

**Mining Large Graphs and
Time Sequences:
Patterns, Anomalies, and Fraud
Detection**

Christos Faloutsos

CMU

Thank you!

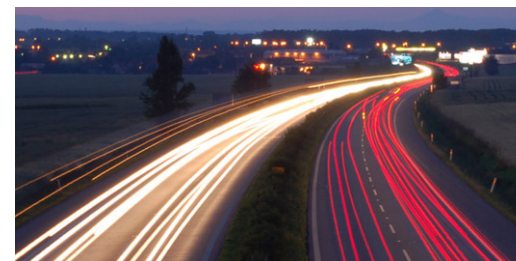
- Alkis Polyzotis



- Denise Olivera

Roadmap

- ➔ • Introduction – Motivation
 - Why study (big) graphs?
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
- Part#3: time sequences
- Conclusions

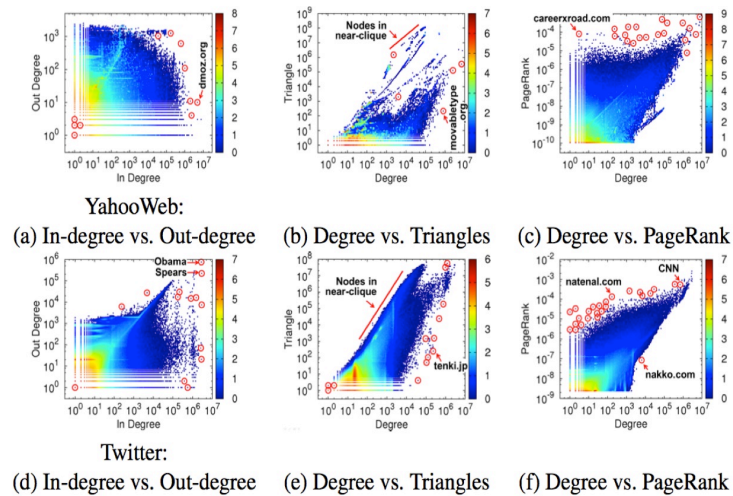


Graphs - why should we care?



>\$10B; ~1B users



Graphs - why should we care?



~1B nodes (web sites)
 ~6B edges (http links)
 'YahooWeb graph'

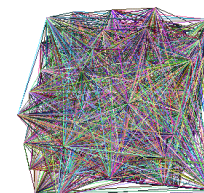
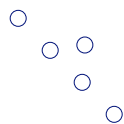
U Kang, Jay-Yoon Lee, Danai Koutra, and Christos Faloutsos. *Net-Ray: Visualizing and Mining Billion-Scale Graphs* PAKDD 2014, Tainan, Taiwan.

Graphs - why should we care?

- web-log ('blog') news propagation 
- computer network security: email/IP traffic and anomaly detection
- Recommendation systems 
-
- Many-to-many db relationship -> graph

Motivating problems

- P1: patterns? Fraud detection?



- P2: patterns in time-evolving graphs / tensors

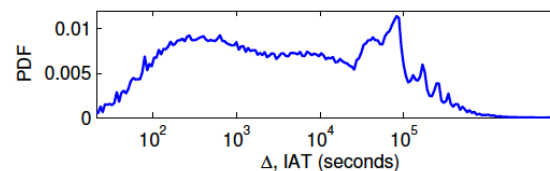
destination



source

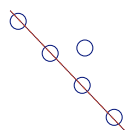
time

- P3: time sequences



Motivating problems

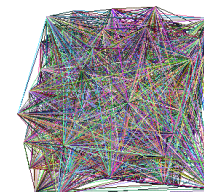
- P1: patterns? Fraud detection?



Patterns



anomalies



- P2: patterns in time-evolving graphs / tensors

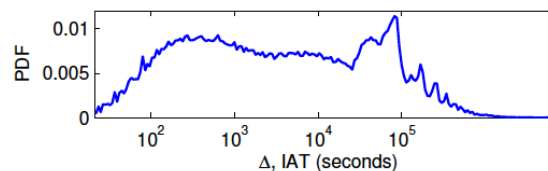
destination



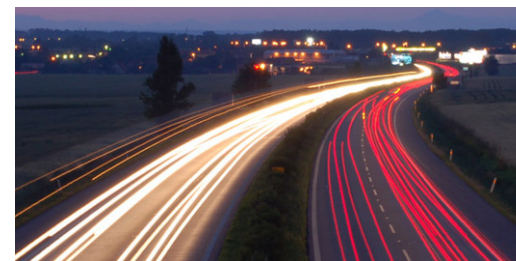
source

time

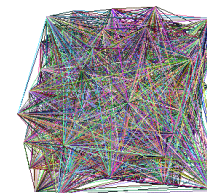
- P3: time sequences



Roadmap



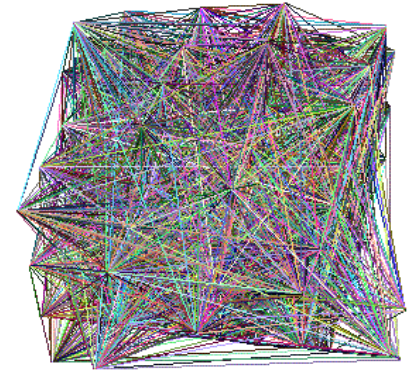
- Introduction – Motivation
 - Why study (big) graphs?
- ➔ • Part#1: Patterns & fraud detection
- Part#2: time-evolving graphs; tensors
- Conclusions



Part 1: Patterns, & fraud detection

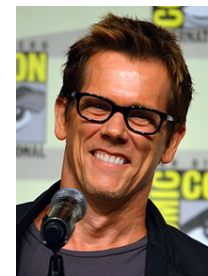
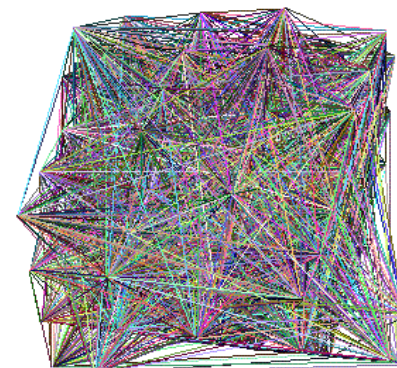
Laws and patterns

- Q1: Are real graphs random?



Laws and patterns

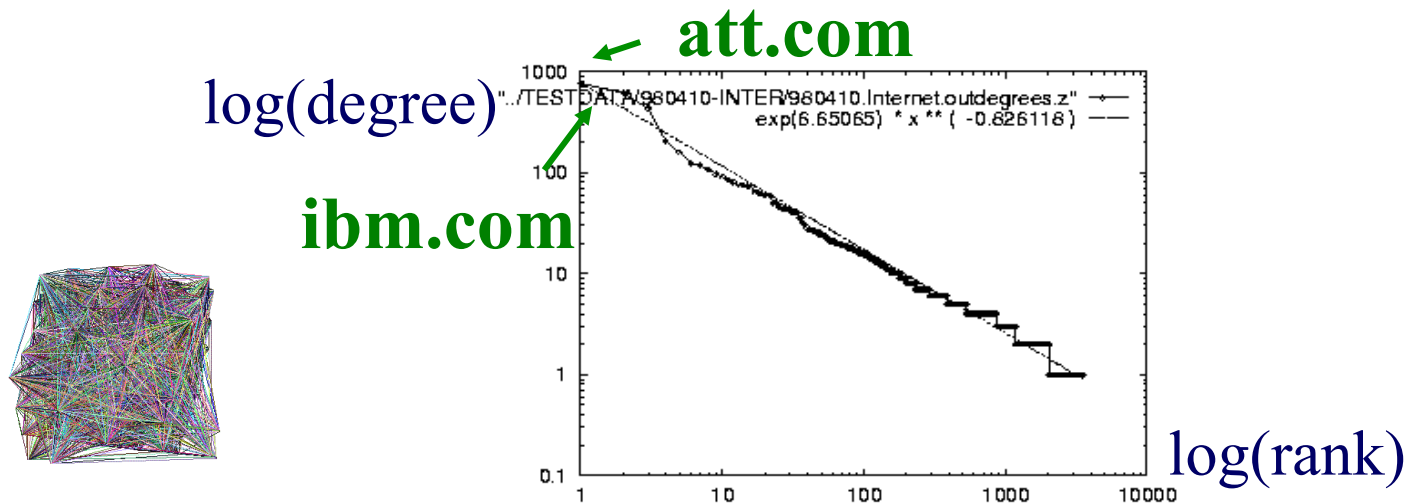
- Q1: Are real graphs random?
- A1: NO!!
 - Diameter ('6 degrees'; 'Kevin Bacon')
 - in- and out- degree distributions
 - other (surprising) patterns
- So, let's look at the data



Solution# S.1

- Power law in the degree distribution [Faloutsos x 3 SIGCOMM99]

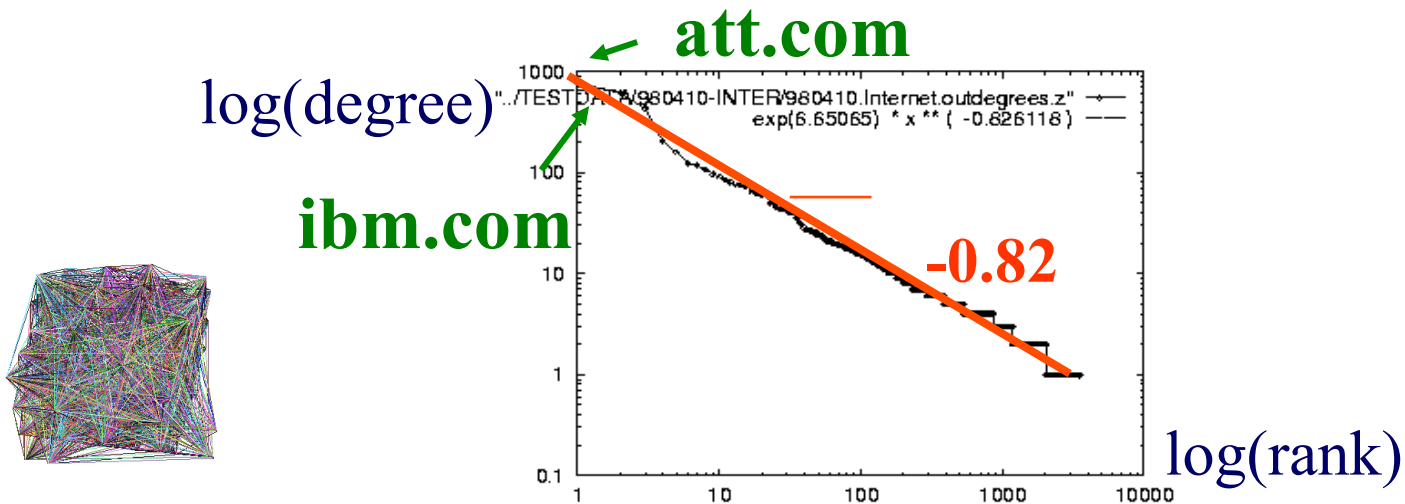
internet domains



Solution# S.1

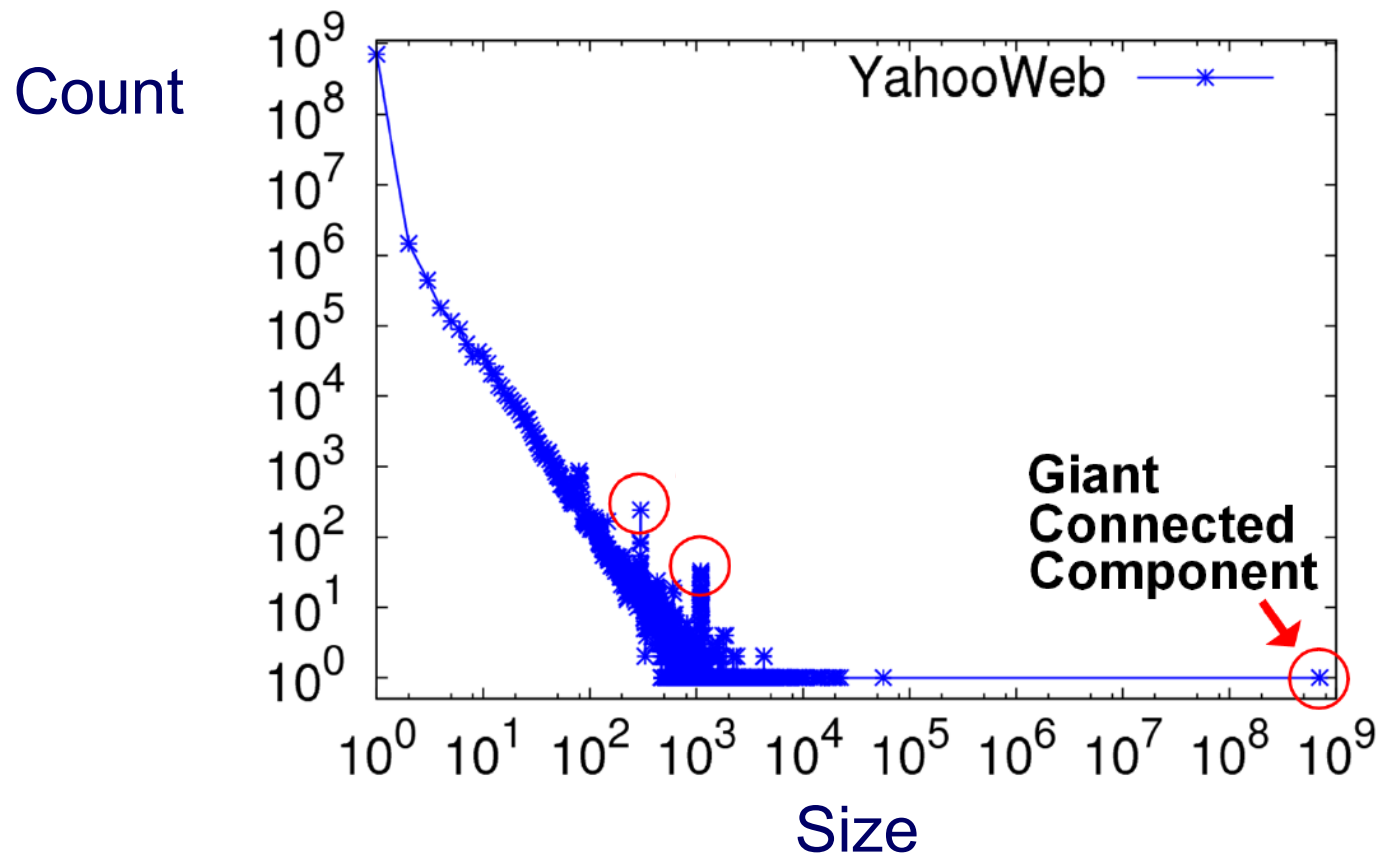
- Power law in the degree distribution [Faloutsos x 3 SIGCOMM99]

internet domains



S2: connected component sizes

- Connected Components – 4 observations:

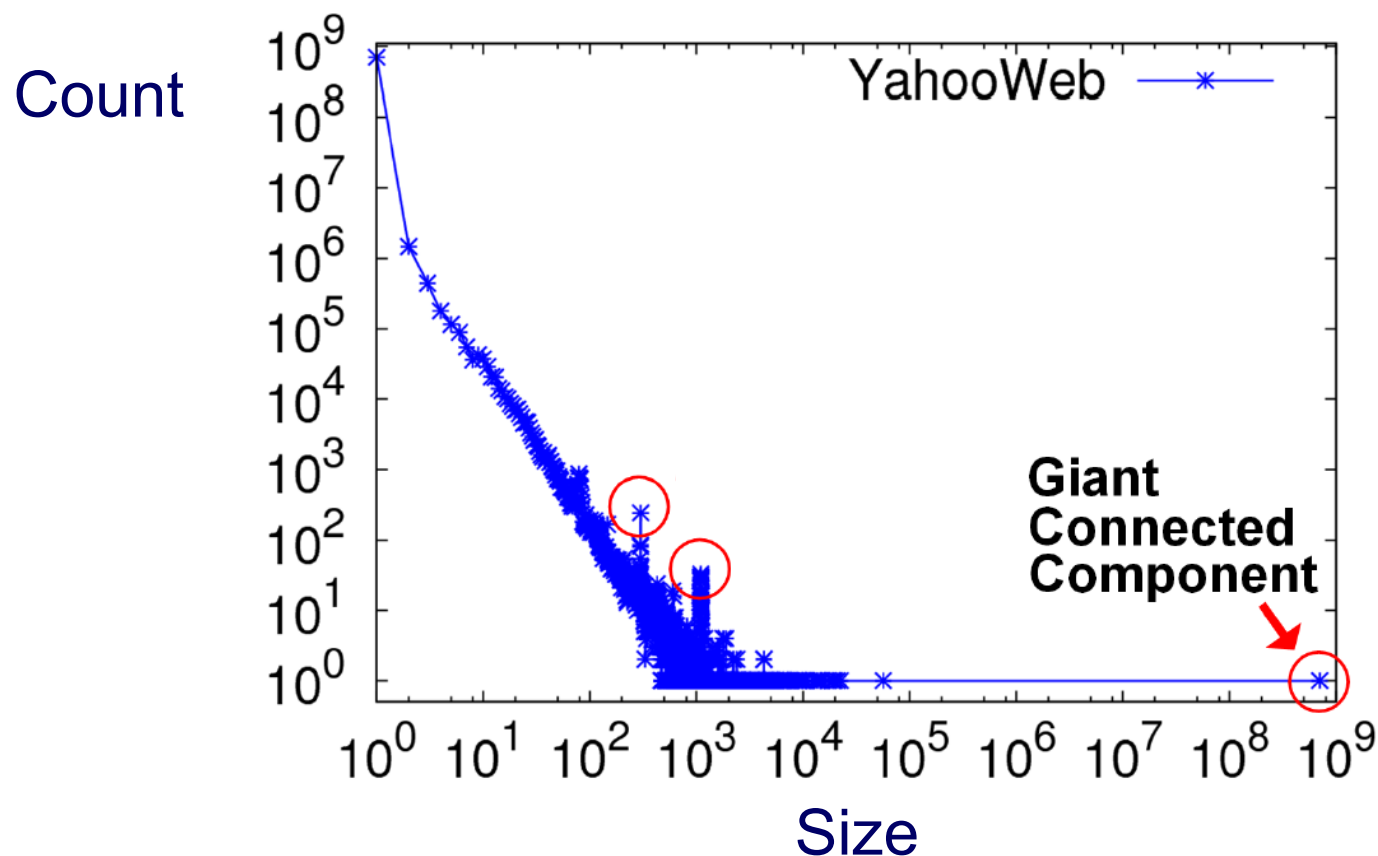


1.4B nodes
6B edges

S2: connected component sizes



- Connected Components

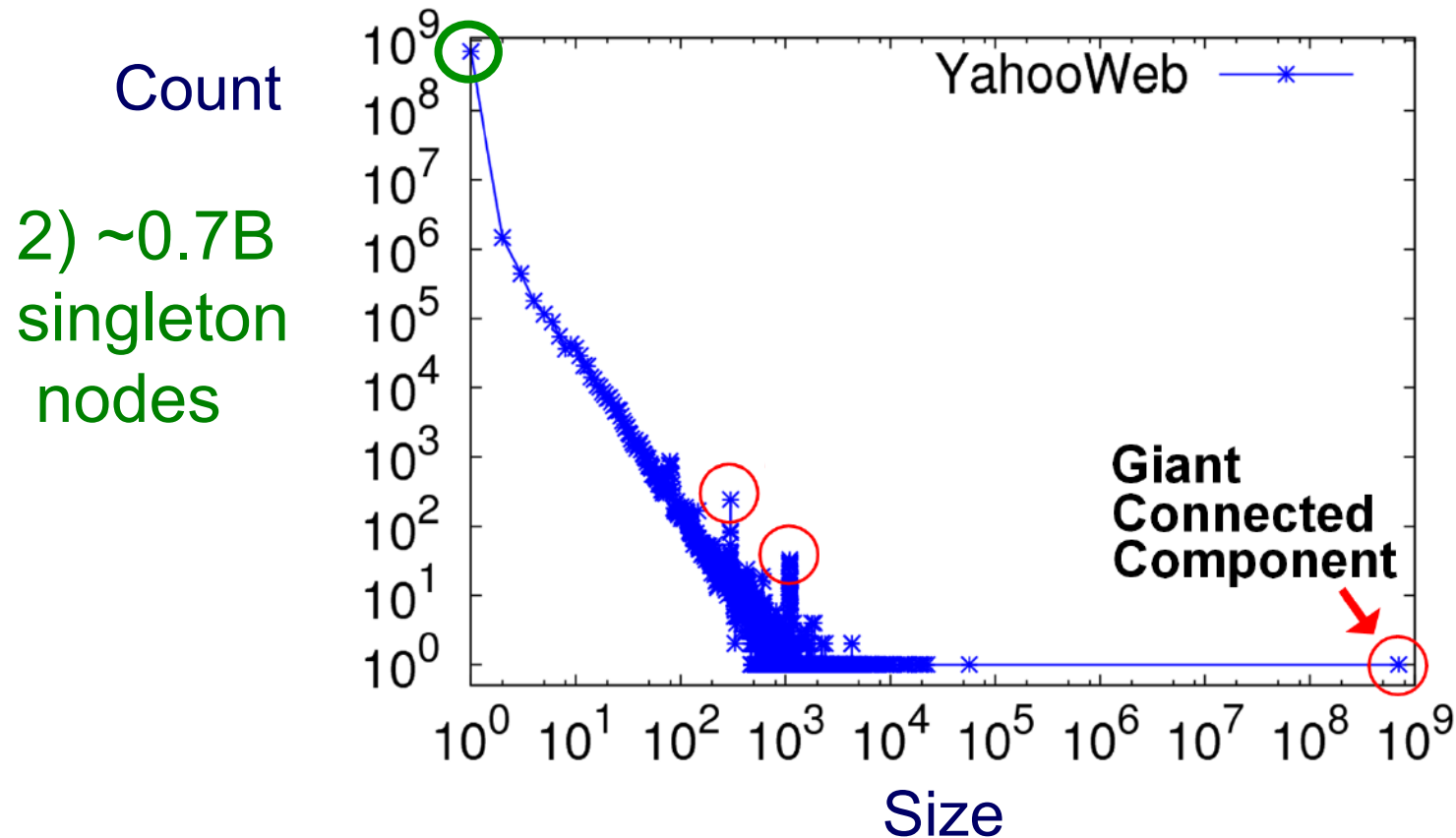


1) 10K x
larger
than next

S2: connected component sizes



- Connected Components

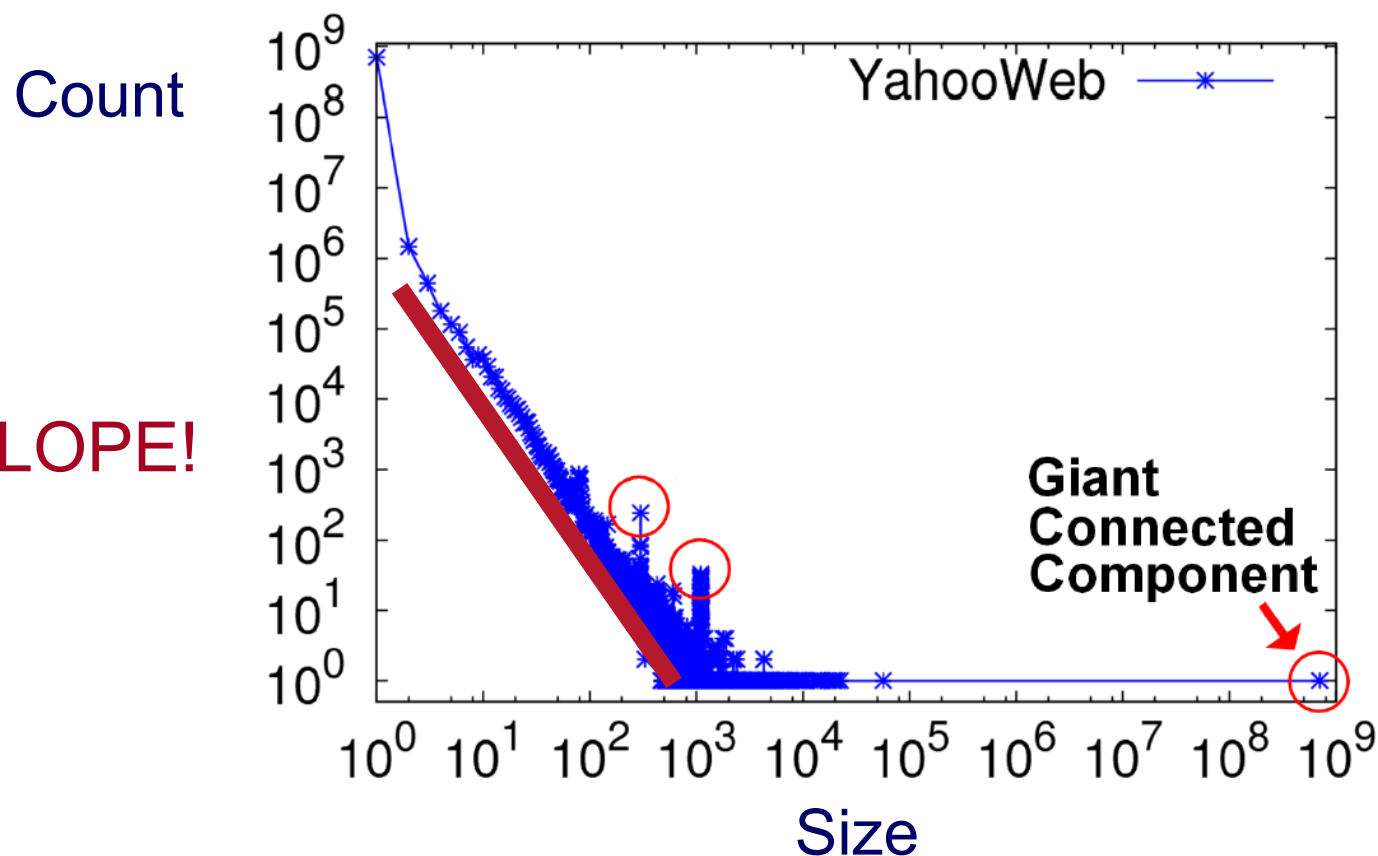


S2: connected component sizes



- Connected Components

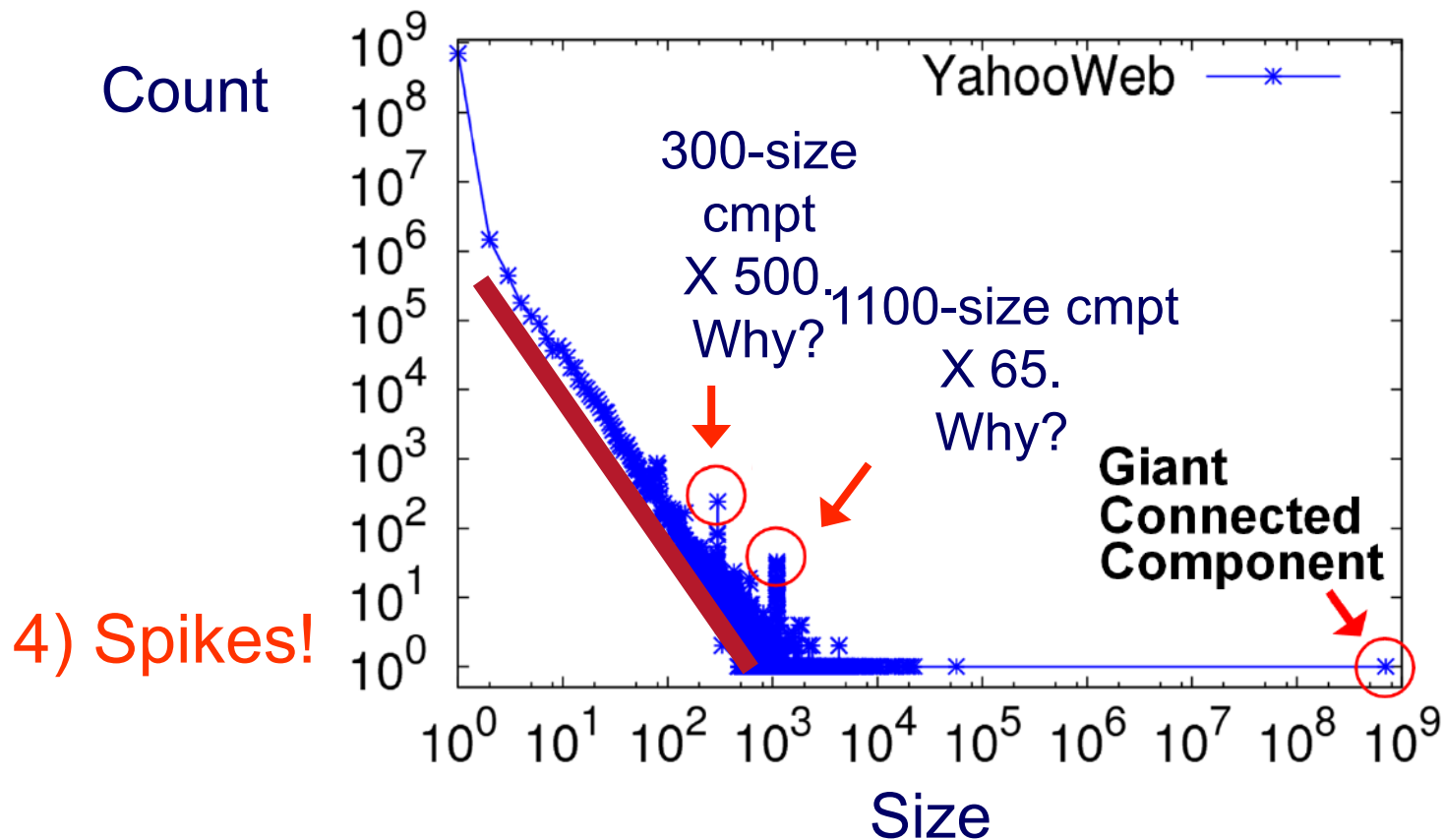
3) SLOPE!



S2: connected component sizes



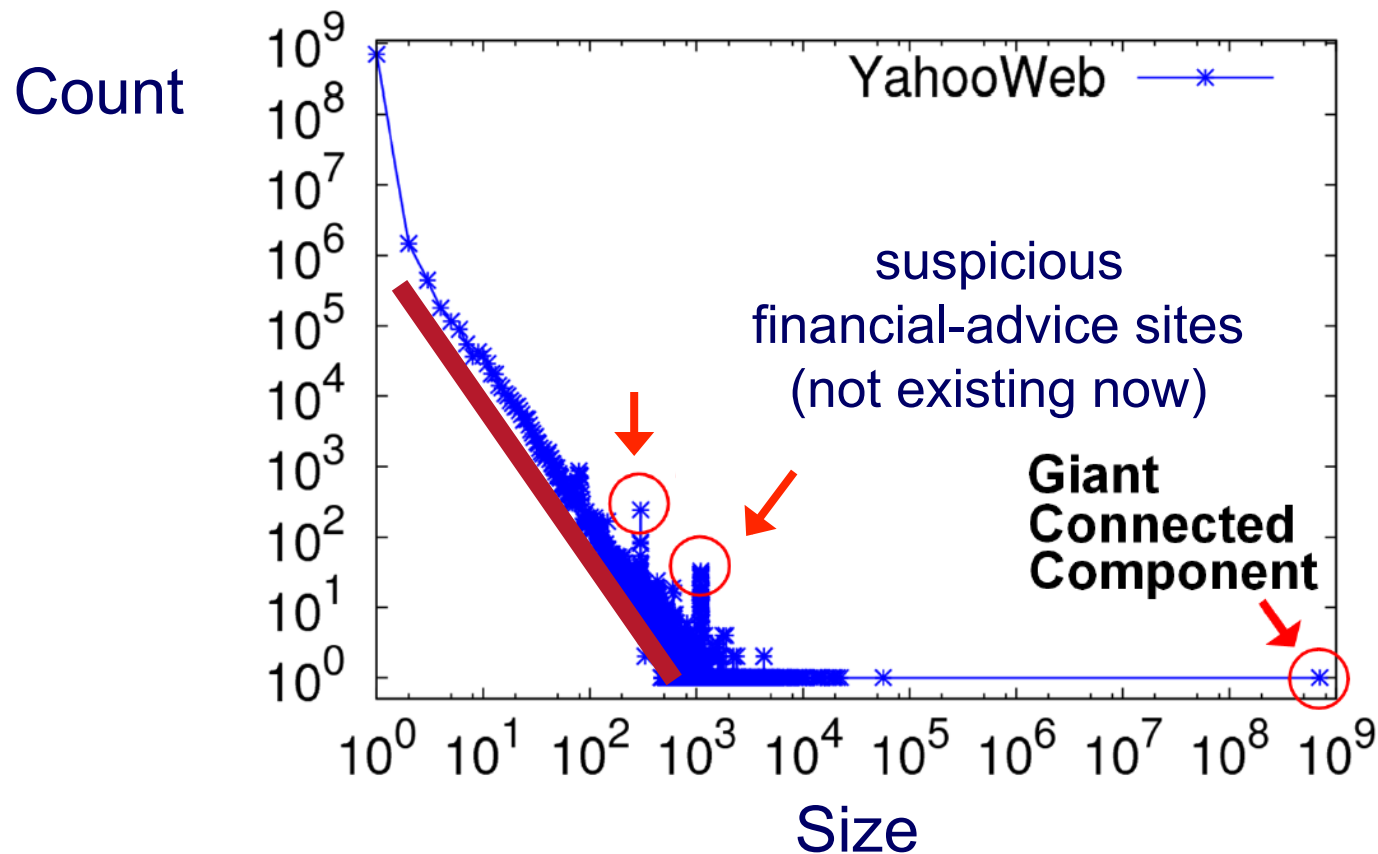
- Connected Components



S2: connected component sizes



- Connected Components

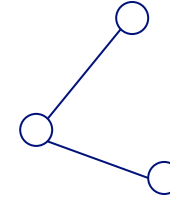


Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
 - ➔ – Patterns: Degree; Triangles
 - Anomaly/fraud detection
- Part#2: time-evolving graphs; tensors
- Part#3: time sequences
- Conclusions

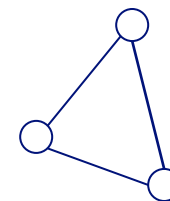


Solution# S.3: Triangle ‘Laws’

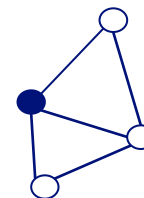


- Real social networks have a lot of triangles

Solution# S.3: Triangle ‘Laws’



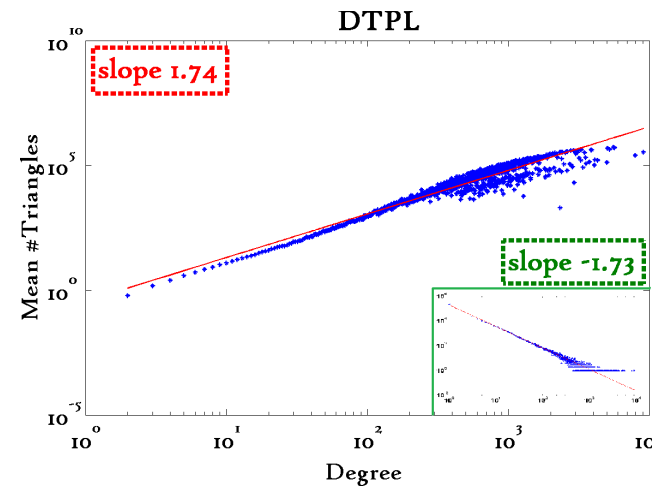
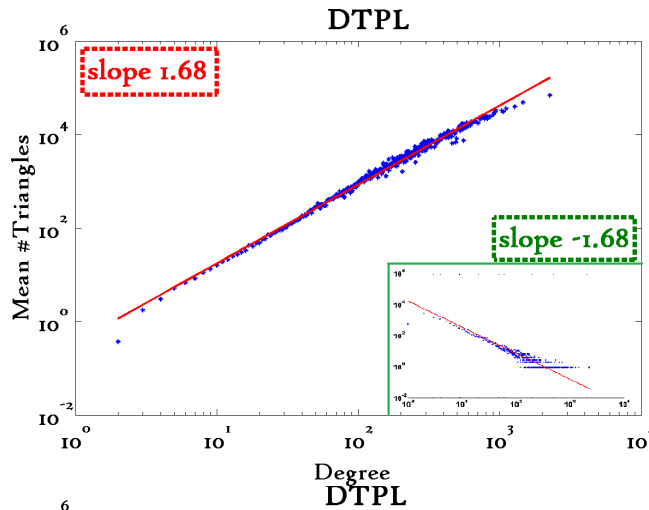
- Real social networks have a lot of triangles
 - Friends of friends are friends
- Any patterns?
 - 2x the friends, 2x the triangles ?



Triangle Law: #S.3

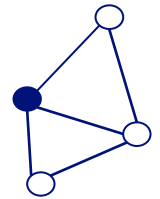
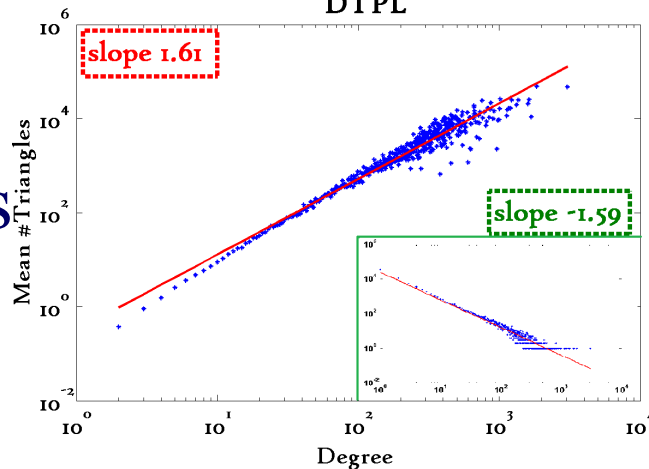
[Tsourakakis ICDM 2008]

Reuters



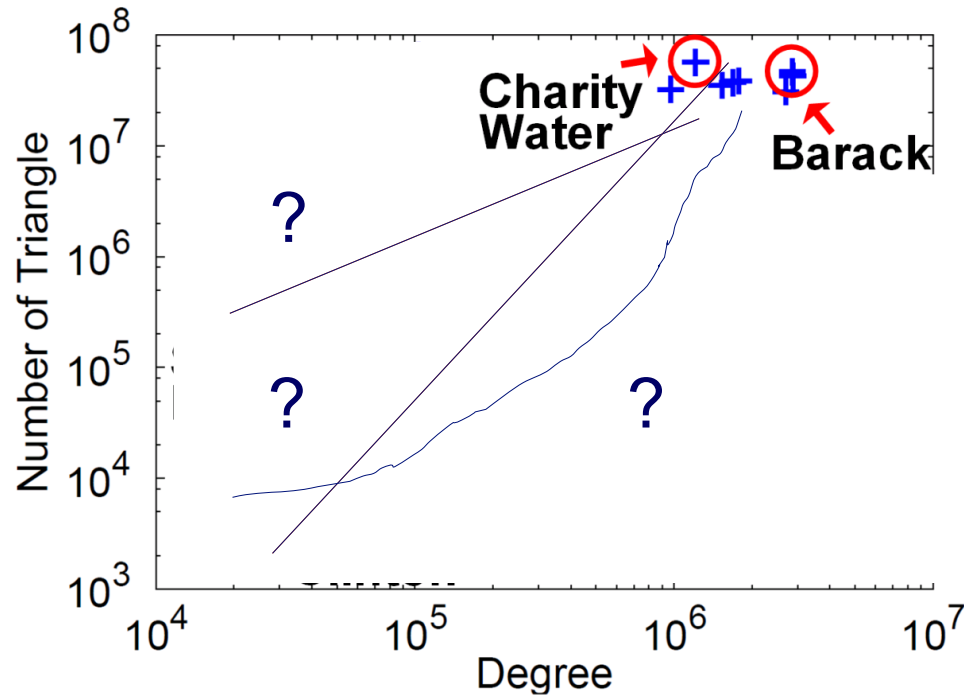
SN

Epinions



X-axis: degree
 Y-axis: mean # triangles
 n friends $\rightarrow \sim n^{1.6}$ triangles

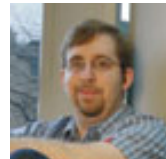
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

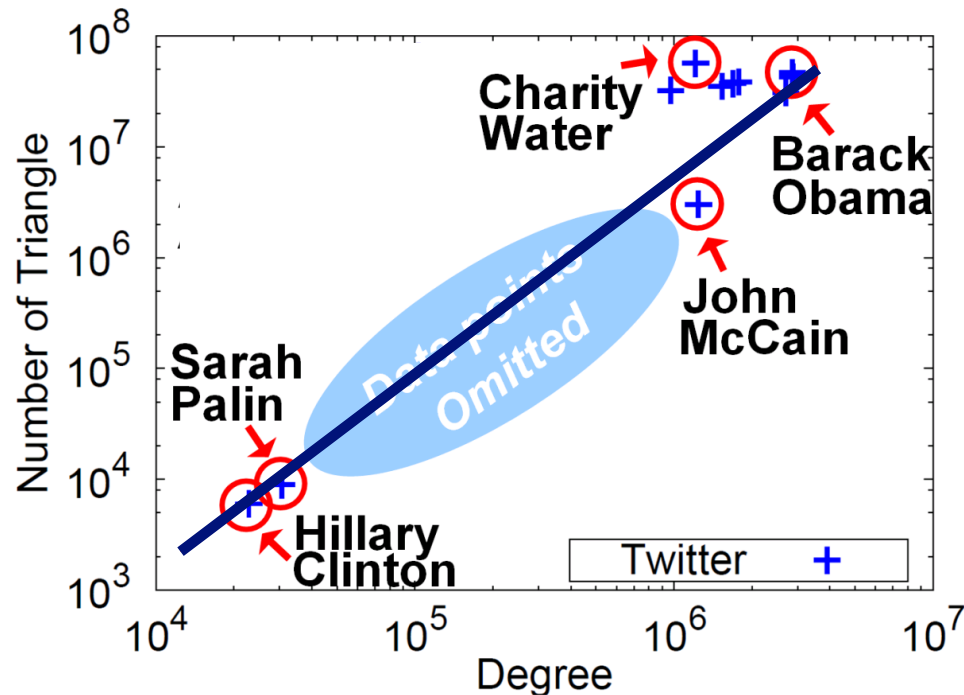
[U Kang, Brendan Meeder, +, PAKDD'11]

Google, Aug '16



(c) 2016, C. Faloutsos

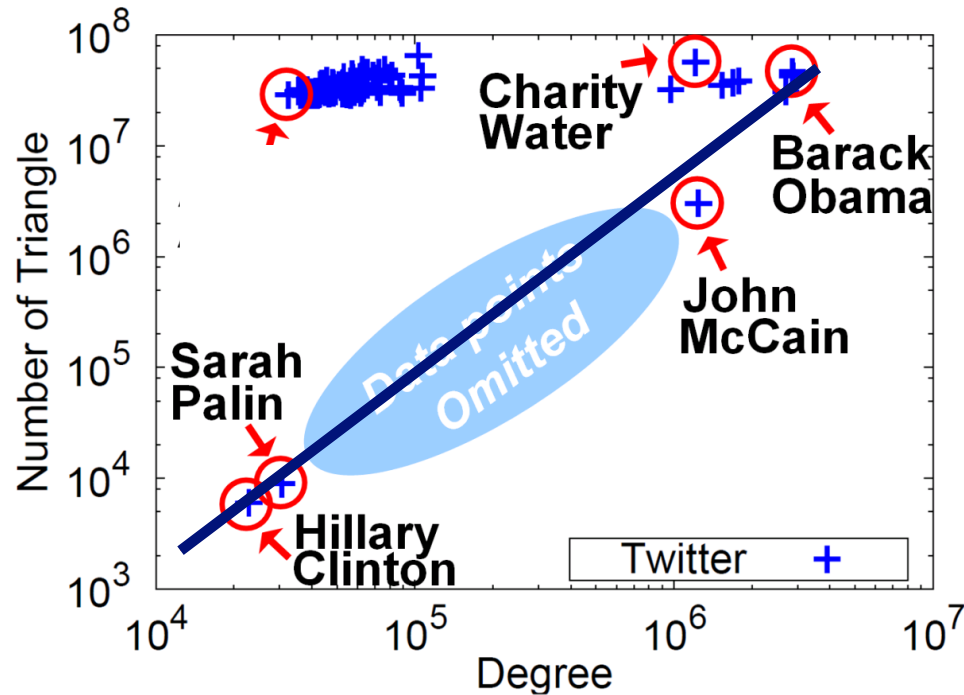
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

[U Kang, Brendan Meeder, +, PAKDD'11]

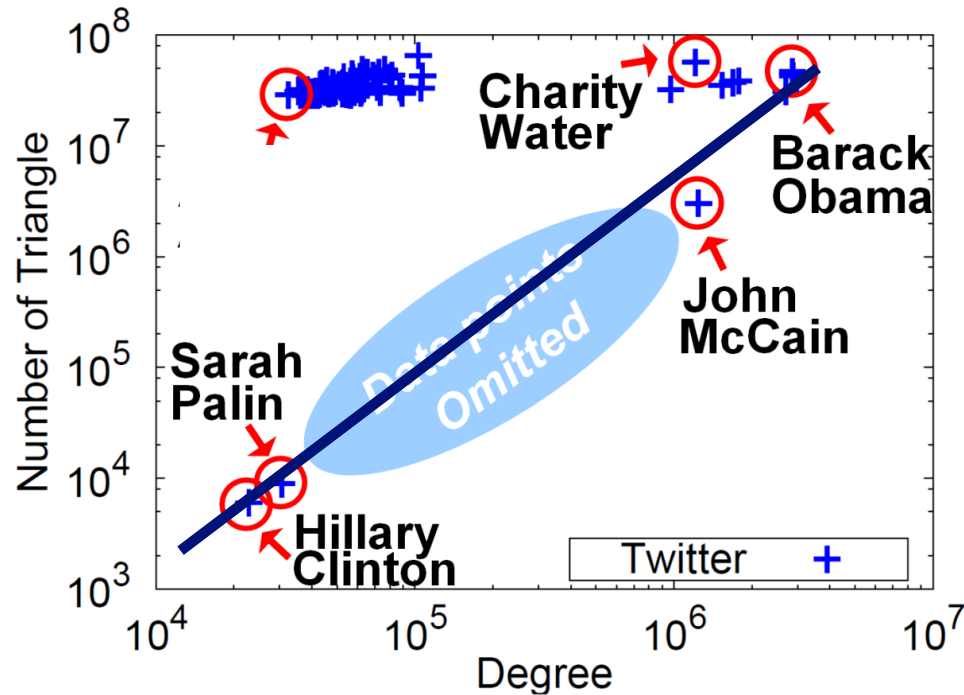
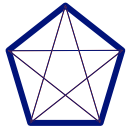
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

[U Kang, Brendan Meeder, +, PAKDD'11]

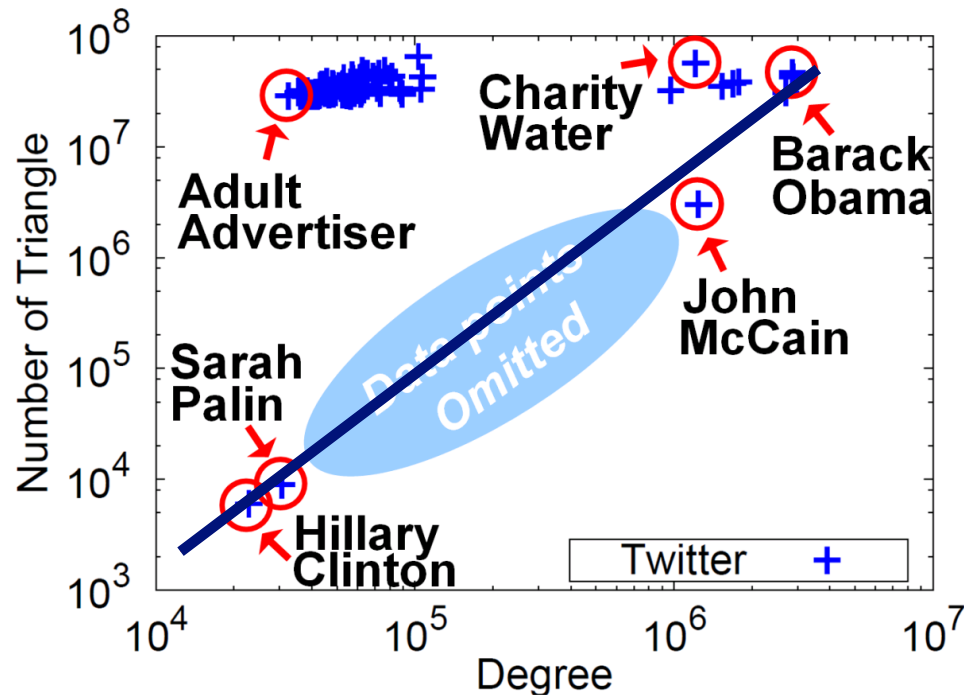
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

[U Kang, Brendan Meeder, +, PAKDD'11]

Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

[U Kang, Brendan Meeder, +, PAKDD'11]

MORE Graph Patterns

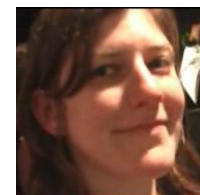
	Unweighted	Weighted
Static	<p>L01. Power-law degree distribution [Faloutsos et al. '99, Kleinberg et al. '99, Chakrabarti et al. '04, Newman '04]</p> <p>L02. Triangle Power Law (TPL) [Tsourakakis '08]</p> <p>L03. Eigenvalue Power Law (EPL) [Siganos et al. '03]</p> <p>L04. Community structure [Flake et al. '02, Girvan and Newman '02]</p>	<p>L10. Snapshot Power Law (SPL) [McGlohon et al. '08]</p>
Dynamic	<p>L05. Densification Power Law (DPL) [Leskovec et al. '05]</p> <p>L06. Small and shrinking diameter [Albert and Barabási '99, Leskovec et al. '05]</p> <p>L07. Constant size 2nd and 3rd connected components [McGlohon et al. '08]</p> <p>L08. Principal Eigenvalue Power Law (λ_1PL) [Akoglu et al. '08]</p> <p>L09. Bursty/self-similar edge/weight additions [Gomez and Santonja '98, Gribble et al. '98, Crovella and</p>	<p>L11. Weight Power Law (WPL) [McGlohon et al. '08]</p>

RTG: A Recursive Realistic Graph Generator using Random Typing Leman Akoglu and Christos Faloutsos. *PKDD'09*.

MORE Graph Patterns

	Unweighted	Weighted
Static	<p>L01. Power-law degree distribution [Faloutsos et al. '99, Kleinberg et al. '99, Chakrabarti et al. '04, Newman '04]</p> <p>L02. Triangle Power Law (TPL) [Tsourakakis '08]</p> <p>L03. Eigenvalue Power Law (EPL) [Siganos et al. '03]</p> <p>L04. Community structure [Flake et al. '02, Girvan and Newman '02]</p>	<p>L10. Snapshot Power Law (SPL) [McGlohon et al. '08]</p>
Dynamic	<p>L05. Densification Power Law (DPL) [Leskovec et al. '05]</p> <p>L06. Small and shrinking diameter [Albert and Barabási '99, Leskovec et al. '05]</p> <p>L07. Constant size 2nd and 3rd connected components [McGlohon et al. '08]</p> <p>L08. Principal Eigenvalue Power Law (λ_1PL) [Akoglu et al. '08]</p> <p>L09. Bursty/self-similar edge/weight additions [Gomez and Santonja '98, Gribble et al. '98, Crovella and Bestavros '99, McGlohon et al. '08]</p>	<p>L11. Weight Power Law (WPL) [McGlohon et al. '08]</p>

- Mary McGlohon, Leman Akoglu, Christos Faloutsos. *Statistical Properties of Social Networks*. in "Social Network Data Analytics" (Ed.: Charu Aggarwal)

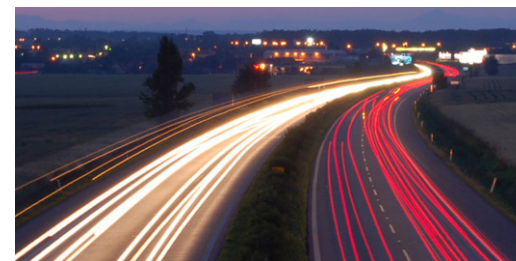


- Deepayan Chakrabarti and Christos Faloutsos, [*Graph Mining: Laws, Tools, and Case Studies*](#) Oct. 2012, Morgan Claypool.



Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
 - Patterns
 - – Anomaly / fraud detection
 - Spectral methods ('fBox')
 - Belief Propagation
- Part#2: time-evolving graphs; tensors
- Conclusions



Problem: Social Network Link Fraud

Target: find “stealthy” attackers missed by other algorithms

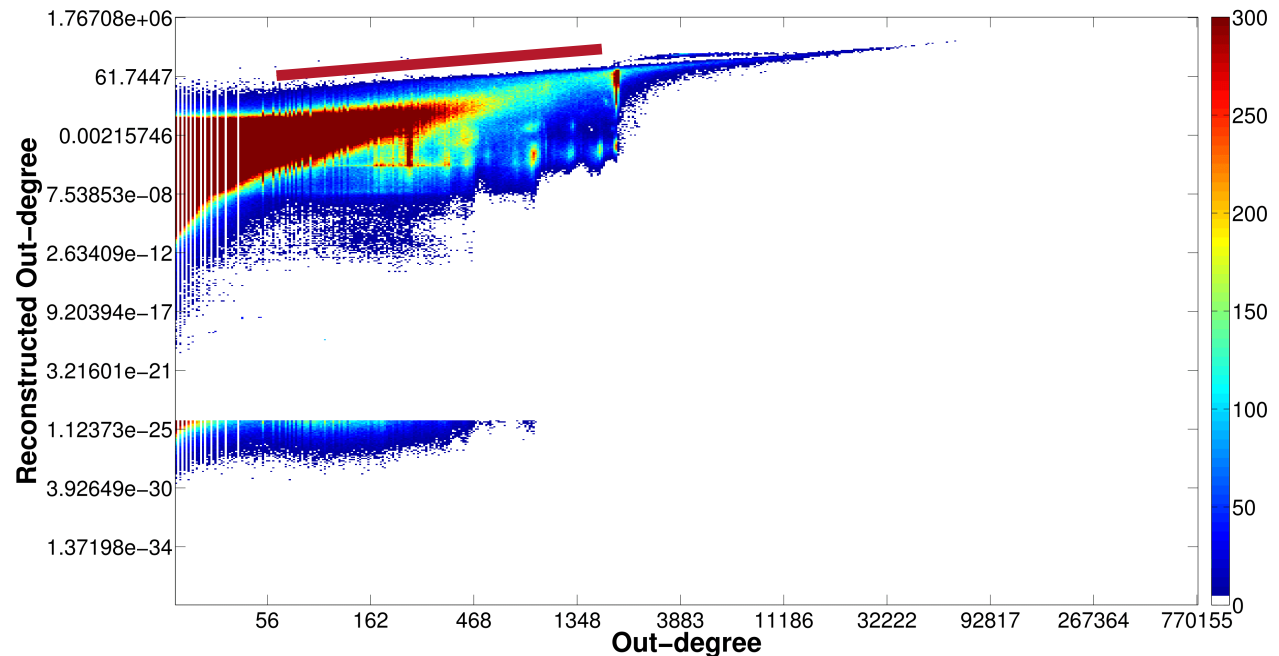


Clique

41.7M nodes
1.5B edges



Bipartite
core



Problem: Social Network Link Fraud

Target: find “stealthy” attackers missed by other algorithms



Lekan Olawole Lowe @loweinc

26 Jul 09

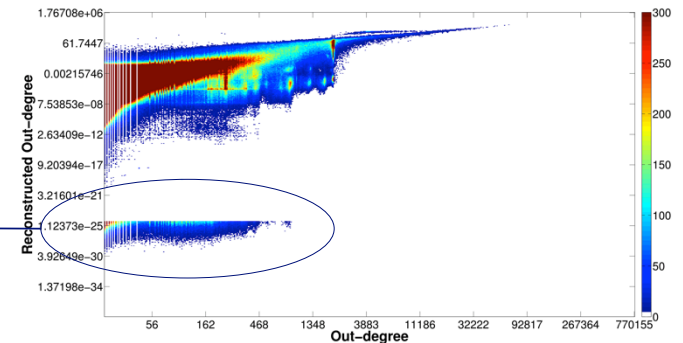
Sign up free and Get 400 followers a day using <http://tweeteradder.com>



Lekan Olawole Lowe @loweinc

26 Jul 09

Get 400 followers a day using <http://www.tweeterfollow.com>



Takeaway: use *reconstruction error* between true/latent representation!



Neil Shah, Alex Beutel, Brian Gallagher and Christos Faloutsos. *Spotting Suspicious Link Behavior with fBox: An Adversarial Perspective*. ICDM 2014, Shenzhen, China.

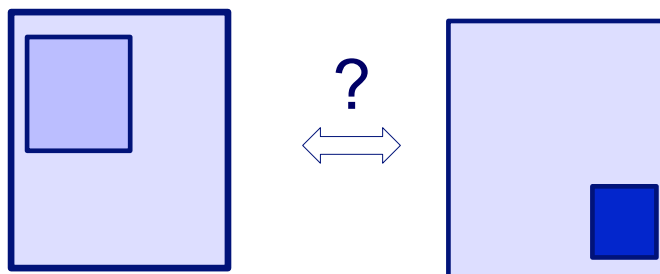
Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
 - Patterns
 - Anomaly / fraud detection
 - CopyCatch
 - Spectral methods ('fBox', **suspiciousness**)
 - Belief Propagation
- Part#2: time-evolving graphs; tensors
- Conclusions

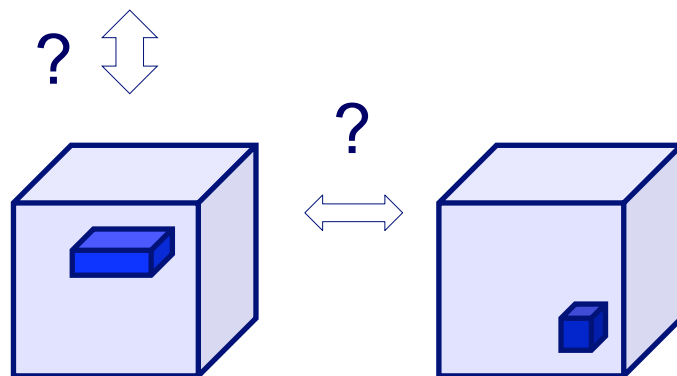


Suspicious Patterns in Event Data

2-modes



n -modes



A General Suspiciousness Metric for Dense Blocks in Multimodal Data, Meng Jiang, Alex Beutel, Peng Cui, Bryan Hooi, Shiqiang Yang, and Christos Faloutsos, *ICDM*, 2015.

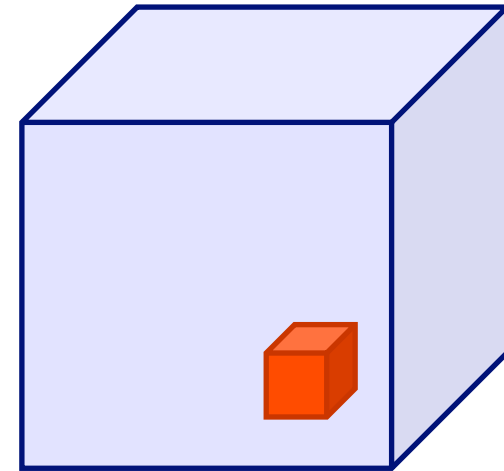
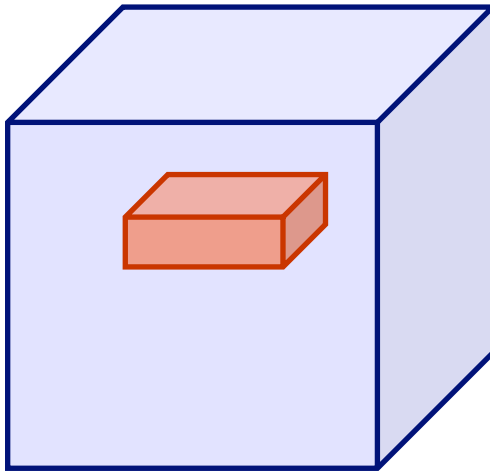
Suspicious Patterns in Event Data

Which is more suspicious?

20,000 Users
Retweeting same 20 tweets
6 times each
All in 10 hours

↔
↔
vs.
↔

225 Users
Retweeting same 1 tweet
15 times each
All in 3 hours
All from 2 IP addresses

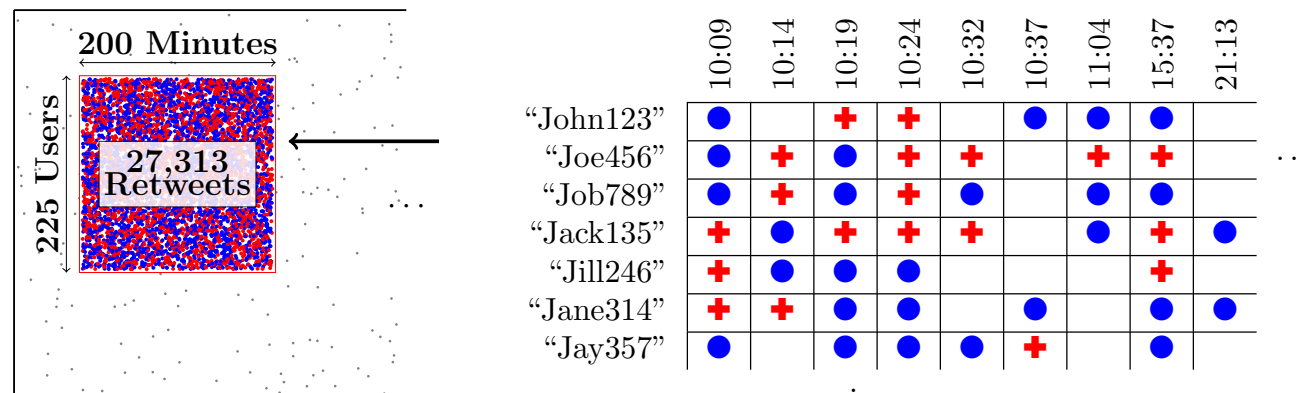


Google

Answer: volume * $D_{KL}(p || p_{background})$

37

Suspicious Patterns in Event Data



Retweeting: "Galaxy Note Dream Project: Happy Happy Life Traveling the World"

	#	User \times tweet \times IP \times minute	Mass c	Suspiciousness
CROSSPOT	1	$14 \times 1 \times 2 \times 1,114$	41,396	1,239,865
	2	$225 \times 1 \times 2 \times 200$	27,313	777,781
	3	$8 \times 2 \times 4 \times 1,872$	17,701	491,323
HOSVD	1	$24 \times 6 \times 11 \times 439$	3,582	131,113
	2	$18 \times 4 \times 5 \times 223$	1,942	74,087
	3	$14 \times 2 \times 1 \times 265$	9,061	381,211

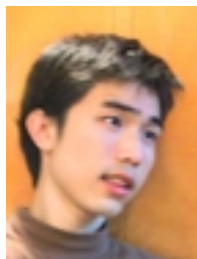
Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
 - Patterns
 - Anomaly / fraud detection
 - Spectral methods ('fBox')
 - High-density sub-matrices
 - Belief propagation

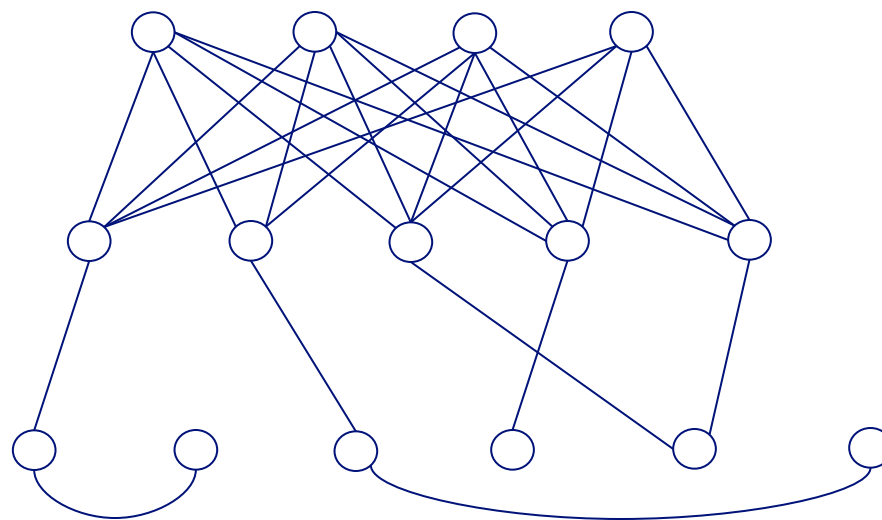


- Part#2: time-evolving graphs; tensors
- Part#3: time sequences
- Conclusions

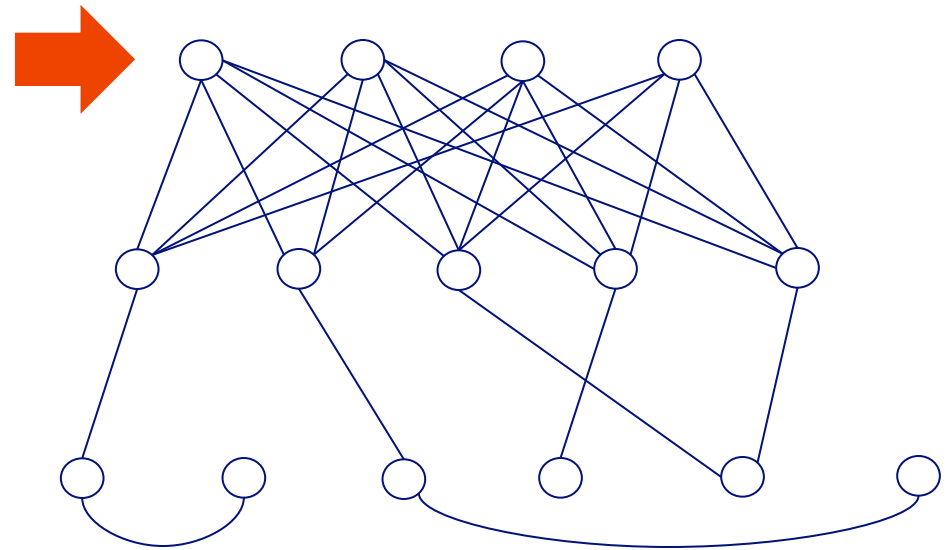
E-bay Fraud detection



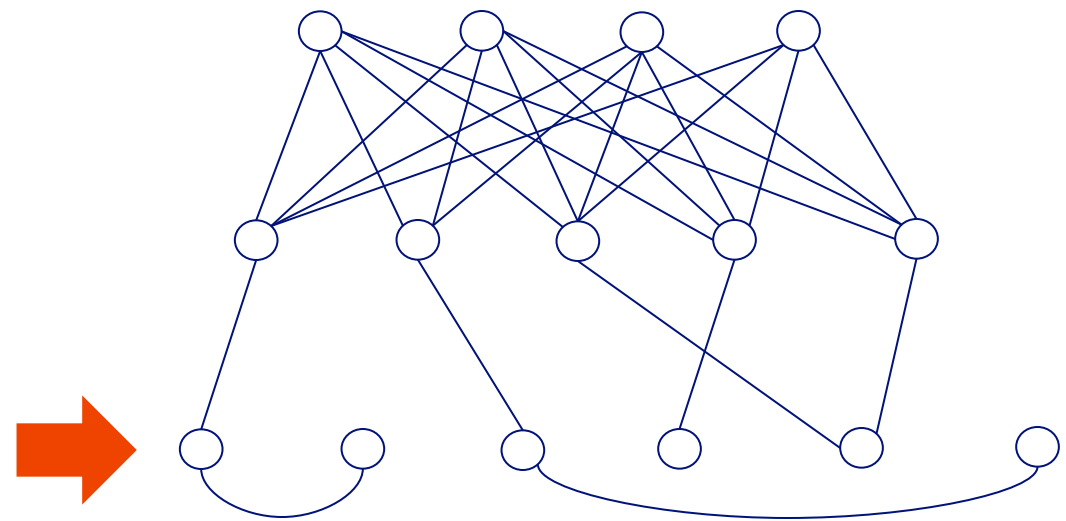
w/ Polo Chau &
Shashank Pandit, CMU
[www'07]



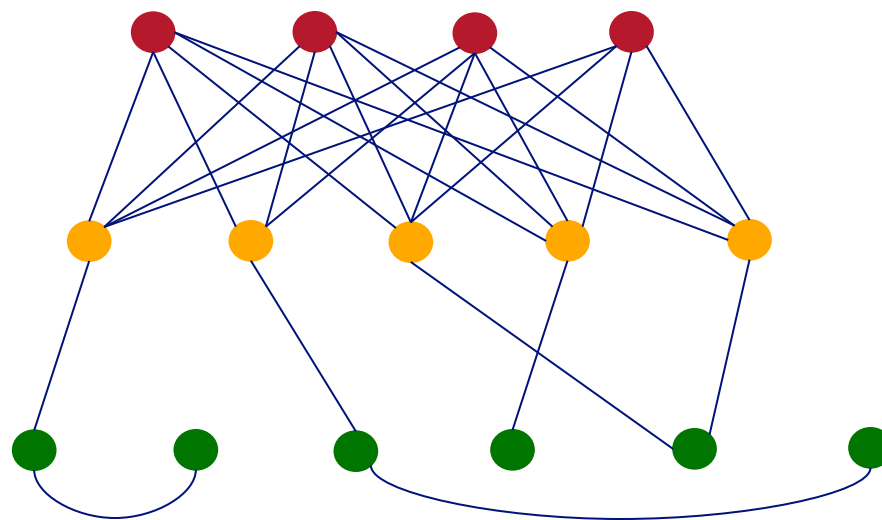
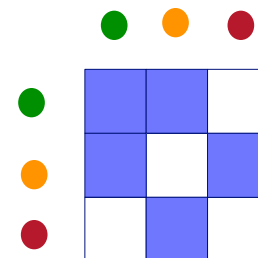
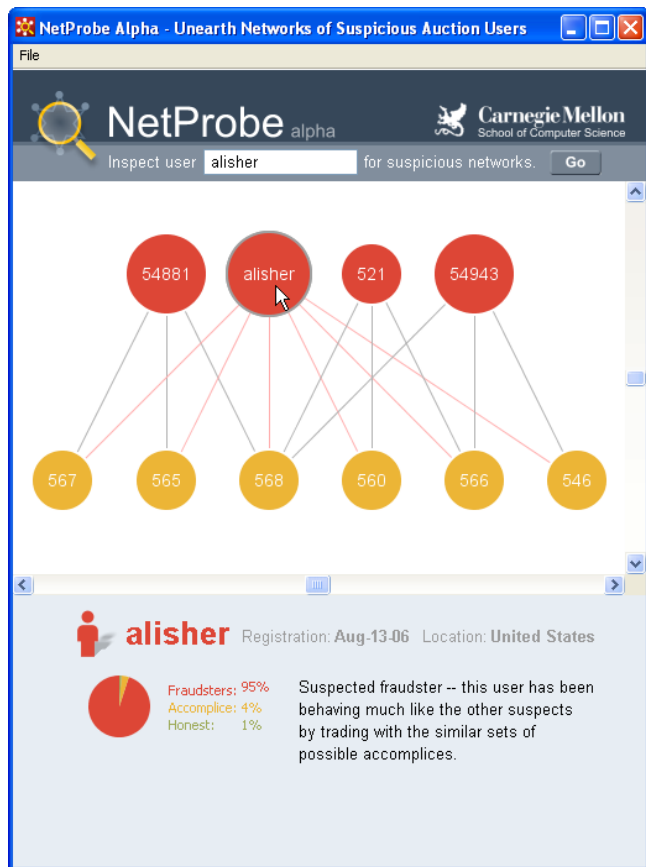
E-bay Fraud detection



E-bay Fraud detection



E-bay Fraud detection - NetProbe



Popular press



The Washington Post

Los Angeles Times

And less desirable attention:

- E-mail from ‘Belgium police’ (‘copy of your code?’)

Roadmap

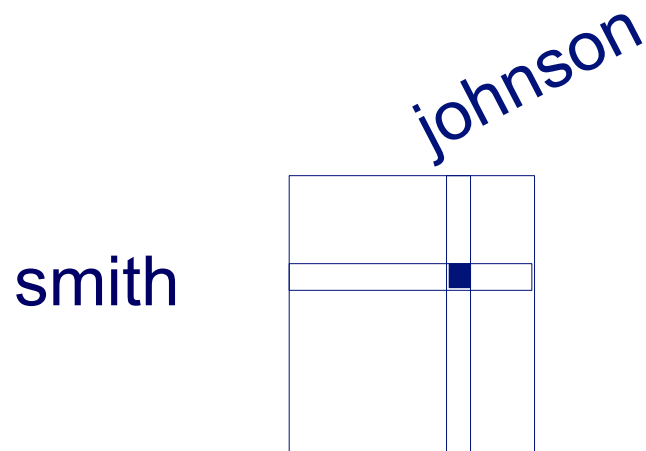


- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
 - ➔ – P2.1: time-evolving graphs
 - [P2.2: with side information (‘coupled’ M.T.F.)
 - Speed]
- Part#3: time sequences
- Conclusions

Part 2: Time evolving graphs; tensors

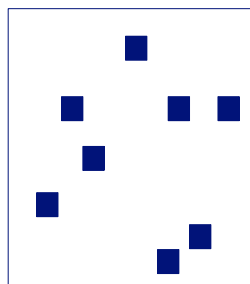
Graphs over time -> tensors!

- Problem #2.1:
 - Given who calls whom, and when
 - Find patterns / anomalies



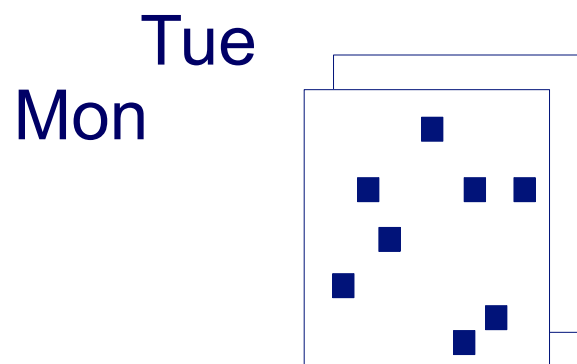
Graphs over time -> tensors!

- Problem #2.1:
 - Given who calls whom, and when
 - Find patterns / anomalies



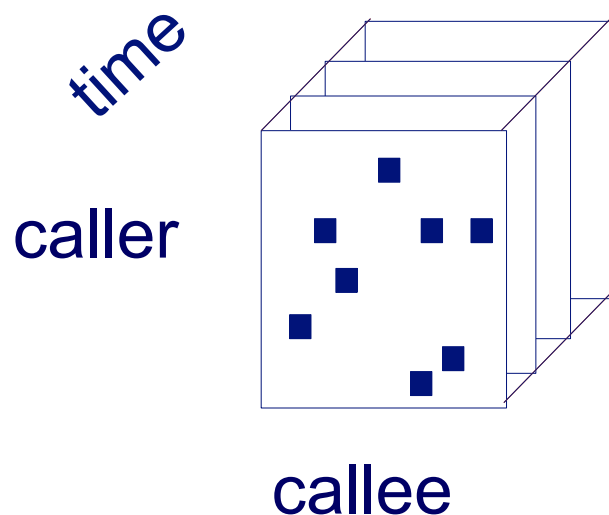
Graphs over time -> tensors!

- Problem #2.1:
 - Given who calls whom, and when
 - Find patterns / anomalies



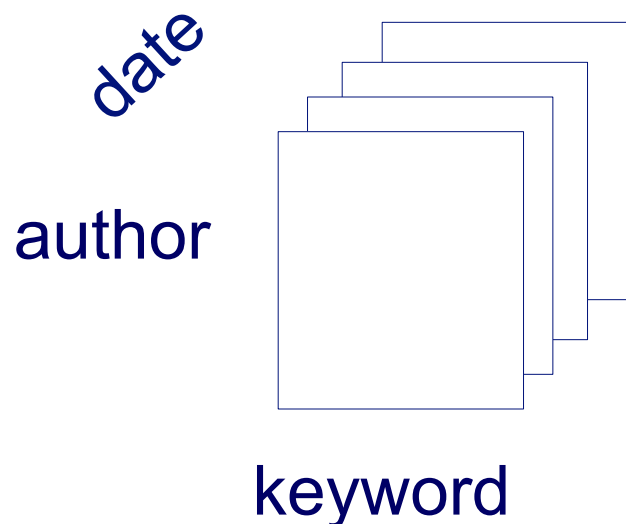
Graphs over time -> tensors!

- Problem #2.1:
 - Given who calls whom, and when
 - Find patterns / anomalies



Graphs over time -> tensors!

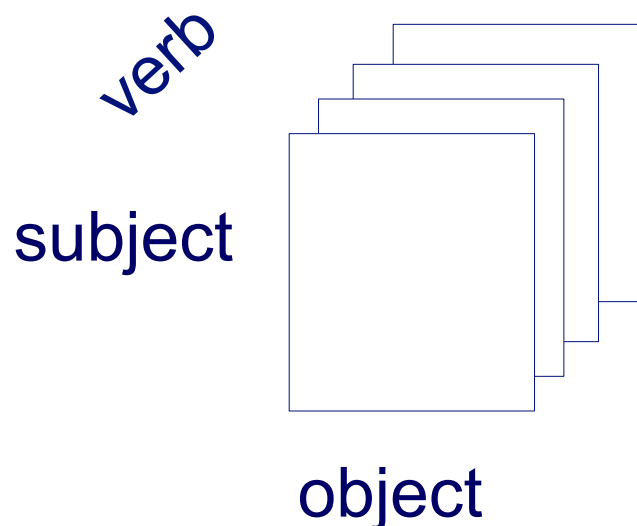
- Problem #2.1':
 - Given author-keyword-date
 - Find patterns / anomalies



MANY more settings,
with >2 'modes'

Graphs over time -> tensors!

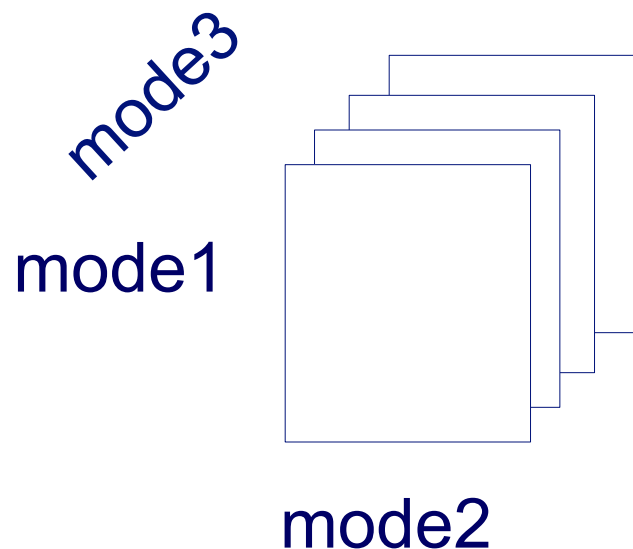
- Problem #2.1’’:
 - Given subject – verb – object facts
 - Find patterns / anomalies



MANY more settings,
with >2 ‘modes’

Graphs over time -> tensors!

- Problem #2.1''':
 - Given <triplets>
 - Find patterns / anomalies



MANY more settings,
with >2 'modes'
(and 4, 5, etc modes)

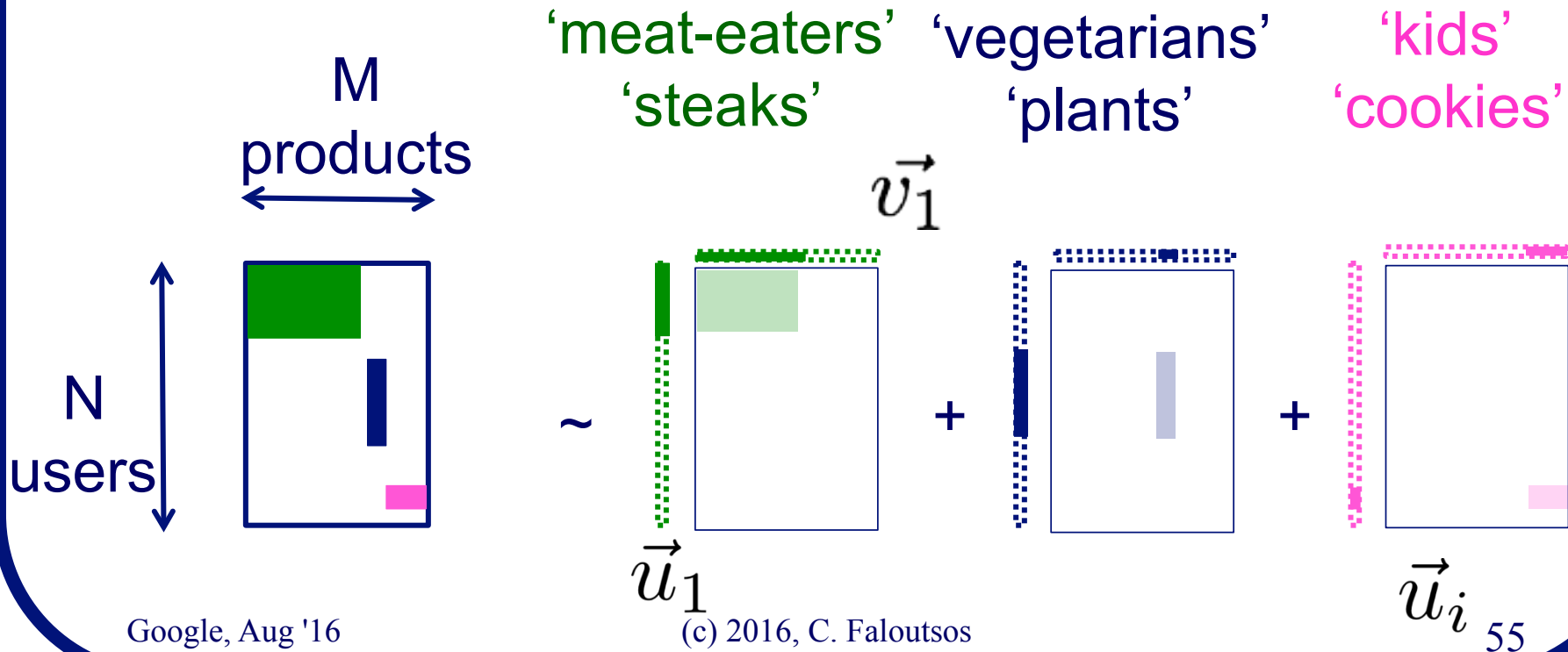
Roadmap



- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
 - ➔ – P2.1: time-evolving graphs
 - [P2.2: with side information ('coupled' M.T.F.)
 - Speed]
- Conclusions

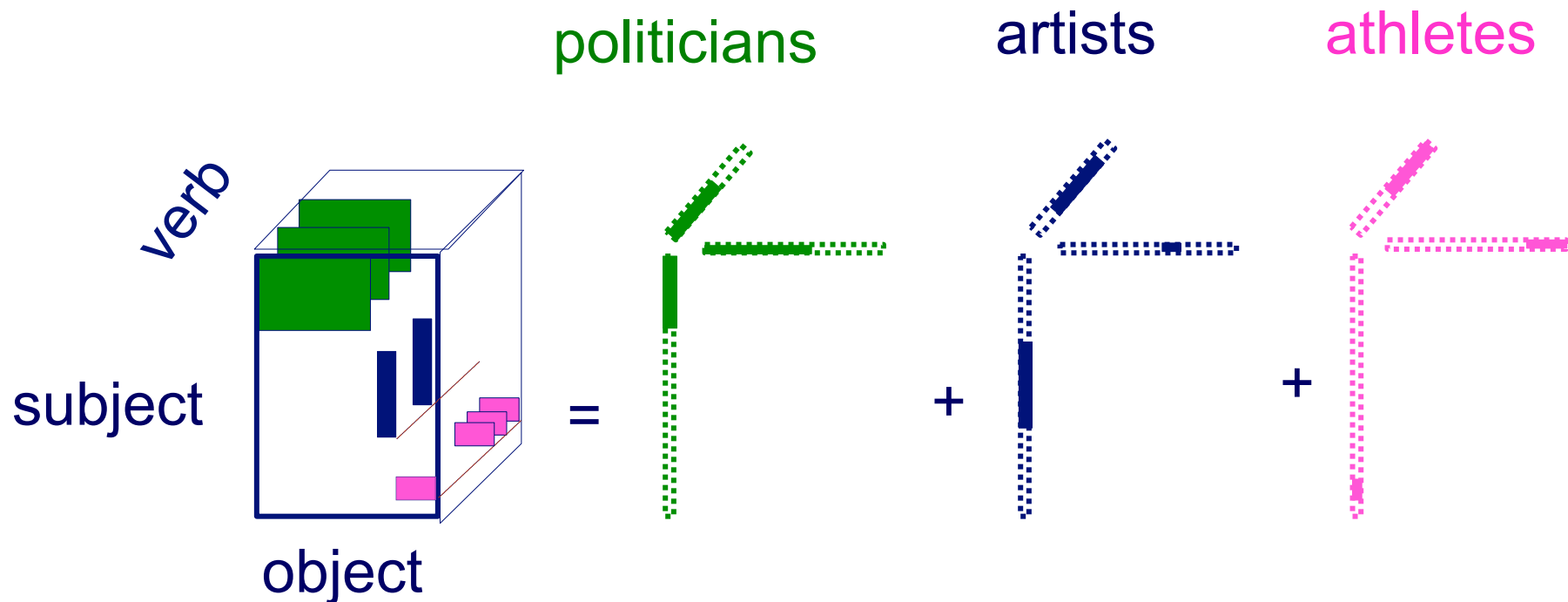
Answer to both: tensor factorization

- Recall: (SVD) matrix factorization: finds blocks



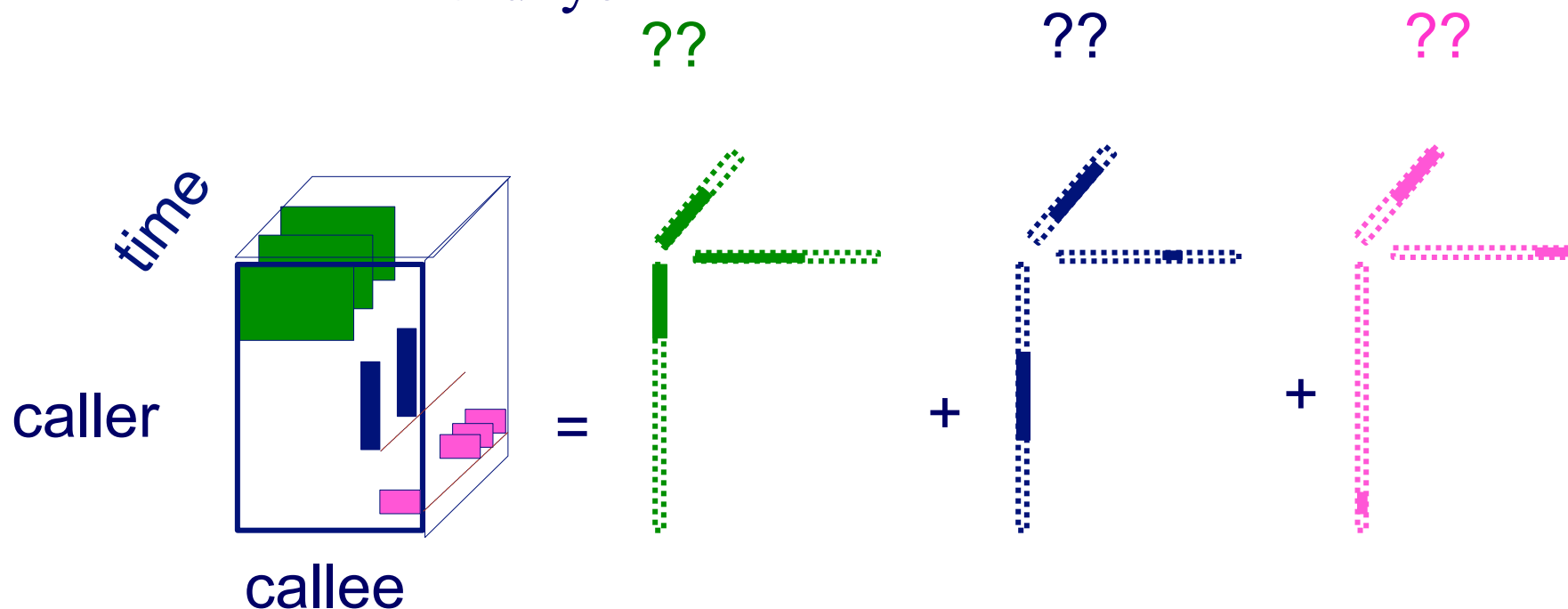
Answer to both: tensor factorization

- PARAFAC decomposition

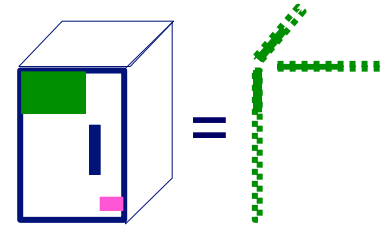


Answer: tensor factorization

- PARAFAC decomposition
- Results for who-calls-whom-when
 - 4M x 15 days

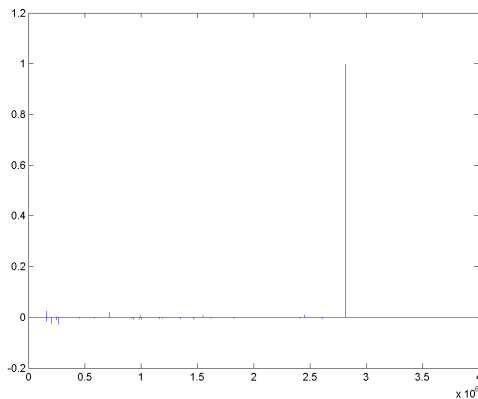


Anomaly detection in time-evolving graphs

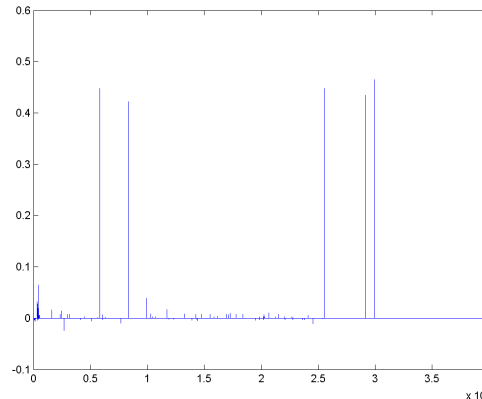


- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks

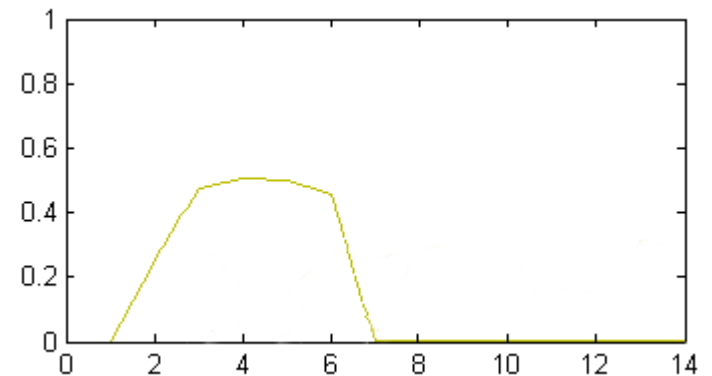
1 caller



5 receivers

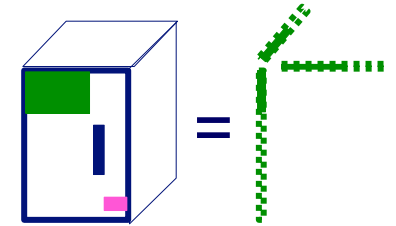


4 days of activity



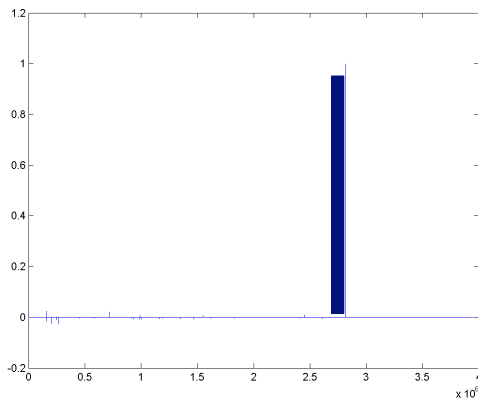
~200 calls to EACH receiver on EACH day!

Anomaly detection in time-evolving graphs

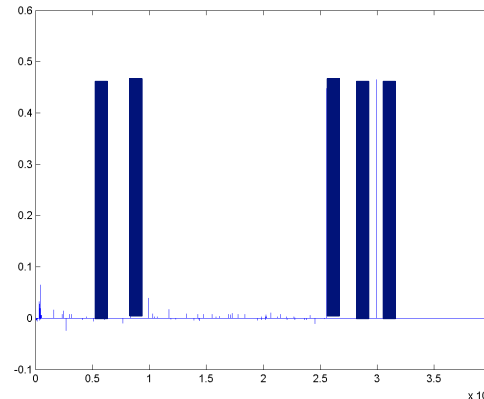


- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks

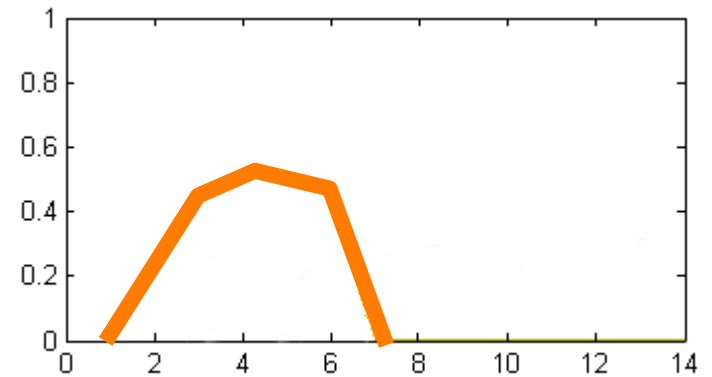
1 caller



5 receivers

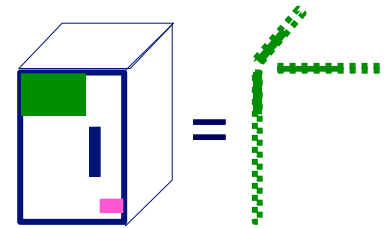


4 days of activity



~200 calls to EACH receiver on EACH day!

Anomaly detection in time-evolving graphs



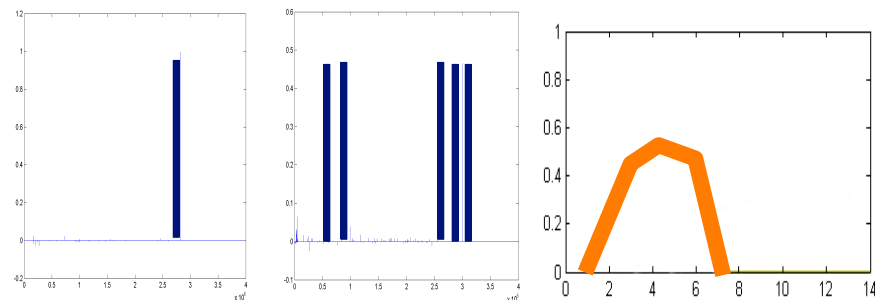
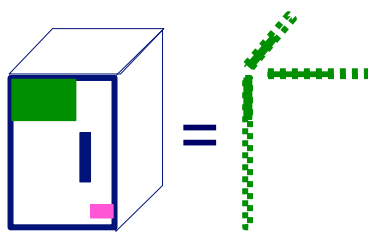
- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks



Miguel Araujo, Spiros Papadimitriou, Stephan Günnemann, Christos Faloutsos, Prithwish Basu, Ananthram Swami, Evangelos Papalexakis, Danai Koutra. *Com2: Fast Automatic Discovery of Temporal (Comet) Communities.* PAKDD 2014, Tainan, Taiwan.

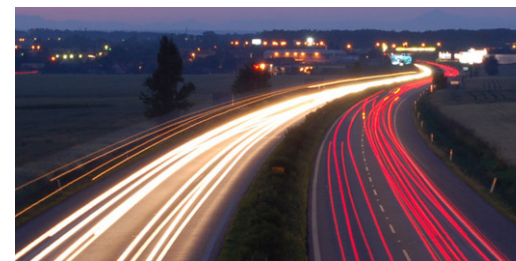
Part 2: Conclusions

- Time-evolving / heterogeneous graphs \rightarrow tensors
- PARAFAC finds patterns
- (GigaTensor/HaTen2 \rightarrow fast & scalable)



Part 3: Time sequences

Roadmap



- Introduction – Motivation
 - Why study (big) graphs?
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
- ➔ • Part#3: time sequences
- Acknowledgements and Conclusions

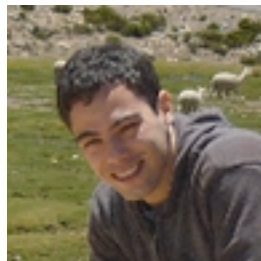
Carnegie Mellon



Carnegie
Mellon
University

KDD 2015 – Sydney,
Australia

RSC: Mining and Modeling Temporal Activity in Social Media



Alceu F. Costa* Yuto Yamaguchi Agma J. M. Traina

Caetano Traina Jr. Christos Faloutsos

*alceufc@icmc.usp.br

Pattern Mining: Datasets

Reddit Dataset

Time-stamp from comments
21,198 users
20 Million time-stamps

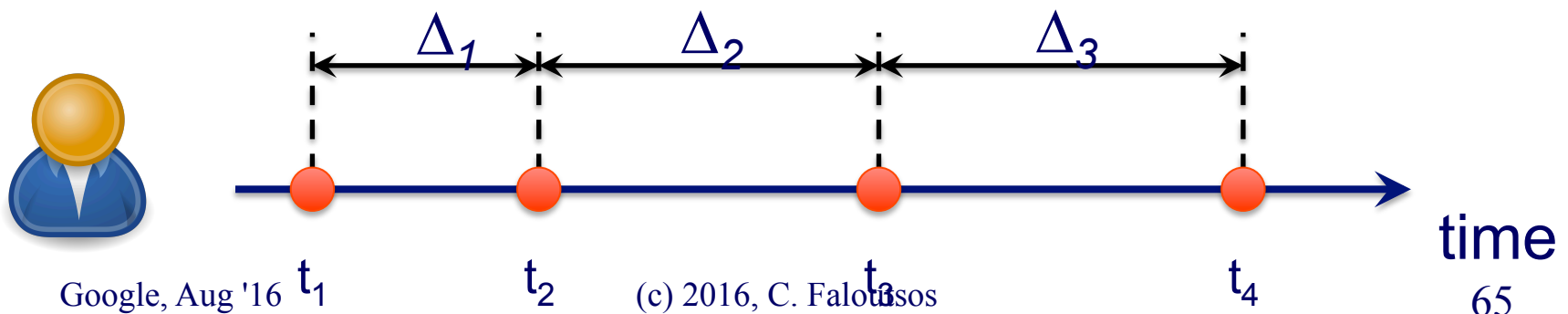
Twitter Dataset

Time-stamp from tweets
6,790 users
16 Million time-stamps

For each user we have:

Sequence of postings time-stamps: $T = (t_1, t_2, t_3, \dots)$

Inter-arrival times (IAT) of postings: $(\Delta_1, \Delta_2, \Delta_3, \dots)$

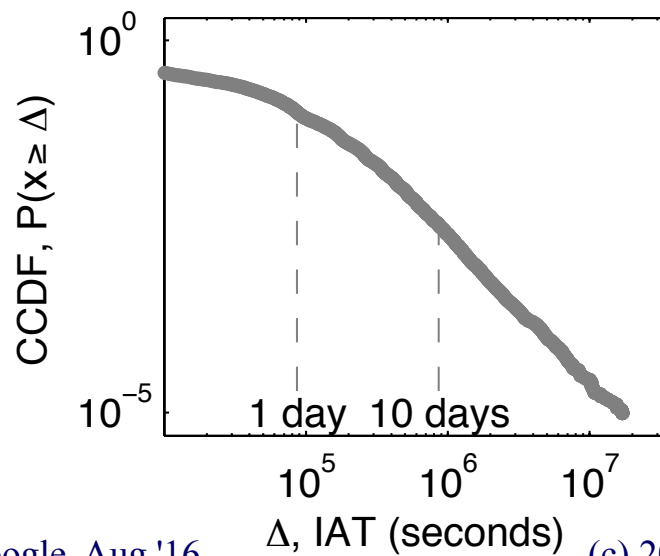


Pattern Mining

Pattern 1: Distribution of IAT is heavy-tailed

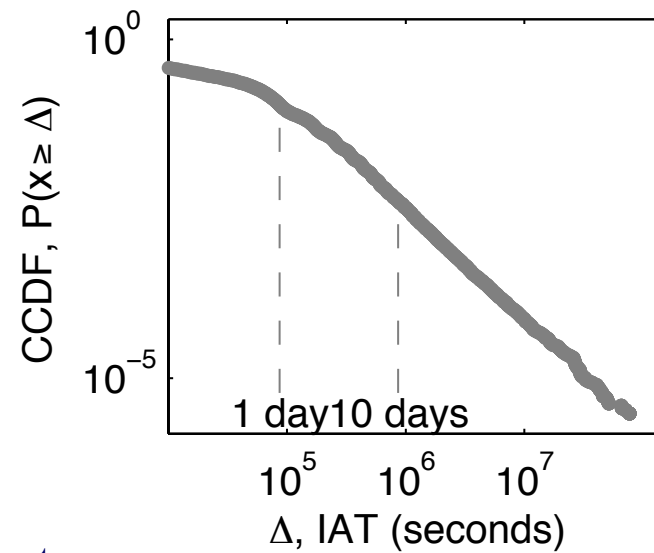
Users can be inactive for long periods of time before making new postings

IAT Complementary Cumulative Distribution Function (CCDF)
(log-log axis)



Google, Aug '16

Reddit Users



Twitter Users

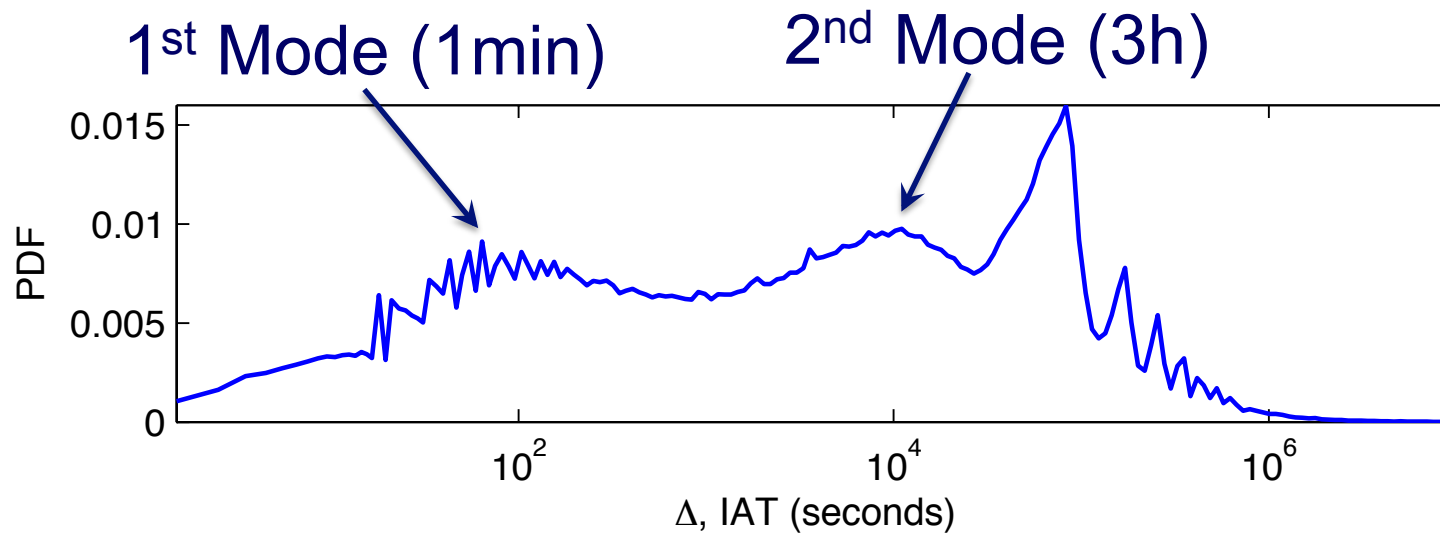
(c) 2016, C. Faloutsos

Pattern Mining

Pattern 2: Bimodal IAT distribution

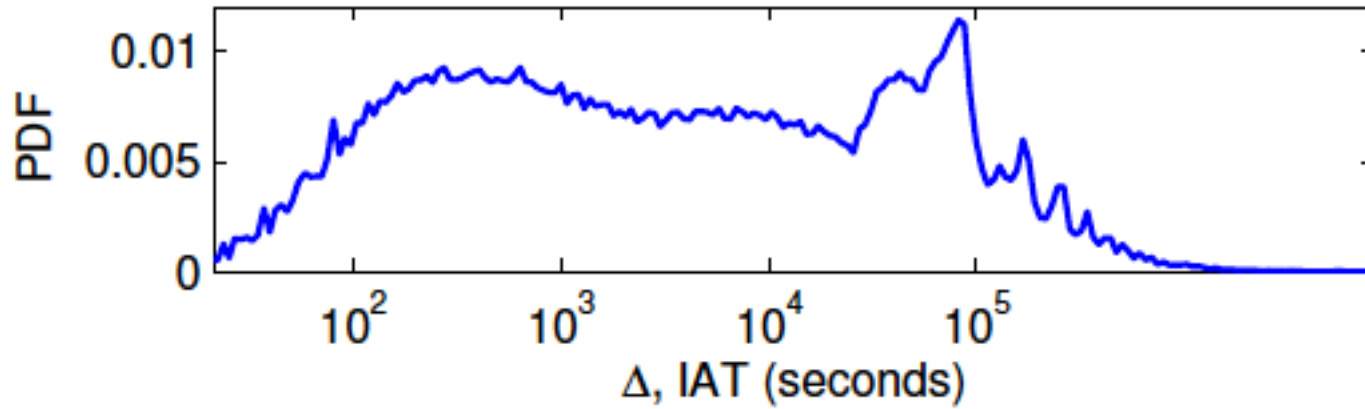
‘Active’/ ‘resting’ periods

Log-binned histogram of postings IAT

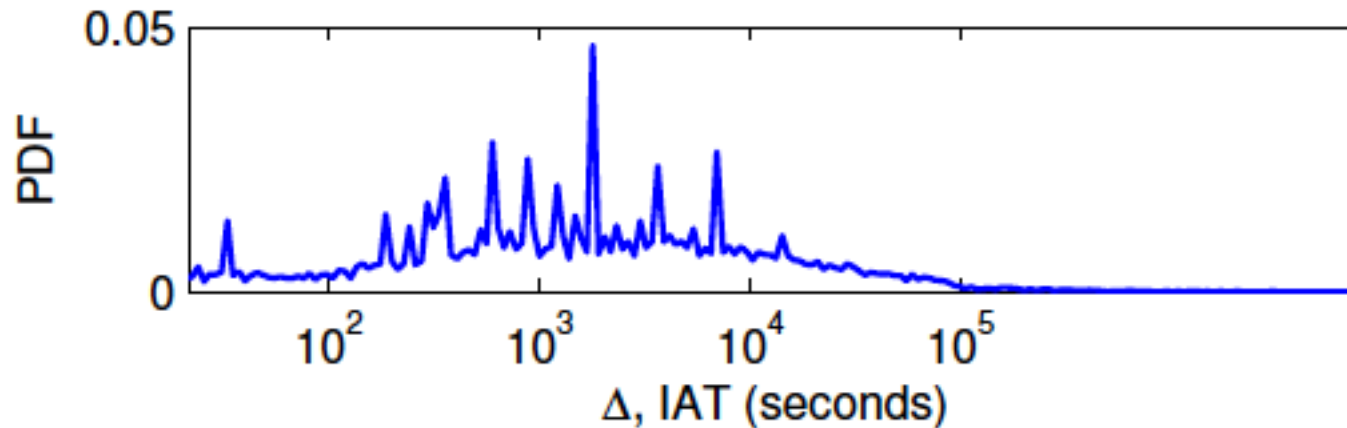


Human? Robots?

linear



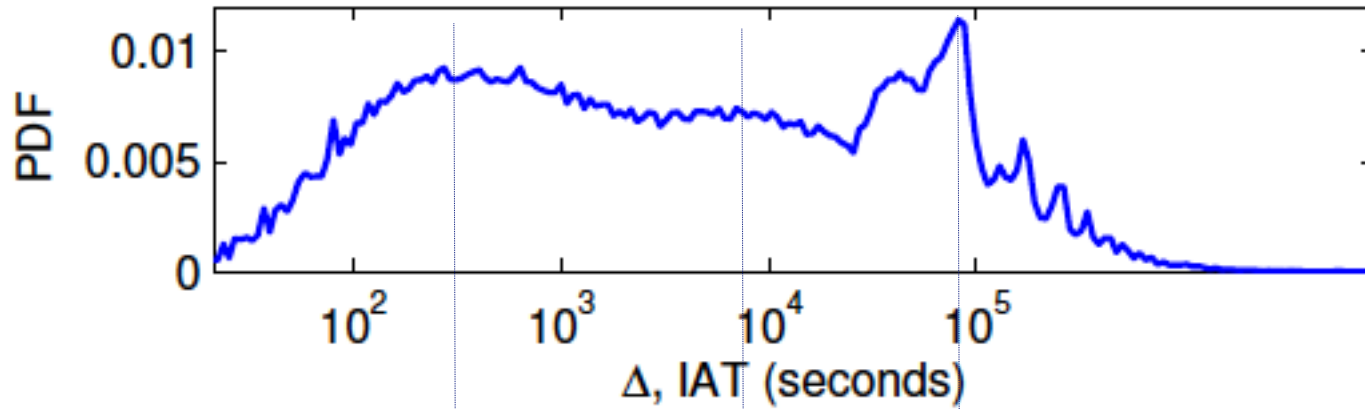
log



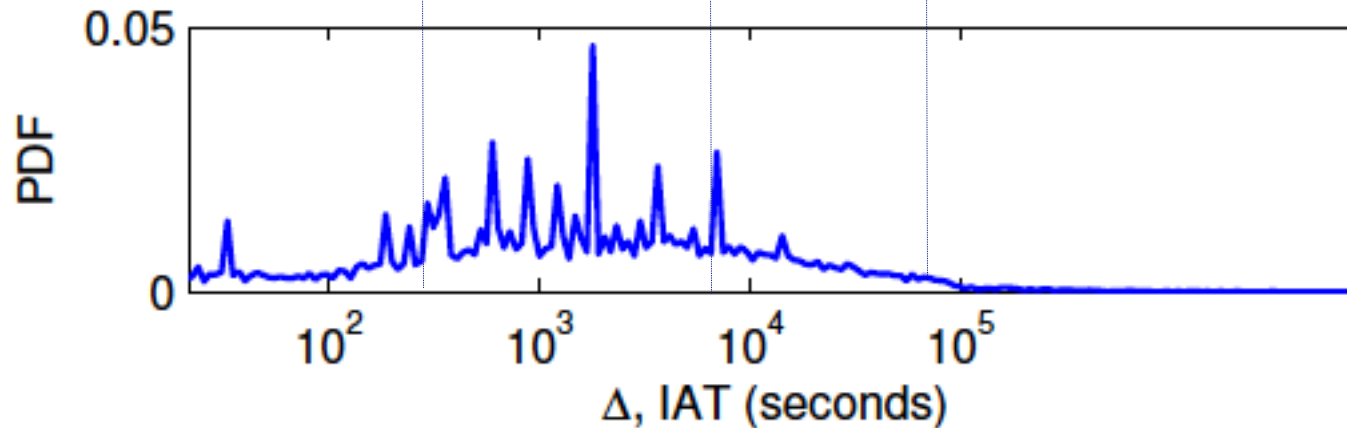
Human? Robots?

2' 3h 1day

linear



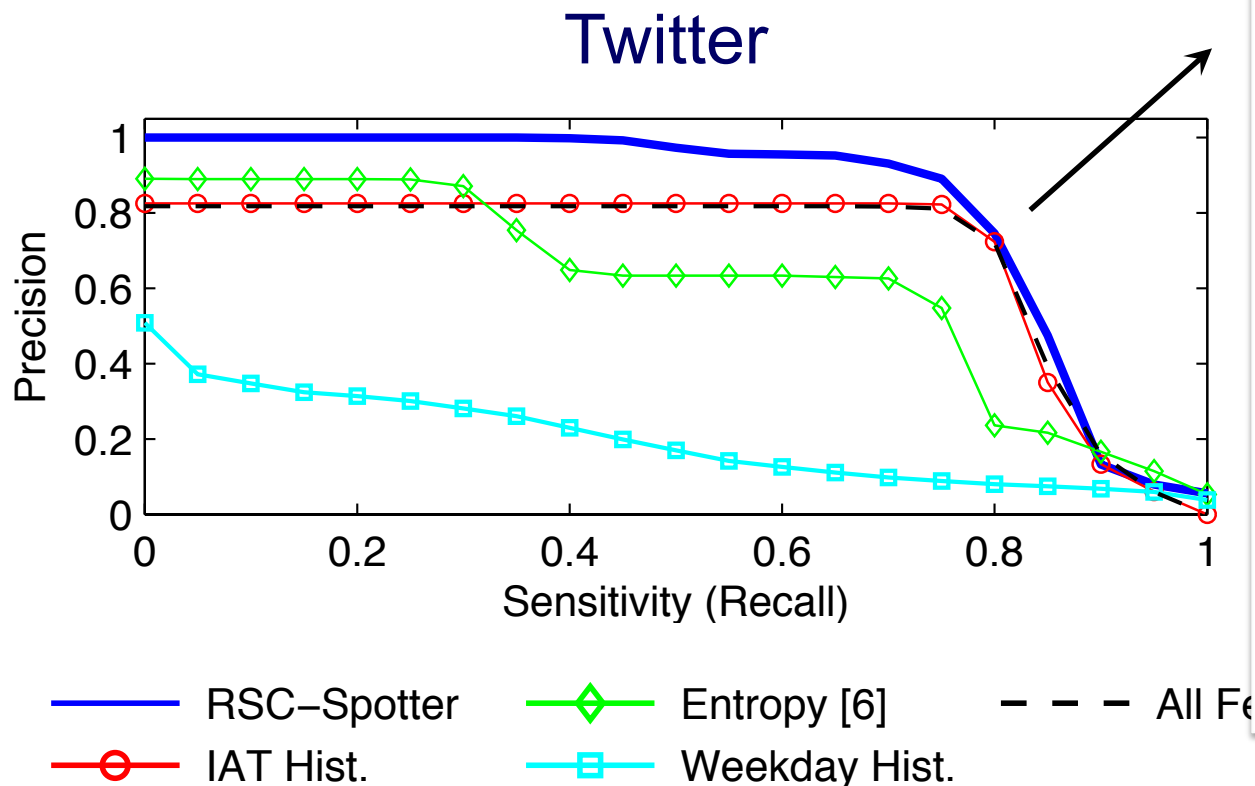
log



Experiments: Can RSC-Spotter Detect Bots?

Precision vs. Sensitivity Curves

Good performance: curve close to the top



Precision > 94%
Sensitivity > 70%

With strongly imbalanced datasets

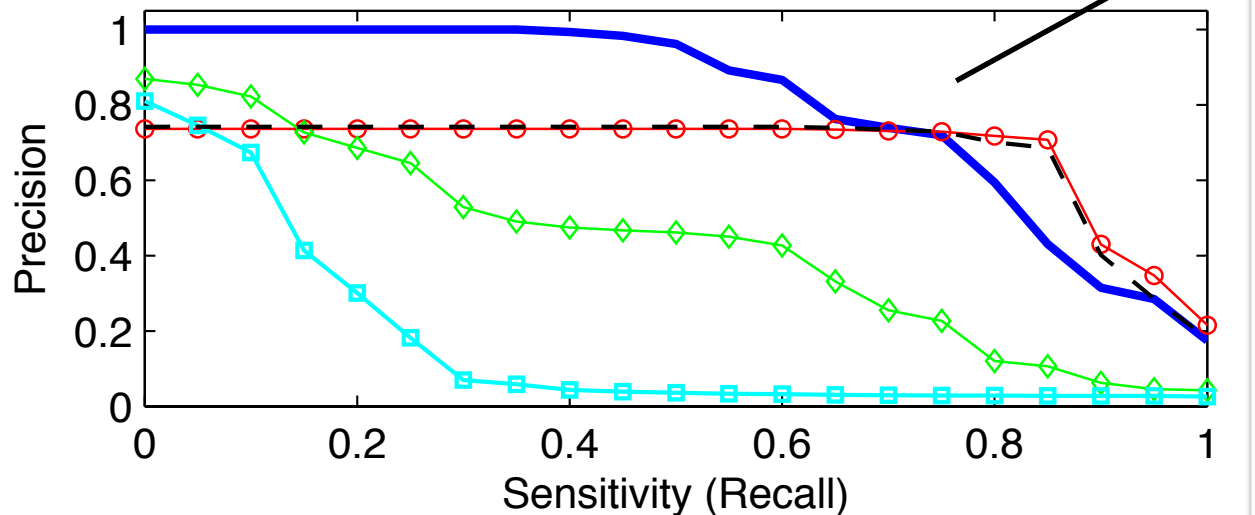
humans \gg # bots

Experiments: Can RSC-Spotter Detect Bots?

Precision vs. Sensitivity Curves

Good performance: curve close to the top

Reddit

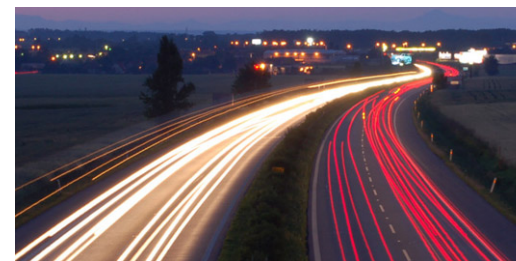


Precision > 96%
Sensitivity > 47%

With strongly imbalanced datasets
humans >> # bots

- RSC-Spotter
 - ◇— Entropy [6]
 - IAT Hist.
 - Weekday Hist.
 - - - All Fe
- Google, Aug '16 (c) 2016, C. Faloutsos

Roadmap



- Introduction – Motivation
 - Why study (big) graphs?
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
- Part#3: time sequences
- ➔ • Acknowledgements and Conclusions

Thanks



Disclaimer: All opinions are mine; not necessarily reflecting the opinions of the funding agencies

Thanks to: NSF IIS-0705359, IIS-0534205, CTA-INARC; Yahoo (M45), LLNL, IBM, SPRINT, Google, INTEL, HP, iLab

Cast



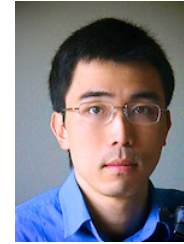
Akoglu,
Leman



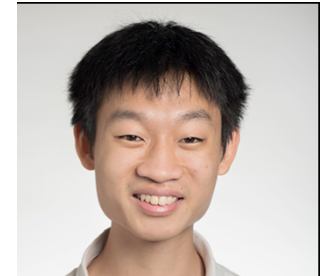
Araujo,
Miguel



Beutel,
Alex



Chau,
Polo



Hooi,
Bryan



Kang, U



Koutra,
Danai



Papalexakis,
Vagelis



Shah,
Neil



Song,
Hyun Ah

Cast



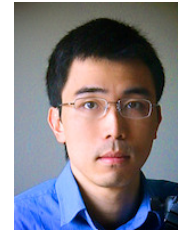
Akoglu,
Leman



Araujo,
Miguel



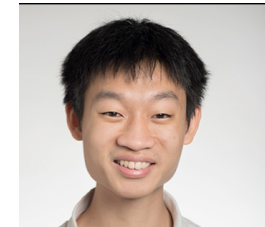
Beutel,
Alex



Chau,
Polo



Eswaran,
Dhivya



Hooi,
Bryan



Kang, U



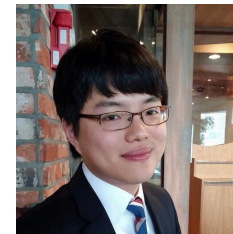
Koutra,
Danai



Papalexakis,
Vagelis



Shah,
Neil




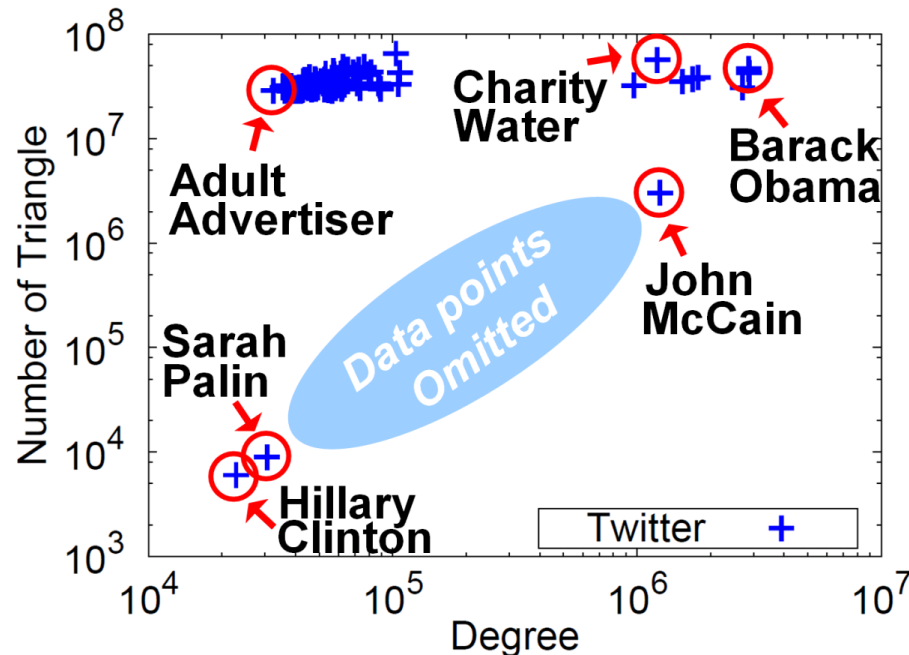
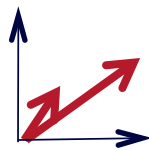
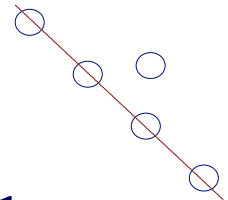
Shin,
Kijung



Song,
Hyun Ah

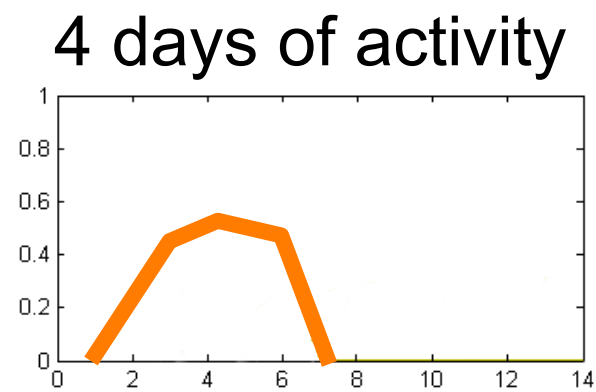
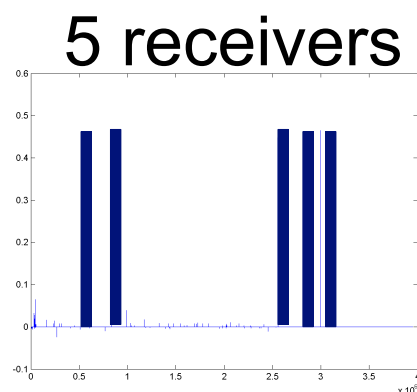
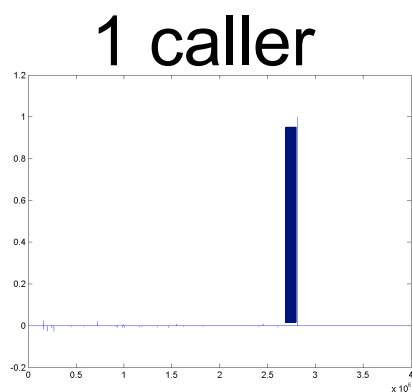
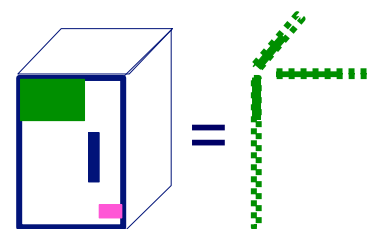
CONCLUSION#1 – Big data

- **Patterns**  **Anomalies**
- **Large datasets reveal patterns/outliers that are invisible otherwise**



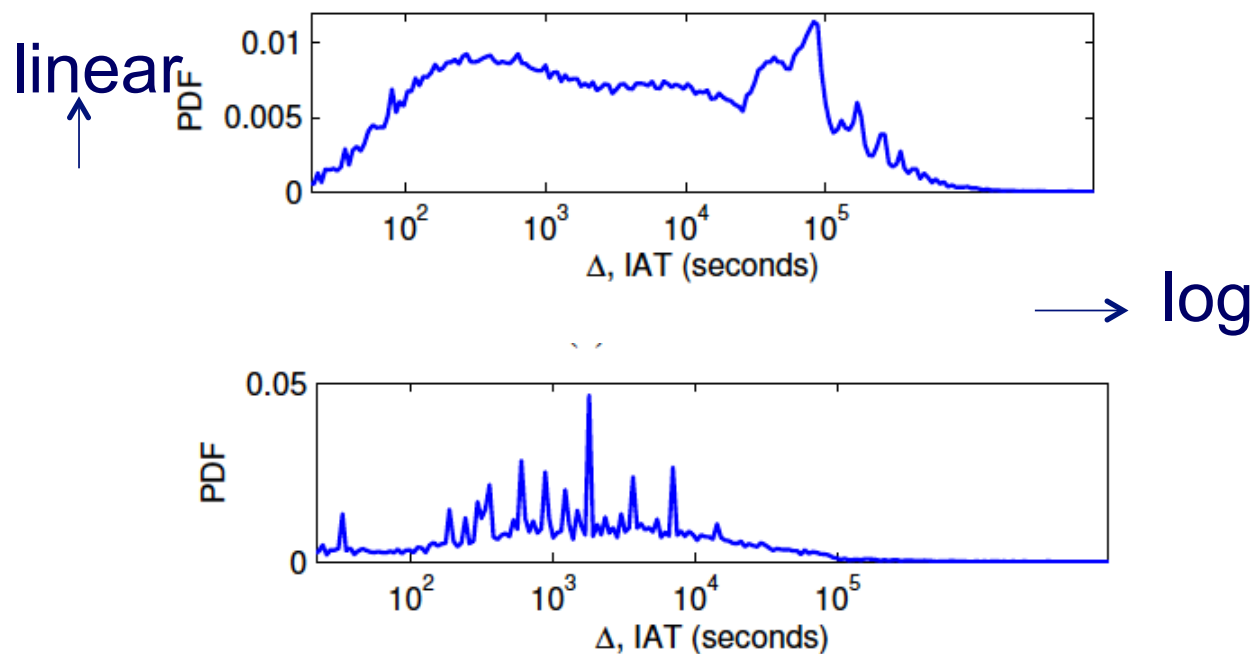
CONCLUSION#2 – tensors

- powerful tool



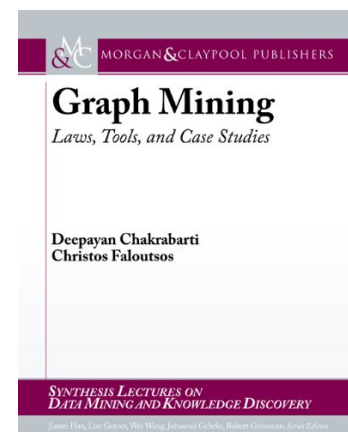
Conclusion#3

- Different footprints of real vs 'robot' users

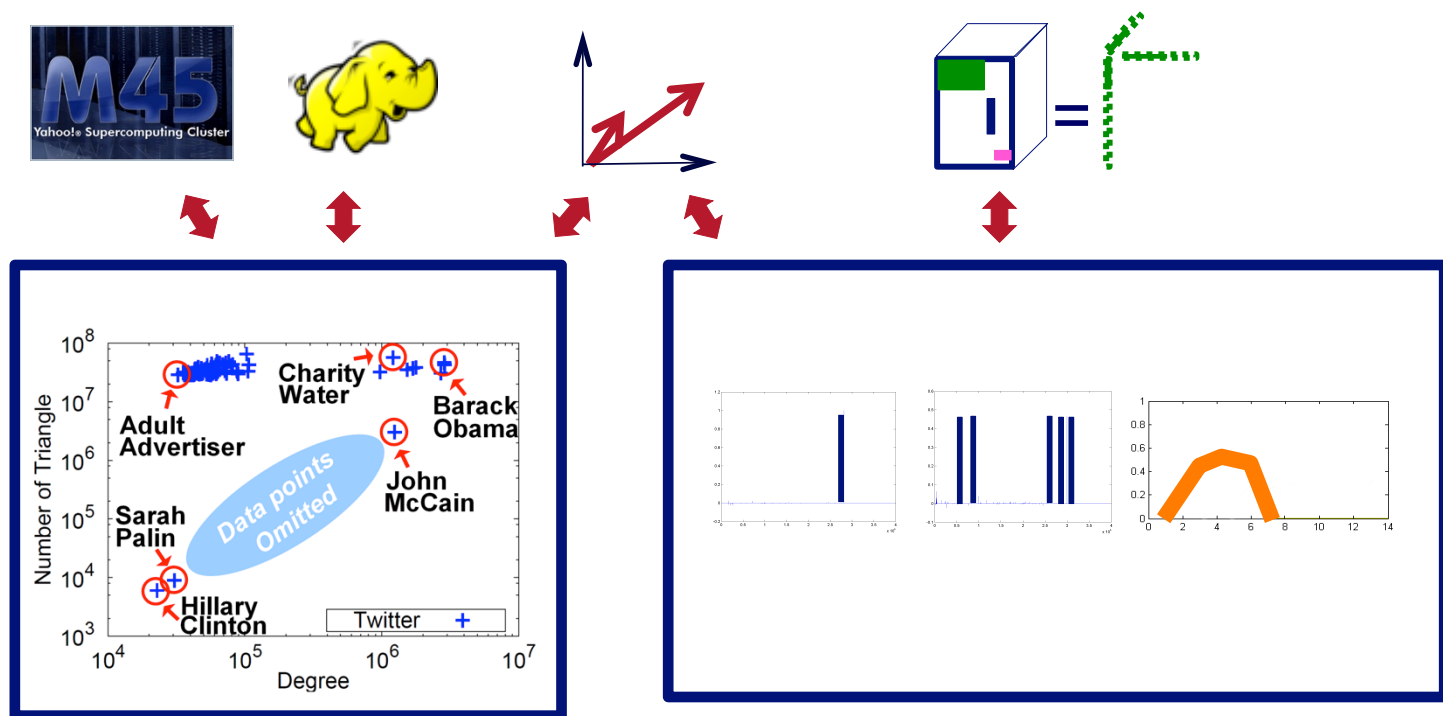


References

- D. Chakrabarti, C. Faloutsos: *Graph Mining – Laws, Tools and Case Studies*, Morgan Claypool 2012
- <http://www.morganclaypool.com/doi/abs/10.2200/S00449ED1V01Y201209DMK006>
- *Graph-based Anomaly Detection and Description: A Survey*, [Leman Akoglu](#), [Hanghang Tong](#), [Danai Koutra](#)
- <http://arxiv.org/abs/1404.4679>

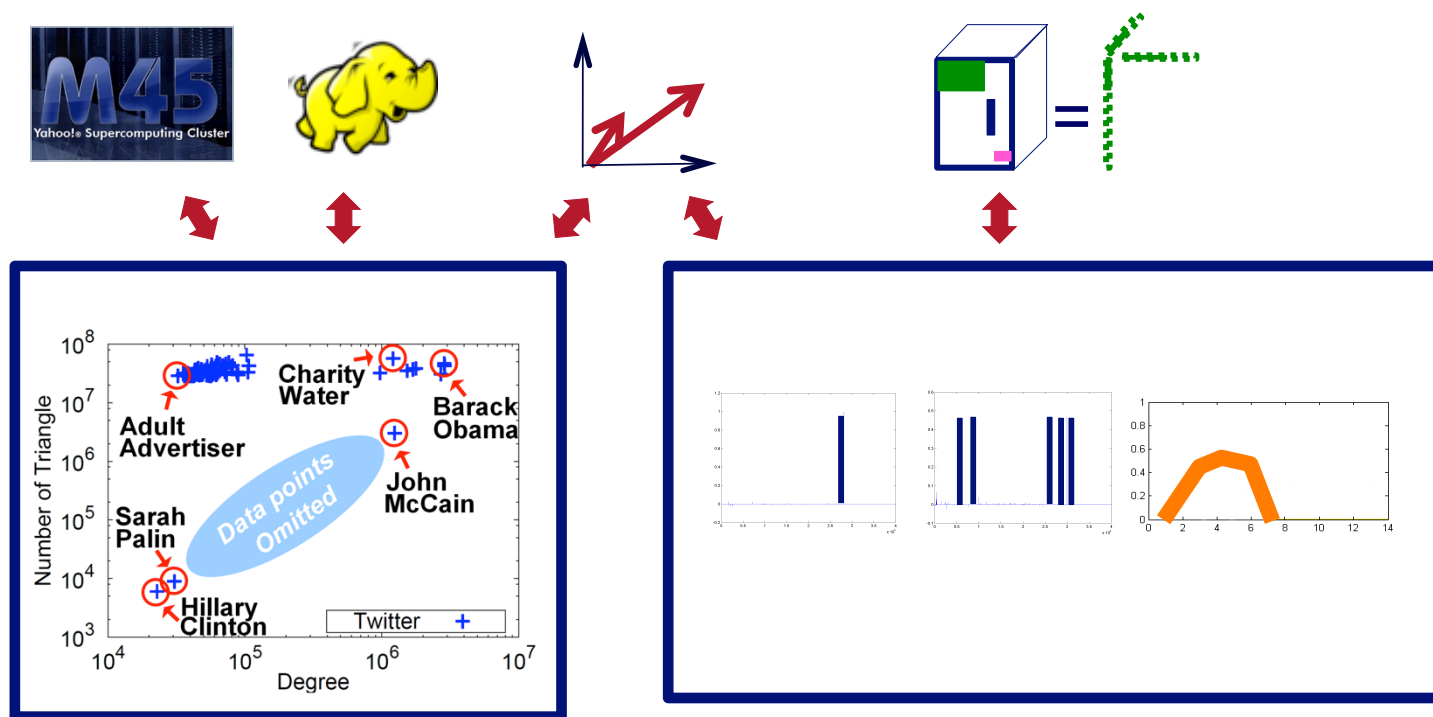


TAKE HOME MESSAGE: Cross-disciplinarity



Thank you!

Cross-disciplinarity





Catchsync: catch synchronized behavior in large directed graphs

Meng Jiang

Joint work with Peng Cui, Alex Beutel,
Christos Faloutsos and Shiqiang Yang

August 26, 2014 – NYC, USA



Fraud Detection: Graph Analysis Problem



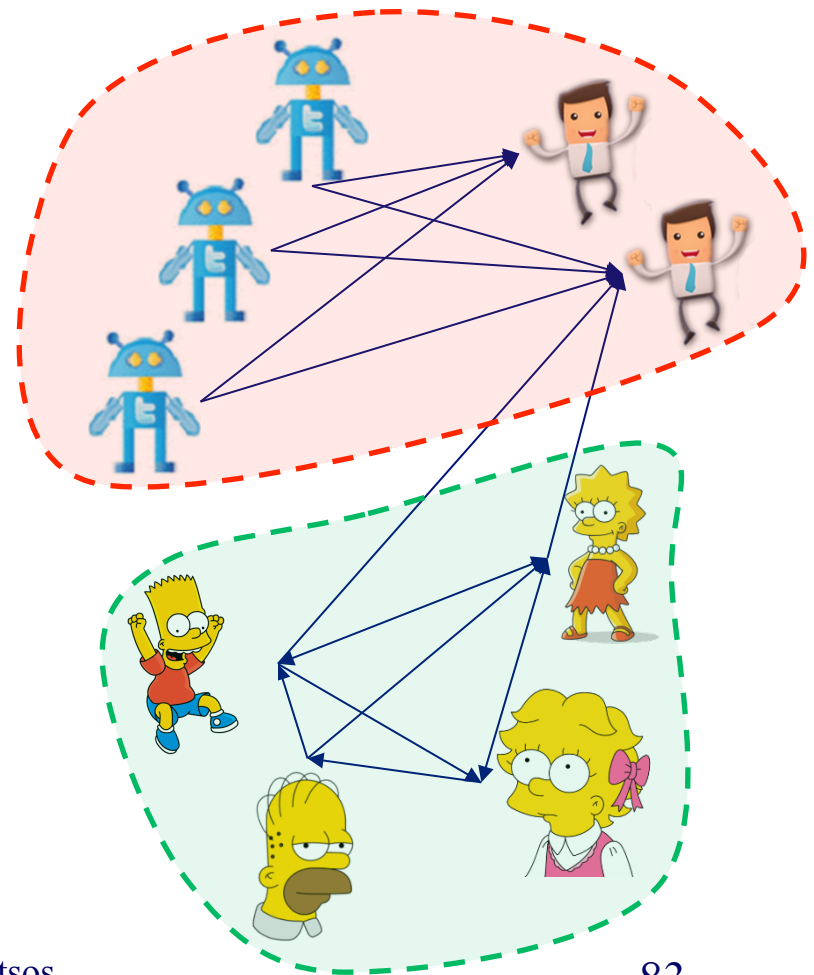
[www.buyfollowz.org]

<p>5,000 FOLLOWERS 400 FREE</p> <p>\$69.99</p> <p>Delivery within 3-4 days</p> <p>Buy Now</p> <p>Save + 3%</p>	<p>2,000 FOLLOWERS 300 FREE</p> <p>\$29.99</p> <p>Delivery within 2-3 days</p> <p>Buy Now</p> <p>Save + 2%</p>	<p>1,000 FOLLOWERS 200 FREE</p> <p>\$15.99</p> <p>Delivery within 1-2 days</p> <p>Buy Now</p>	<p>10,000 FOLLOWERS 500 FREE</p> <p>\$119.99</p> <p>Delivery within 4-5 days</p> <p>Buy Now</p> <p>Save + 14%</p>	<p>20,000 FOLLOWERS 1000 FREE</p> <p>\$229.99</p> <p>Delivery within 5-8 days</p> <p>Buy Now</p> <p>Save + 34%</p>
--	--	---	---	--



[buymorelikes.com]

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

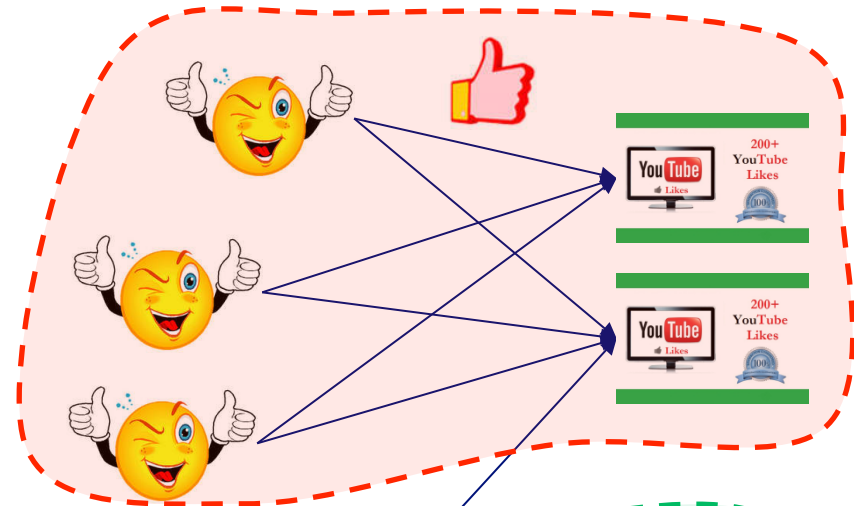


Fraud Detection: Graph Analysis Problem



[buycheaplikes.com]

PACK-1	PACK-2	PACK-3	PACK-4
YOUTUBE \$5	YOUTUBE \$9	YOUTUBE \$13	YOUTUBE \$25
150 Real YouTube Likes	300 Real YouTube Likes	500 Real YouTube Likes	1,000 Real YouTube Likes
\$ 5.00 (USD)	\$ 9.00 (USD)	\$ 13.00 (USD)	\$ 25.00 (USD)
Delivery within 24 hours Enter Your Video URL:	Deliver within 24-48 hours Please Enter Your Video URL:	Delivery within 24-48 hours Enter Your Video URL:	Delivery within 2-3 days Enter Your Video URL:
Current number of likes: <input type="text"/>	Current number of likes: <input type="text"/>	Current number of likes: <input type="text"/>	Current number of likes: <input type="text"/>
Add to Cart	Add to Cart	Add to Cart	Add to Cart



[reviewsteria.com]

It's easy to buy Amazon reviews. Just choose the number of reviews you would like to receive.

High quality reviews that customers love. 100% unique content by native speaking professional writers.

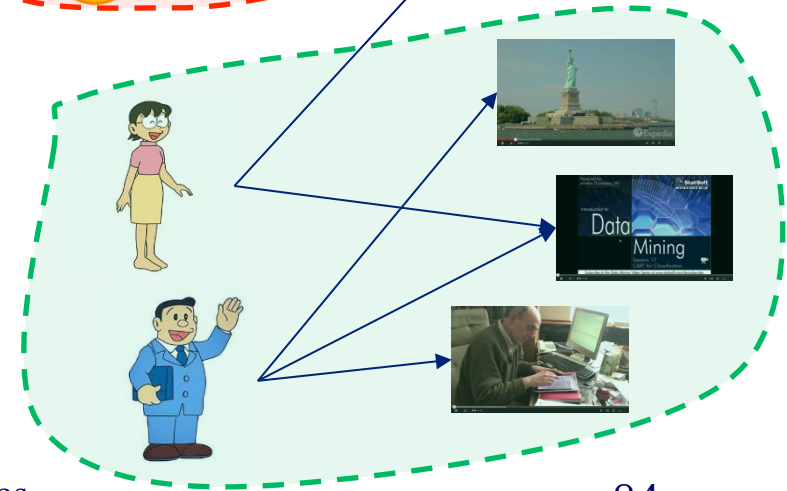
Choose the number of reviews and click Buy Now button to ramp up your Amazon business NOW.

Choose the number of reviews:

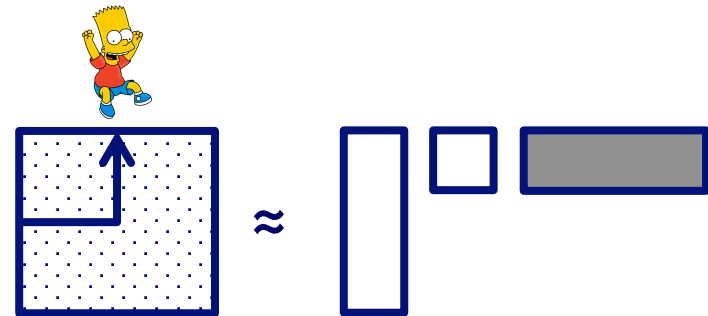
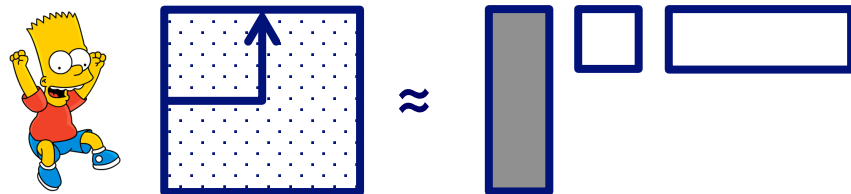
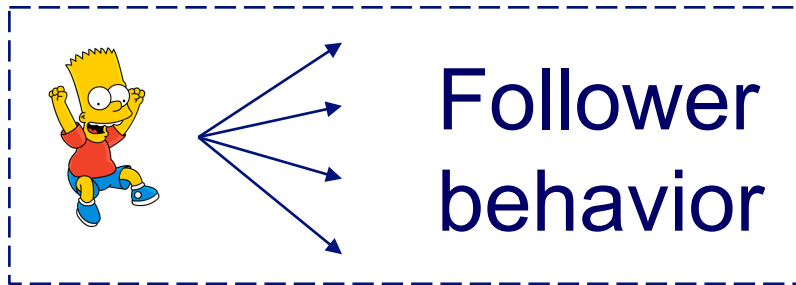
20

Buy Now

MasterCard VISA American Express Discover



Behavior-based Features



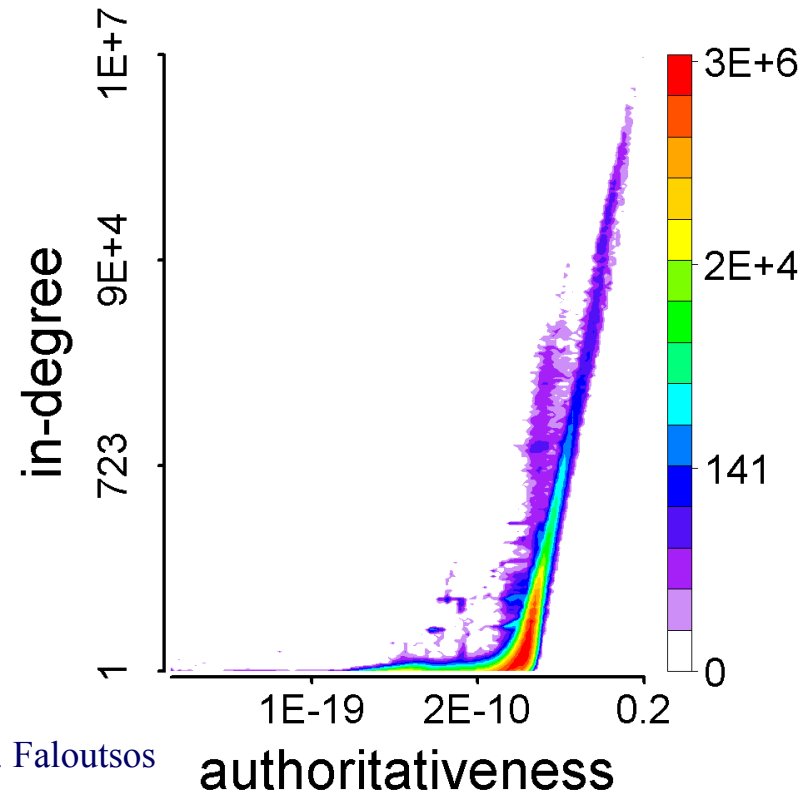
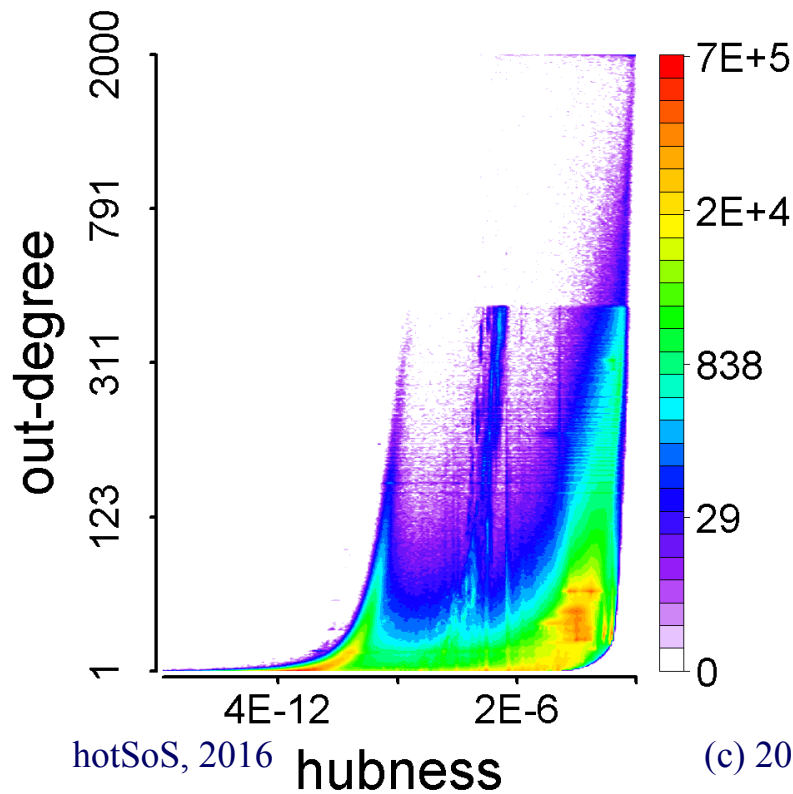
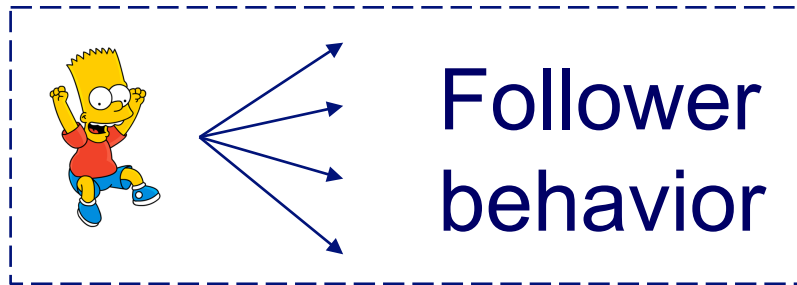
Out-degree
 1st left singular vector
 (**Hubness**)
 2nd left singular vector

.. hotSoS, 2016

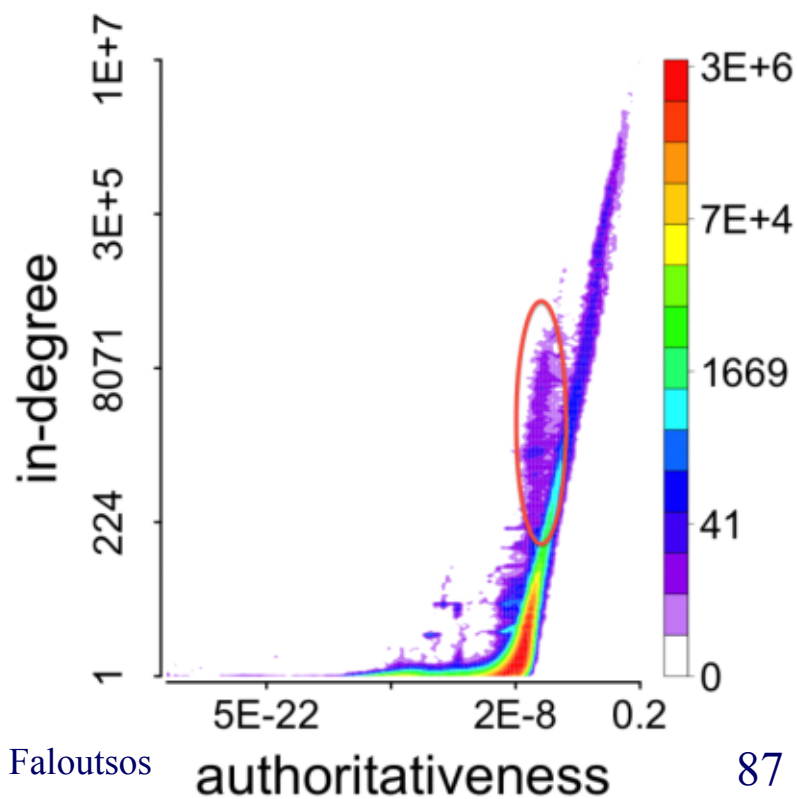
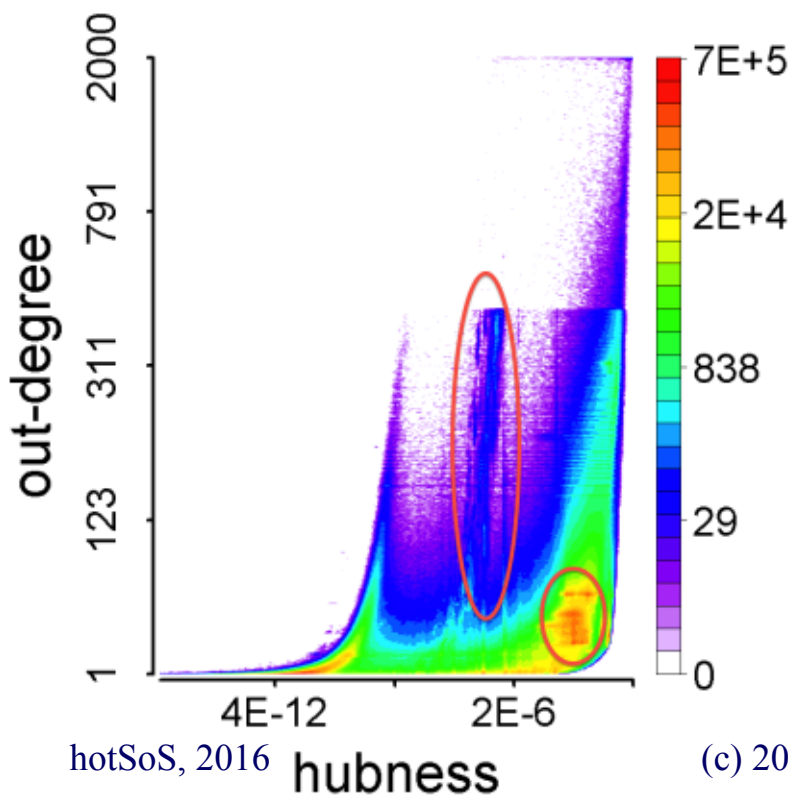
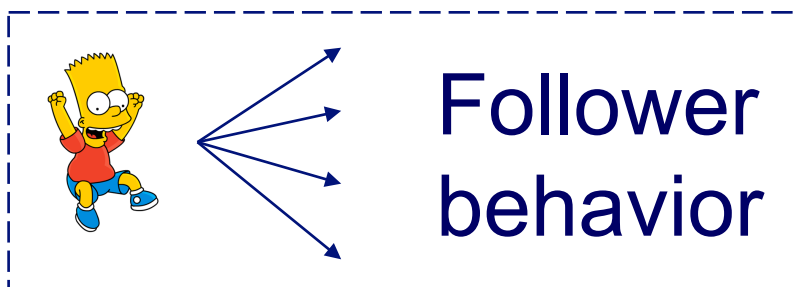
In-degree
 1st right singular vector
 (**Authoritativeness**)
 2nd right singular vector

(c) 2016, C. Faloutsos

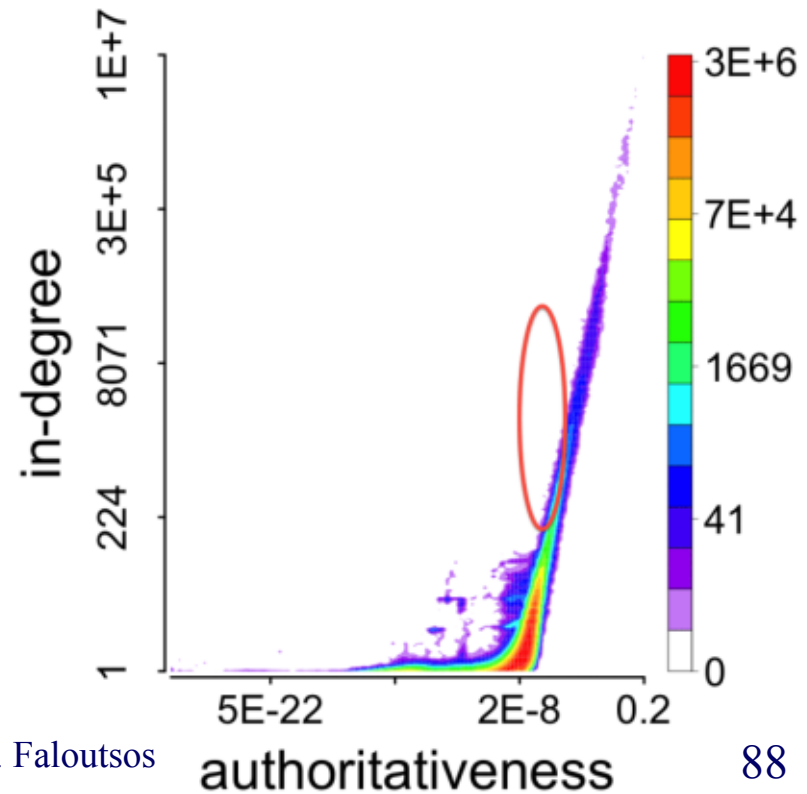
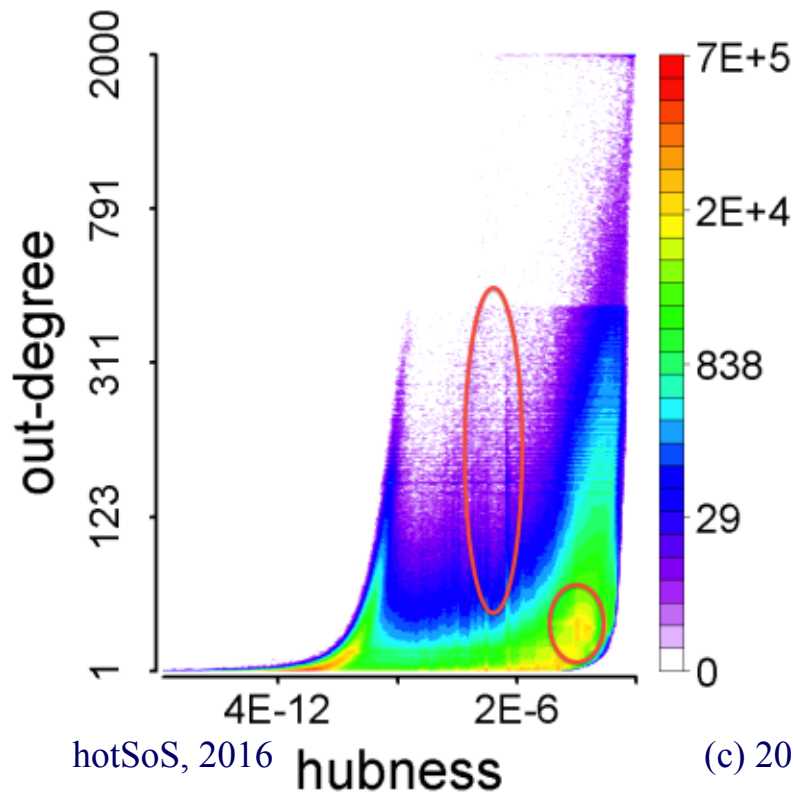
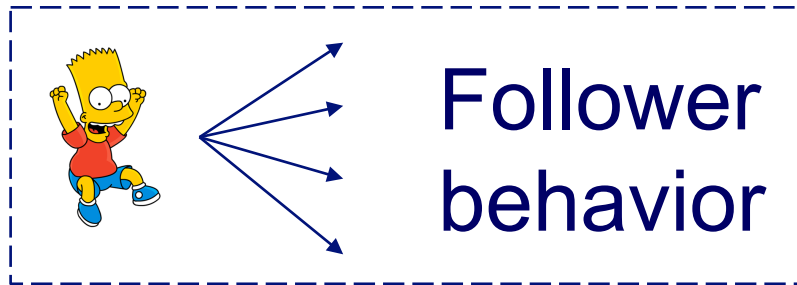
Behavior-based Feature Space



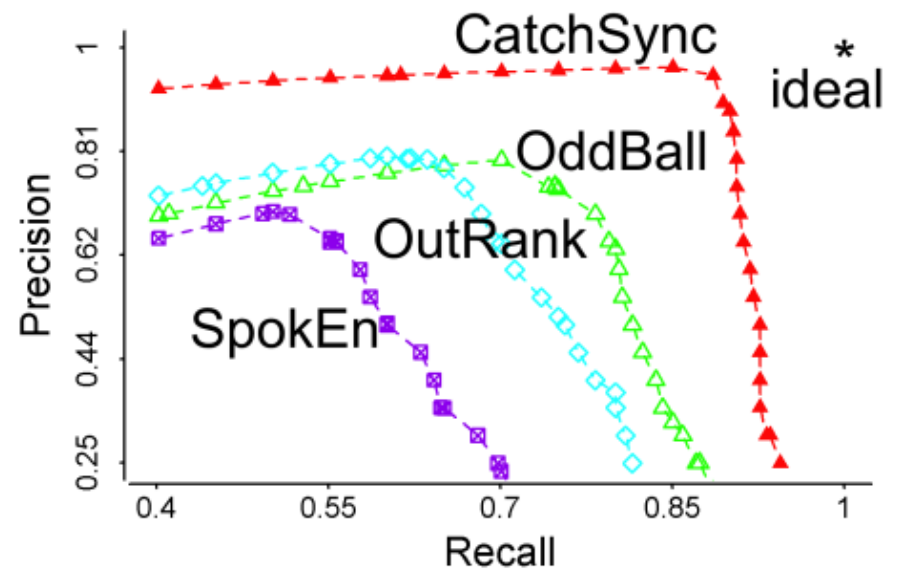
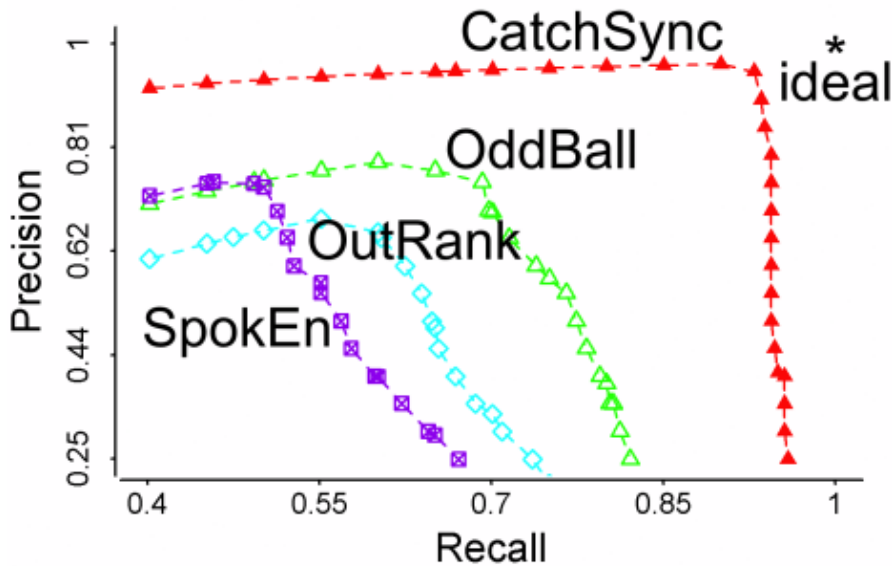
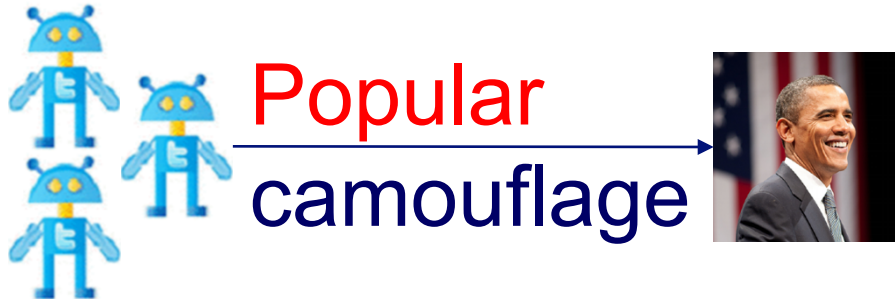
Before CatchSync



After CatchSync

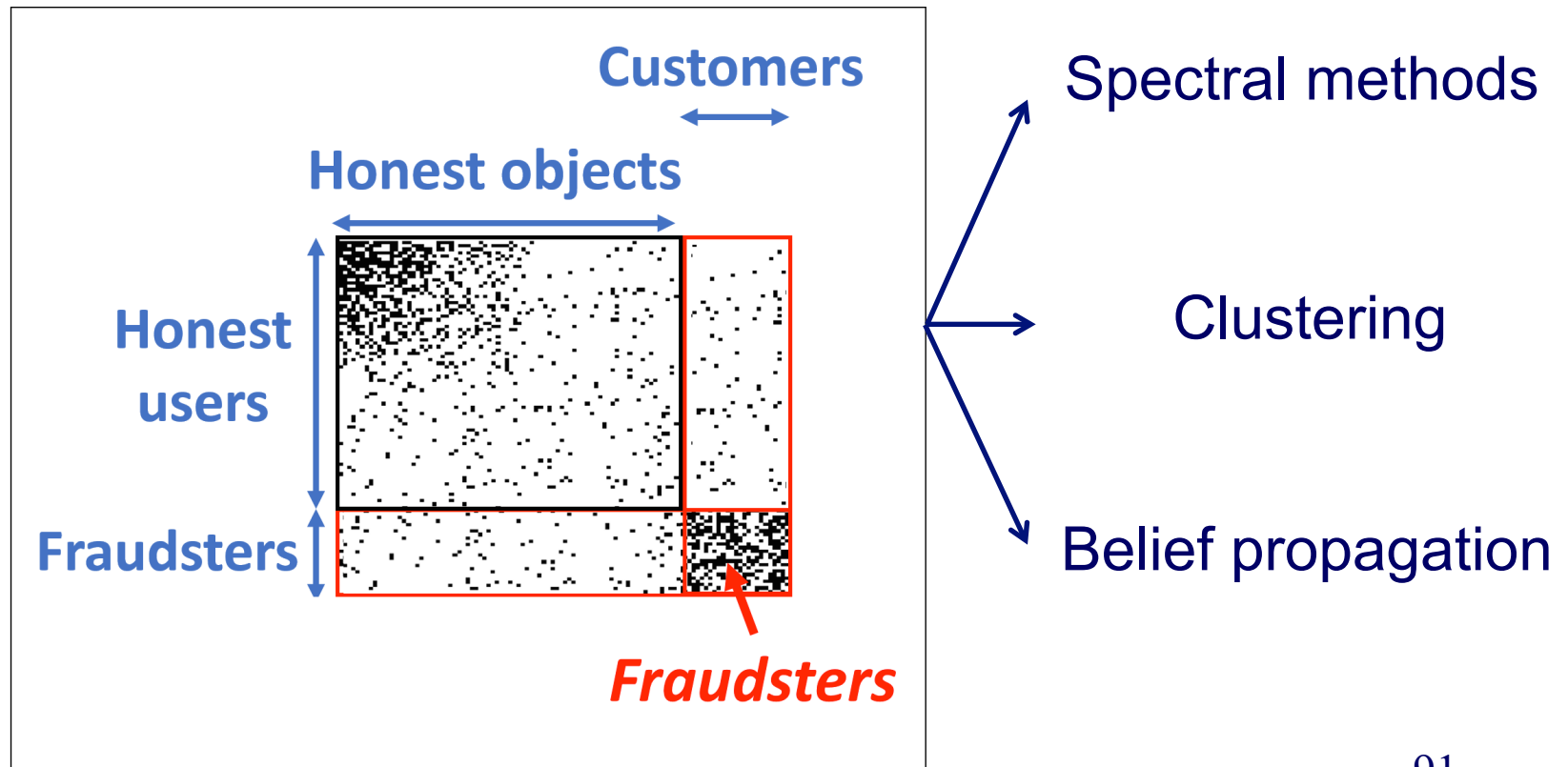


Q3: Is CatchSync Robust?



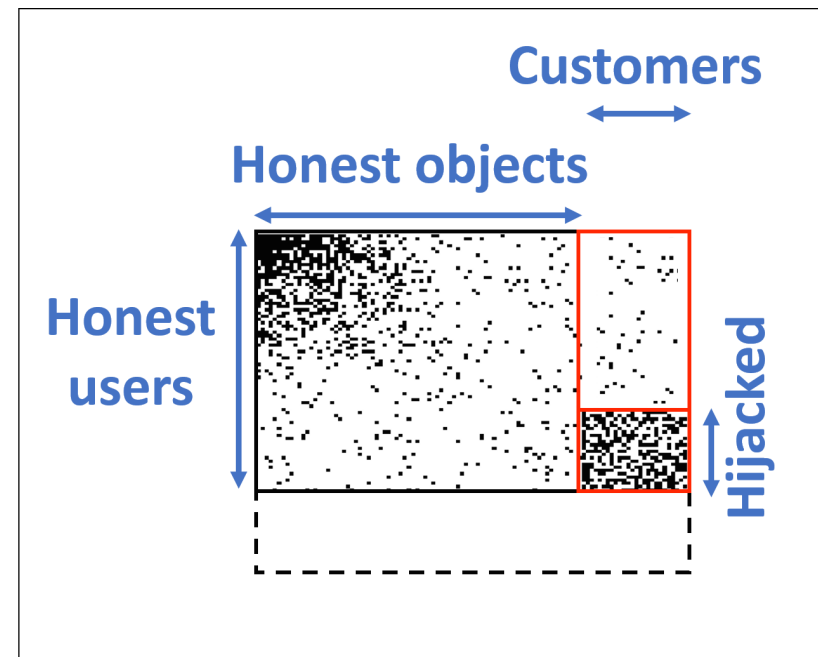
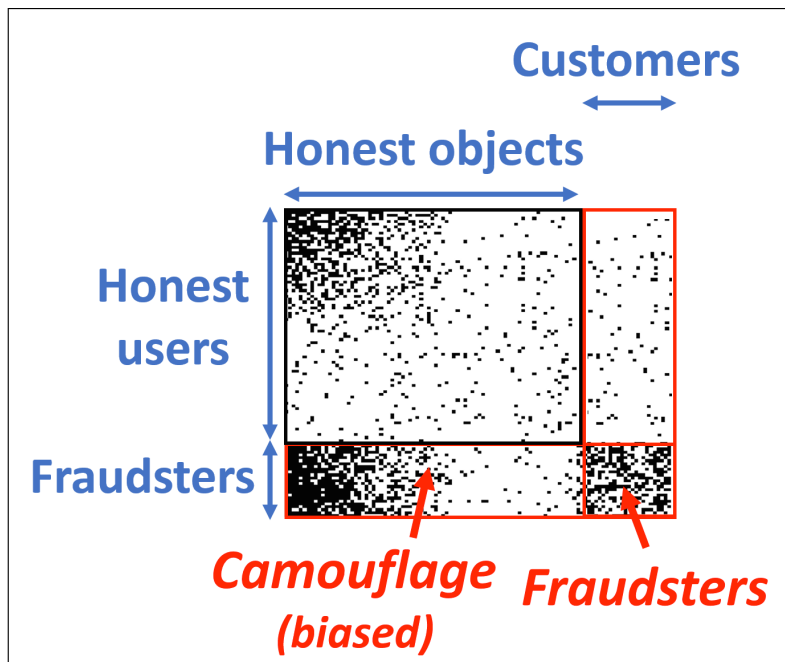
Detecting Review Spam

Many existing methods detect fraudsters using dense subgraph detection.



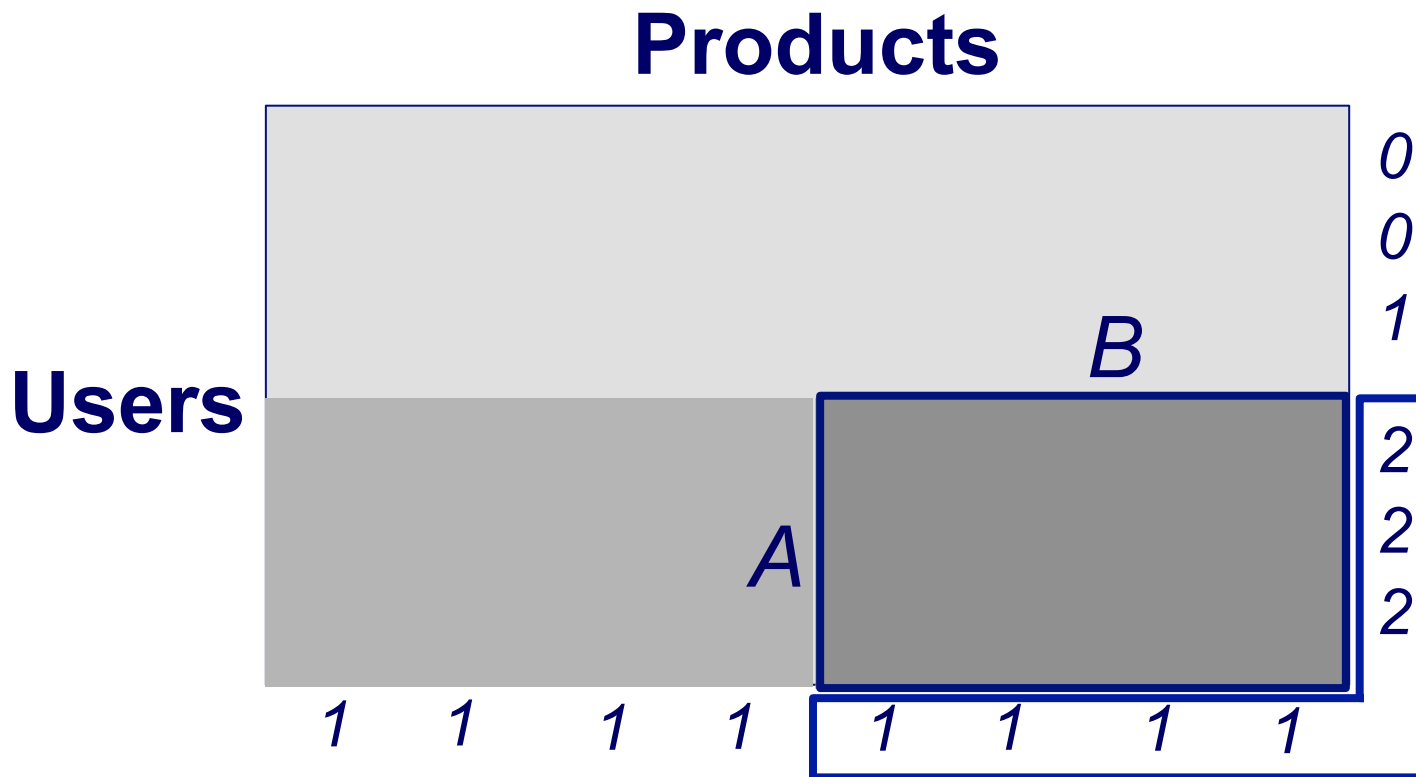
Evading Detection

Attackers can evade detection using *camouflage*.





Node suspiciousness



Node suspiciousness of (A,B)
= 10



Edge suspiciousness

Products

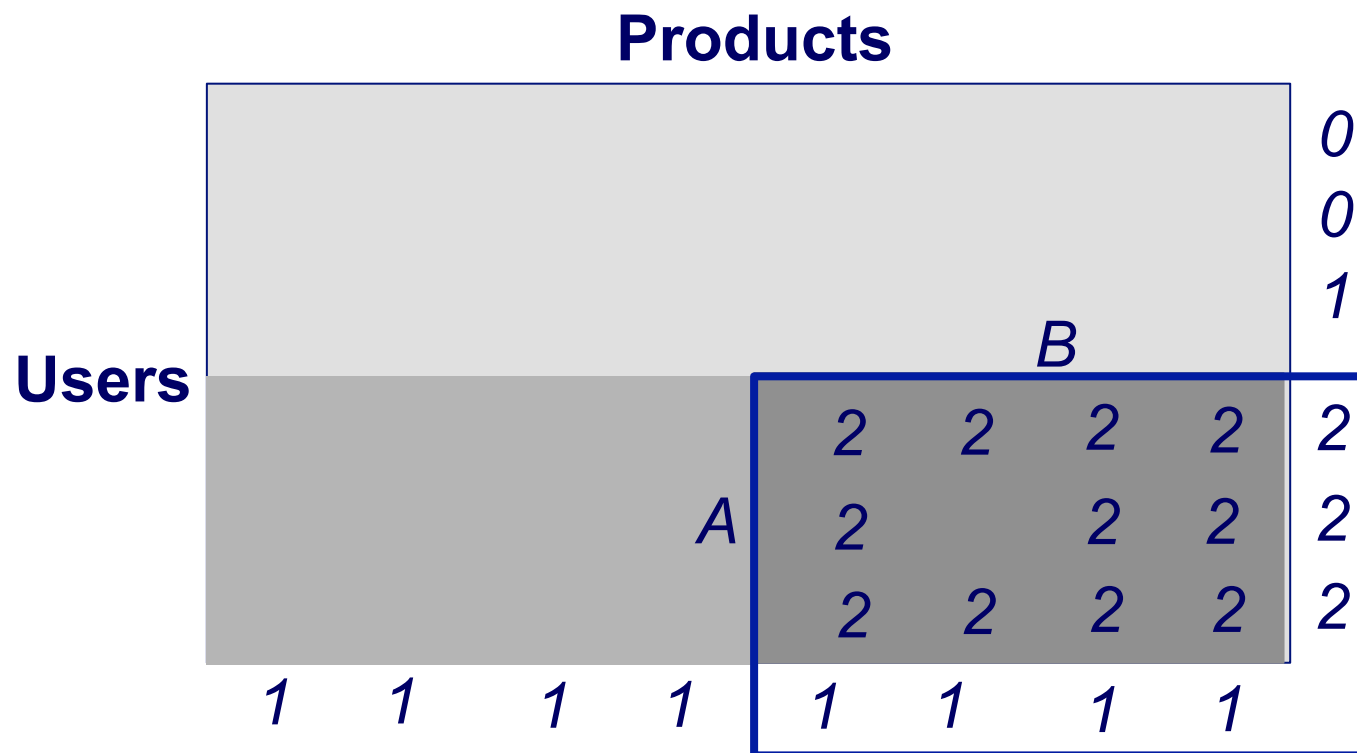
Users

		<i>B</i>			
		2	2	2	2
	<i>A</i>	2		2	2
		2	2	2	2

Edge suspiciousness of (A,B)
= 22



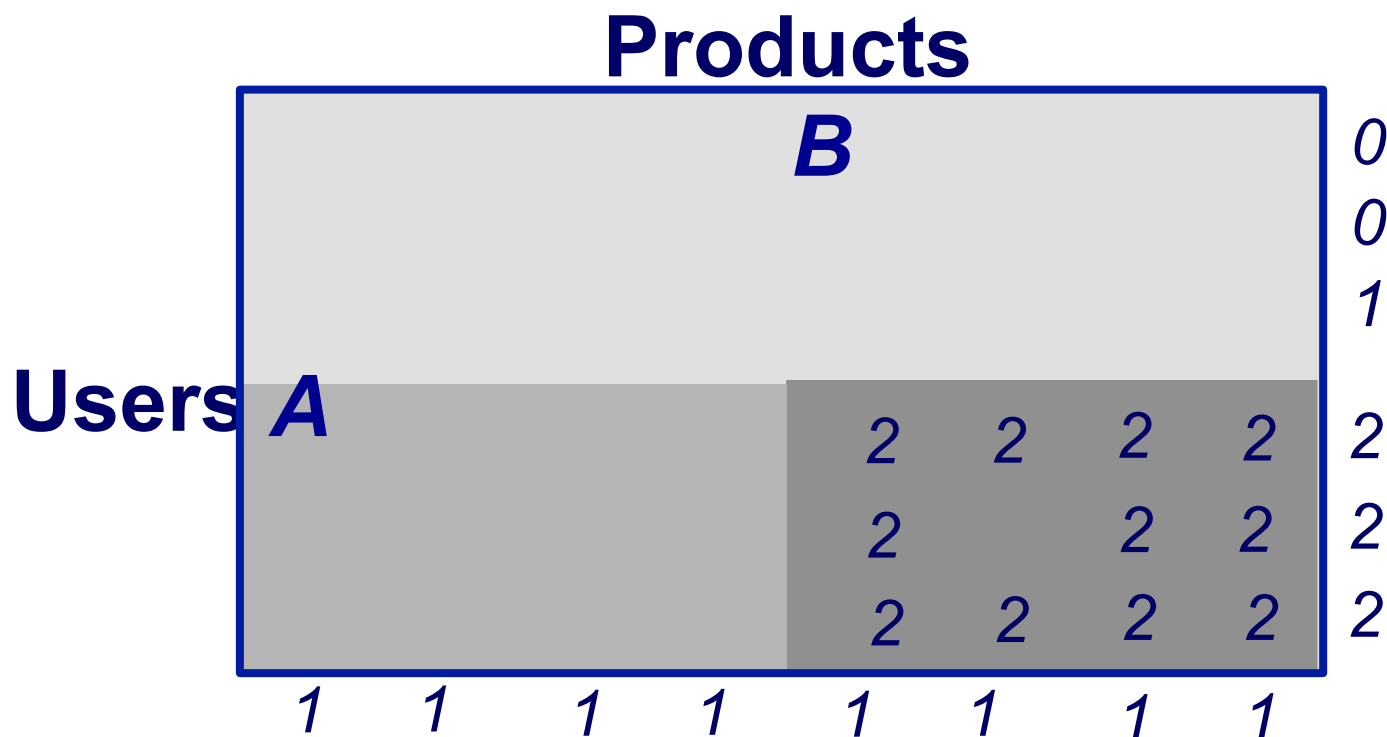
Total suspiciousness



$$f(A,B) = 32$$

$$f(A,B) = (\text{edge susp.}) + (\text{node susp.})$$

Greedy Algorithm



Start with A, B as all users / products



Greedy Algorithm

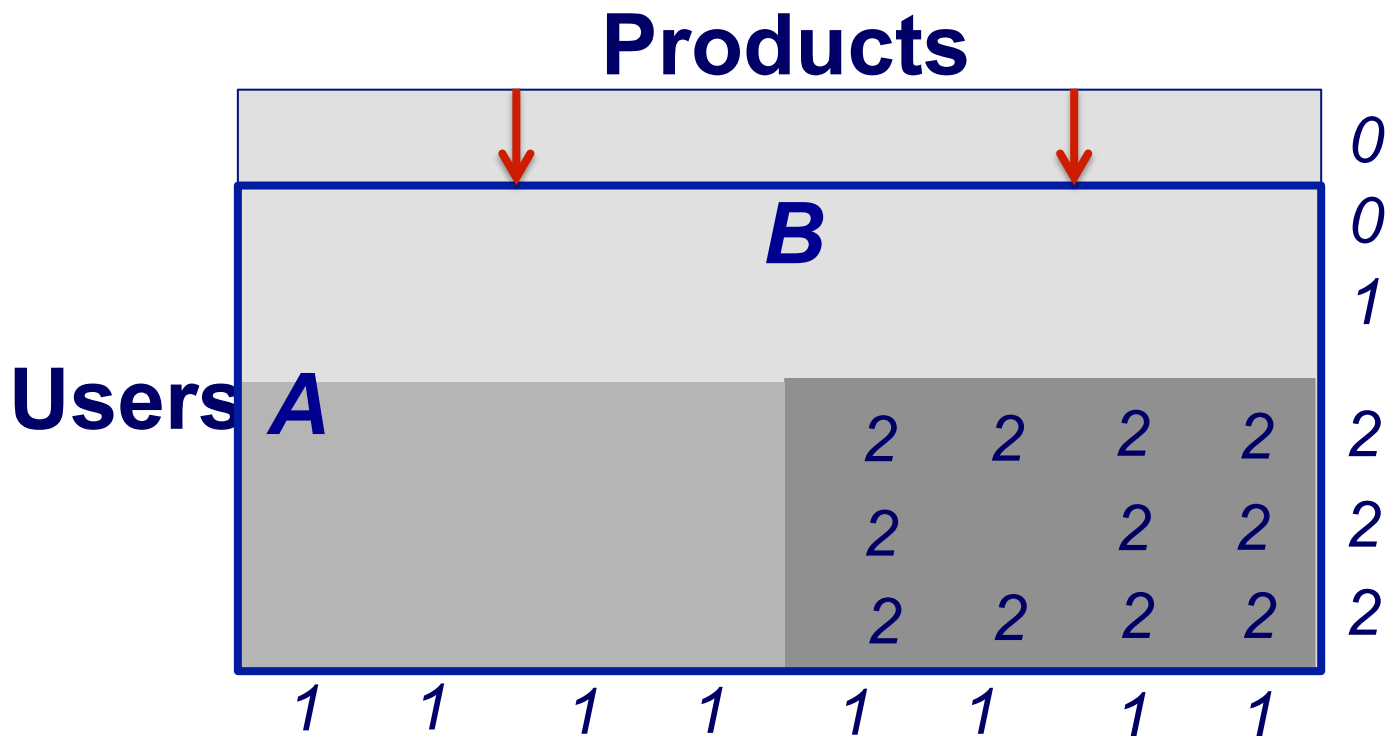
B
Products

								0	
								0	
								1	
Users	A				2	2	2	2	2
					2		2	2	2
					2	2	2	2	2
		1	1	1	1	1	1	1	

Delete rows / columns greedily to maximize g

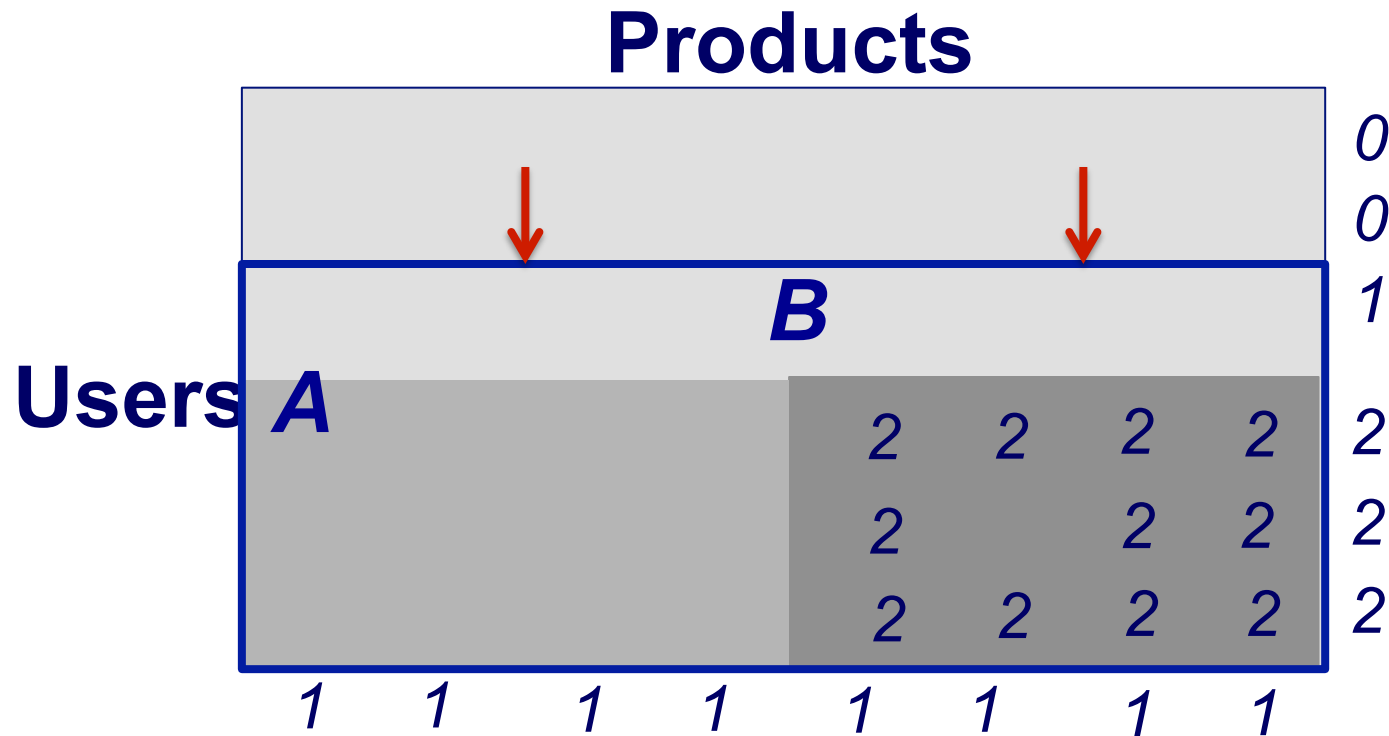


Greedy Algorithm



Delete rows / columns greedily to maximize g

Greedy Algorithm



Delete rows / columns greedily to maximize g



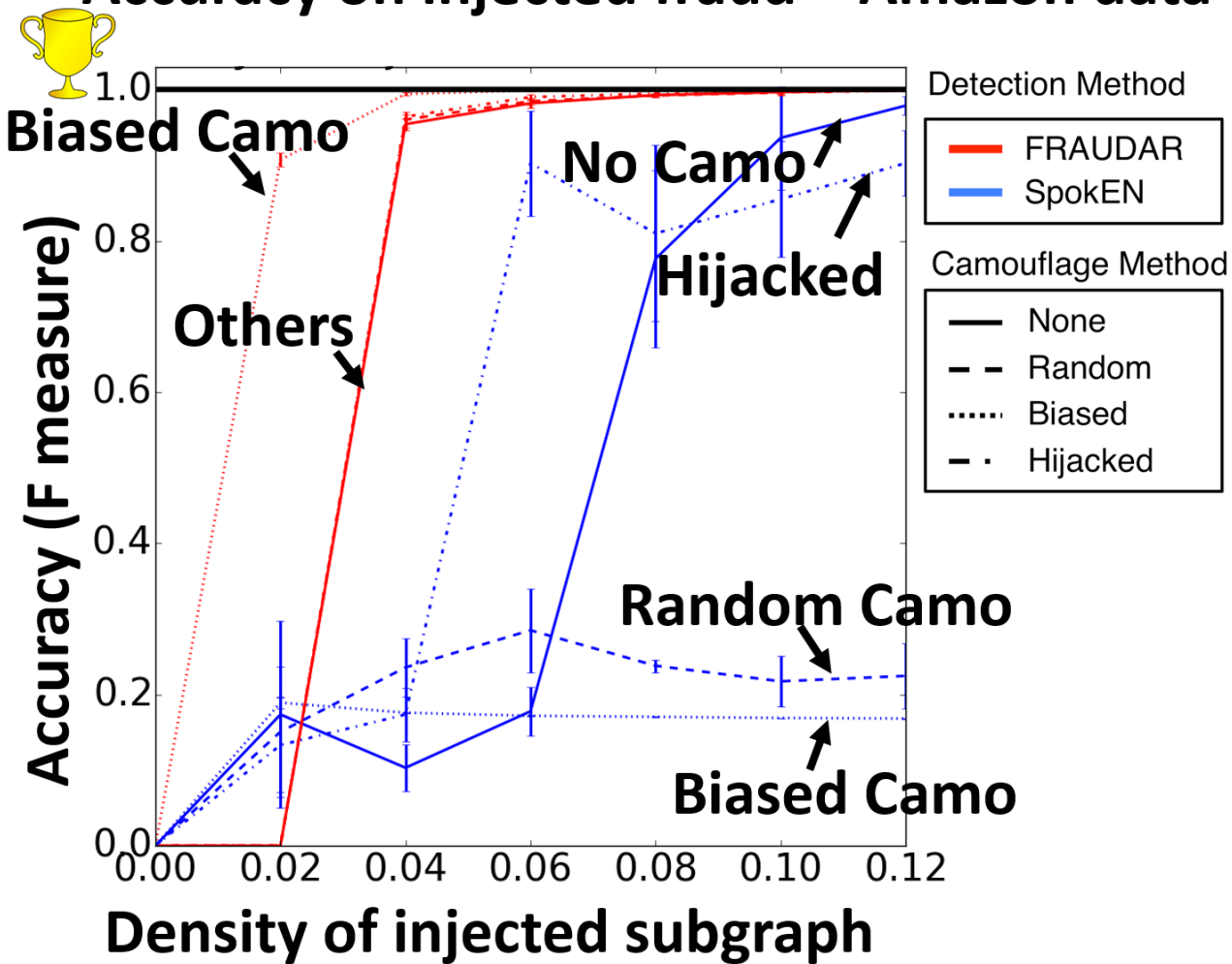
Experiments: Amazon data

- 24K x 4K Amazon review graph
- Injected dense blocks with various types of camouflage
 - None
 - Random camouflage
 - Biased camouflage
 - Hijacked accounts



Experiments: Amazon

Accuracy on injected fraud – Amazon data



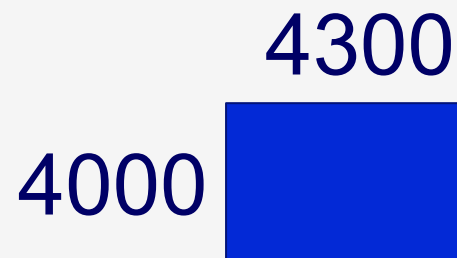


Twitter data

Followees

Density = 4×10^{-7}

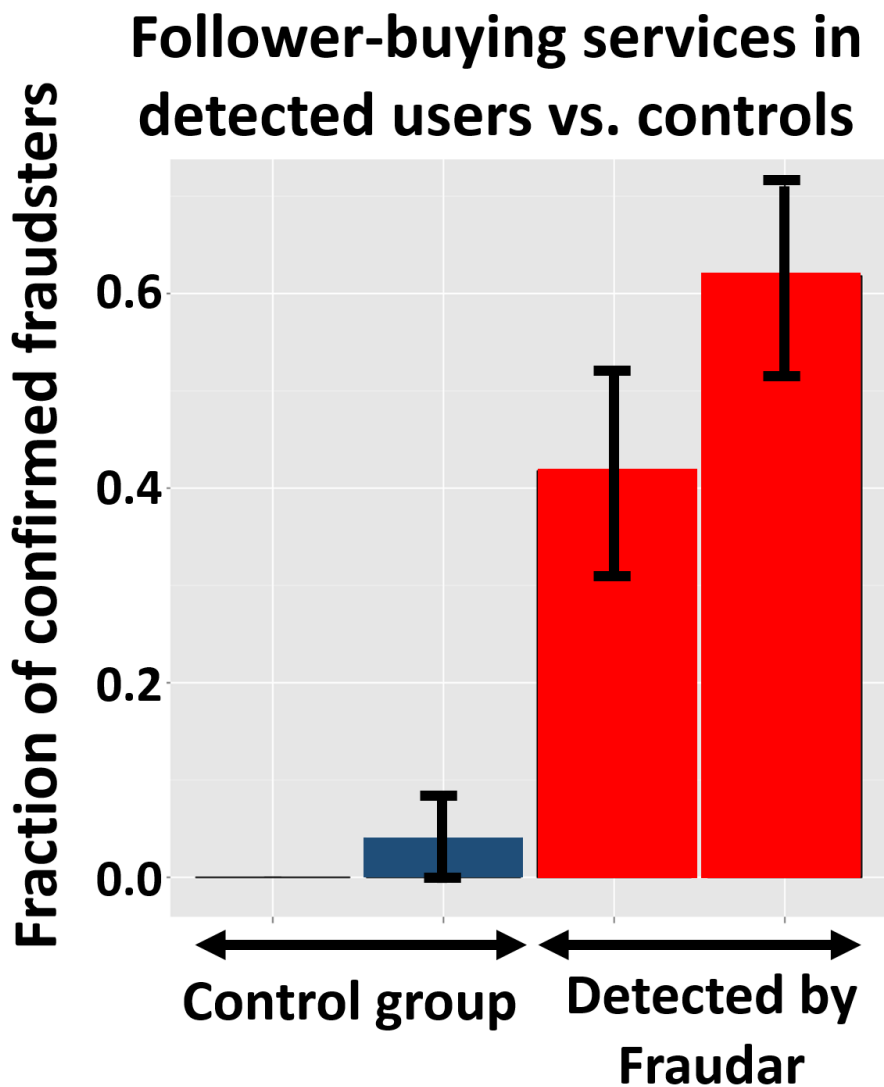
Followers



Density = 0.66



Twitter data



tweetergetter

Get 15,625 New Twitter Followers In 30 Days

"What If You Could Press Just One Button & Automatically Start Getting 1000's Of Legitimate New Twitter Followers On Autopilot... Even If Nobody Knows Who You Are Now?"

Buy A follower \$0.06 PER FOLLOWER \$0.10 PER FAN

FOR FOLLOWERS A \$5.00 START

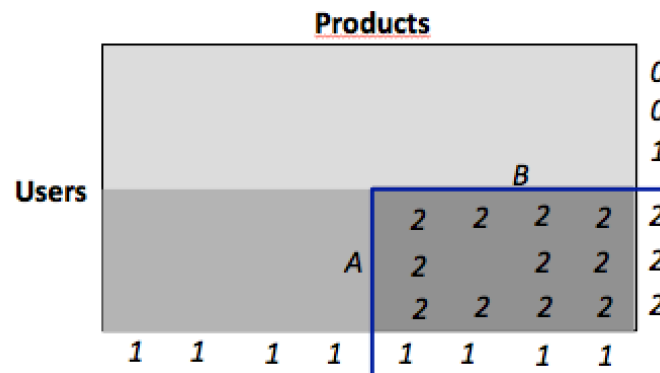
register your Twitter™ account

VALID THROUGH 6/30/10



Conclusion

- Average suspiciousness metric



$$g(A,B) = f(A,B) / (|A| + |B|)$$

$$g(A \cup B) \geq \frac{1}{2} g_{OPT}$$

- Theoretical guarantees

- Effectiveness

