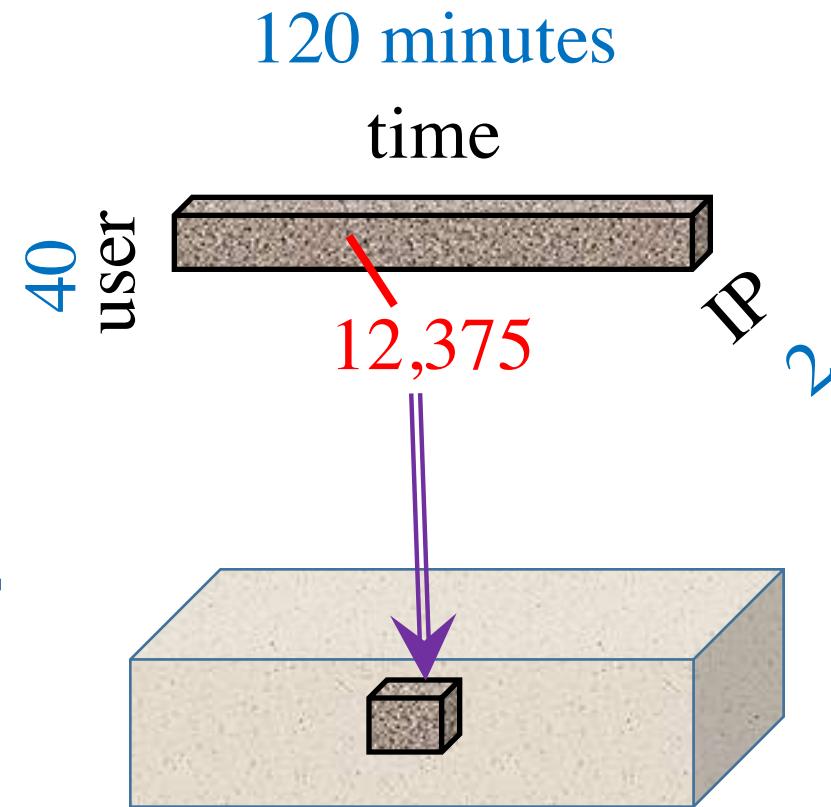
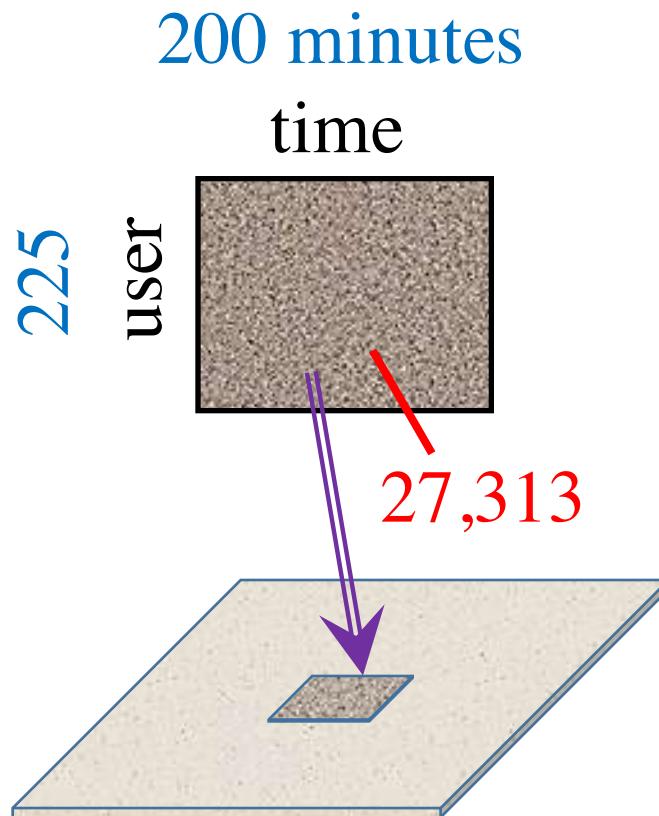
$$f_{K-1} \left([n_k]_{k=1}^{K-1}, c, [N_k]_{k=1}^{K-1}, C \right) = f_K \left(([n_k]_{k=1}^{K-1}, N_K), c, [N_k]_{k=1}^K, C \right)$$

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



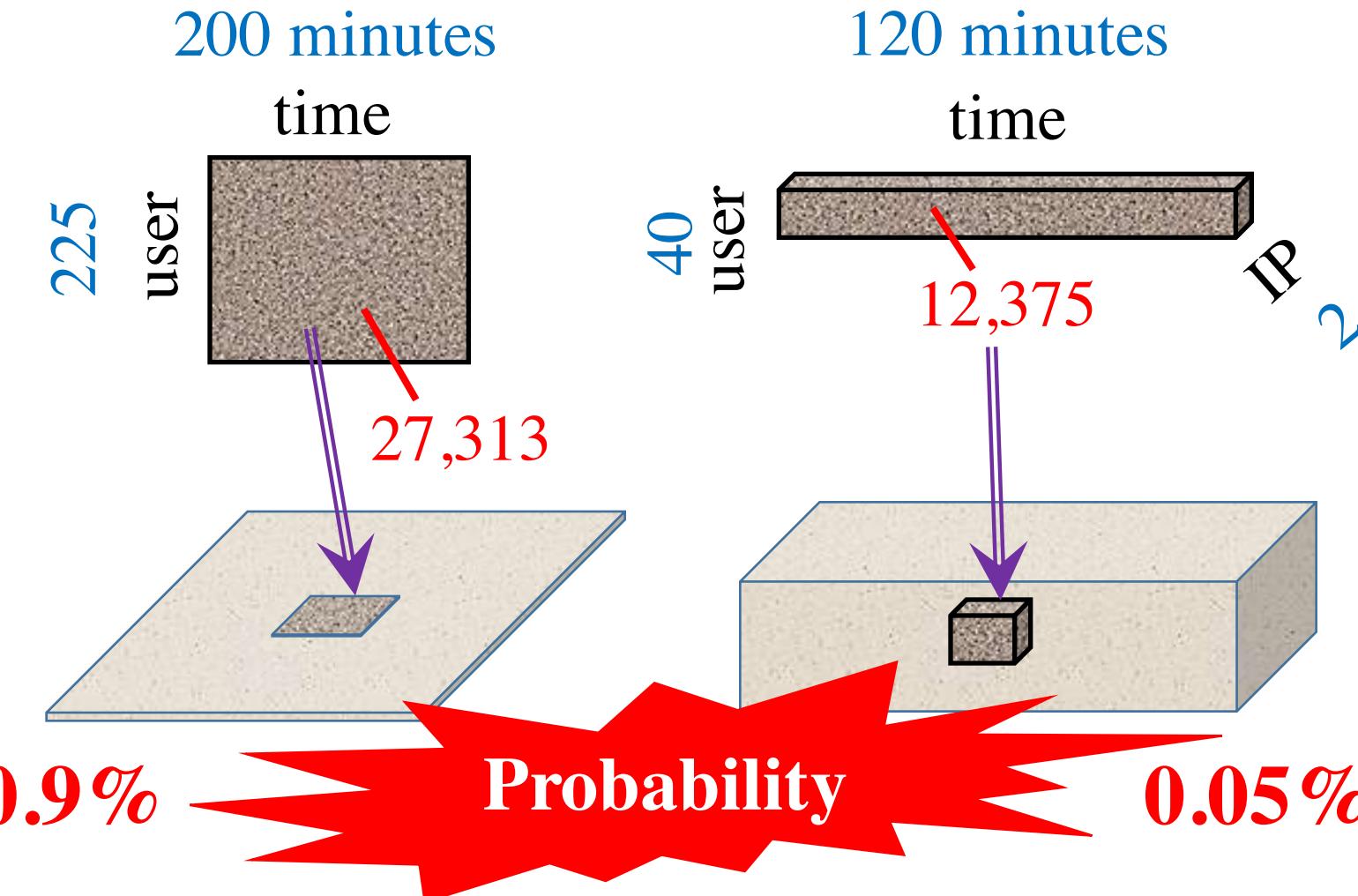
- New information (more modes) can only make our blocks more suspicious

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

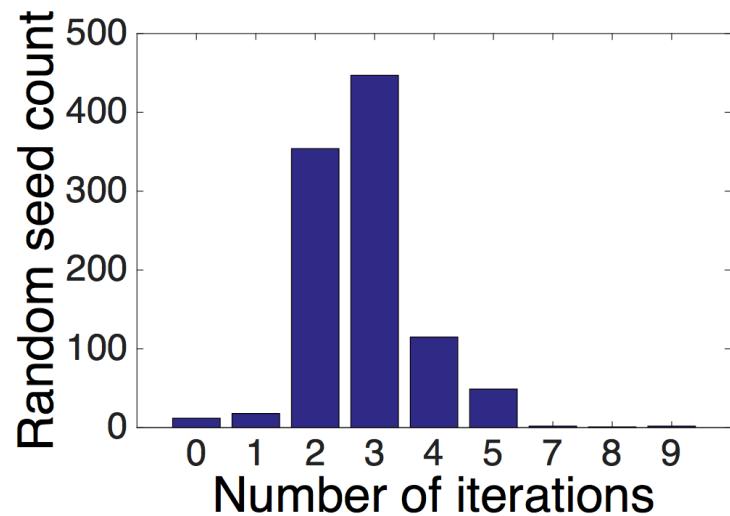
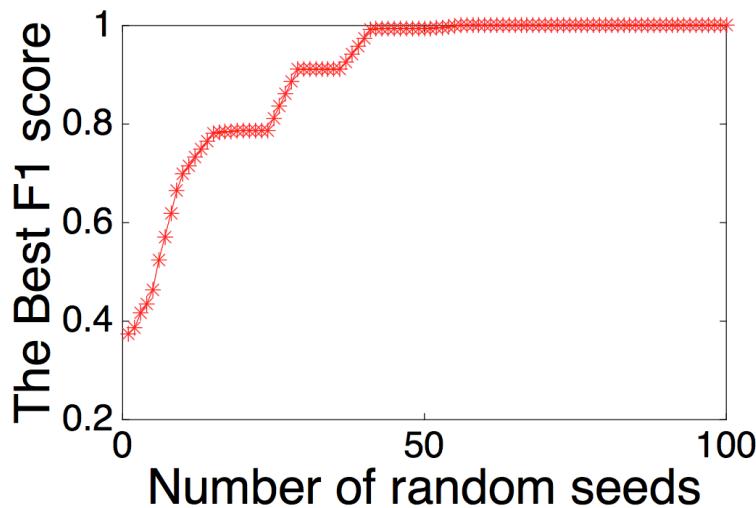
$$f(\mathbf{n}, c, \mathbf{N}, C) = c(\log \frac{c}{C} - 1) + C \prod_{i=1}^K \frac{n_i}{N_i} - c \sum_{i=1}^K \log \frac{n_i}{N_i}$$

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

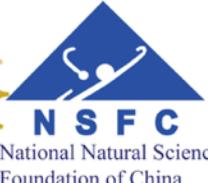


Summary

- ❑ 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



Acknowledgement



National Natural Science
Foundation of China



Carnegie
Mellon
University



Microsoft®
Research
微软亚洲研究院



47



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**