

Anomaly detection in large graphs

Christos Faloutsos

CMU

Thank you!



- Prof. Richard Chbeir



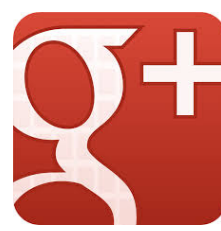
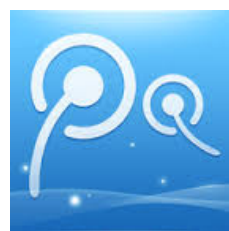
- Prof. Flavius Frasincar

Roadmap

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



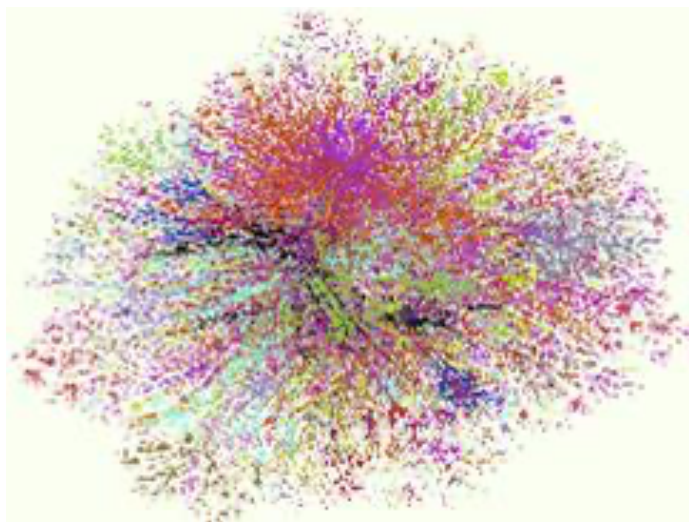
Graphs - why should we care?



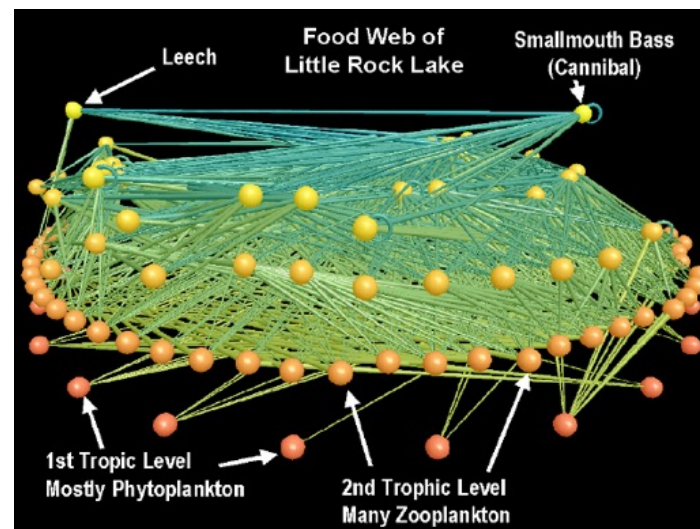
>\$10B; ~1B users



Graphs - why should we care?





Internet Map
[lumeta.com]



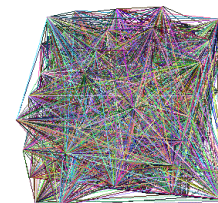
Food Web
[Martinez '91]

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

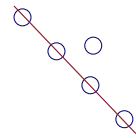


source

time

Motivating problems

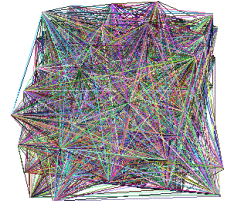
- P1: patterns? Fraud detection?



Patterns



anomalies



- P2: patterns in time-evolving graphs / tensors

destination



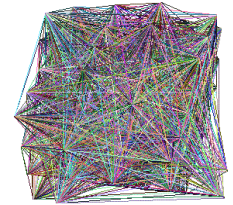
source

time

Roadmap

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

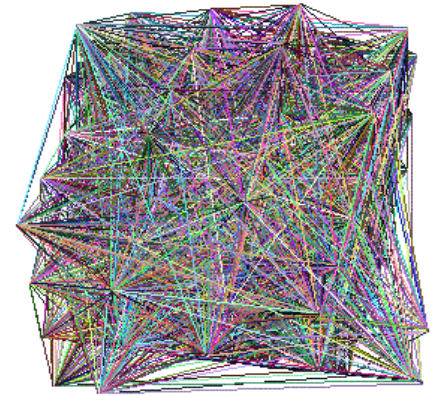




Part 1: Patterns, & fraud detection

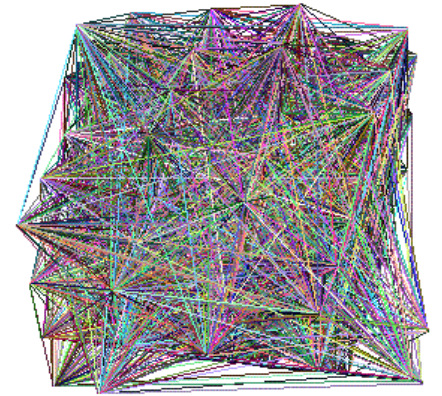
Laws and patterns

- Q1: Are real graphs random?



Laws and patterns

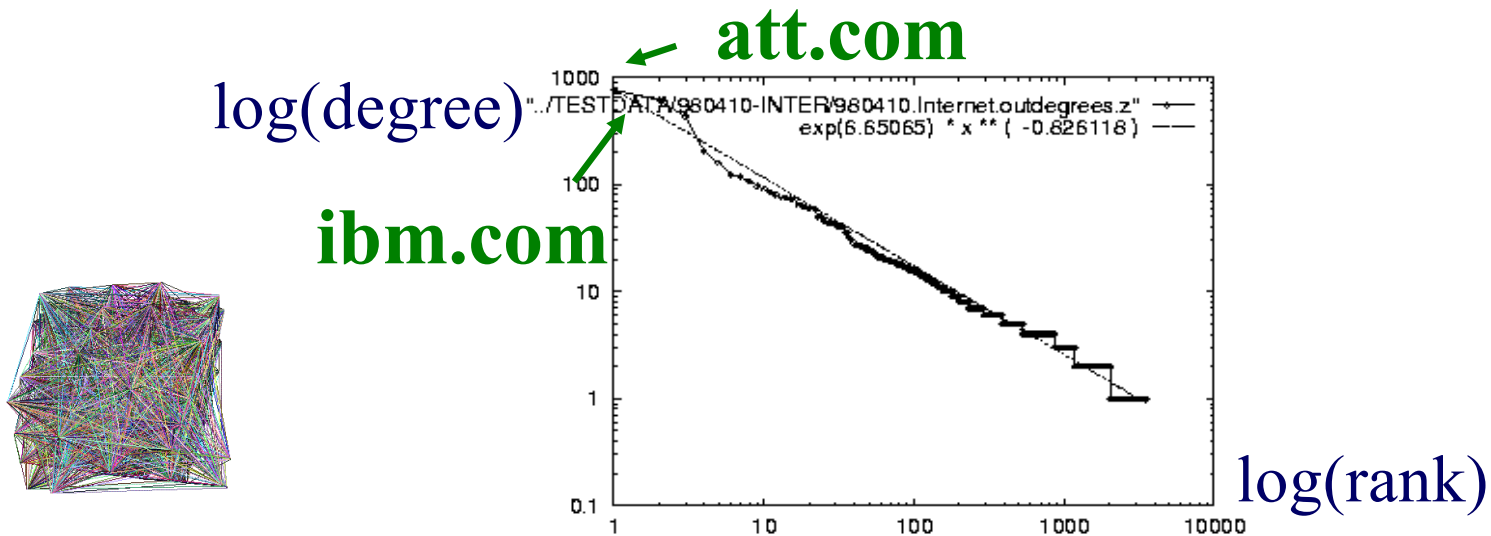
- Q1: Are real graphs random?
- A1: NO!!
 - Diameter ('6 degrees'; 'Kevin Bacon')
 - in- and out- degree distributions
 - other (surprising) patterns
- So, let's look at the data



Solution# S.1

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

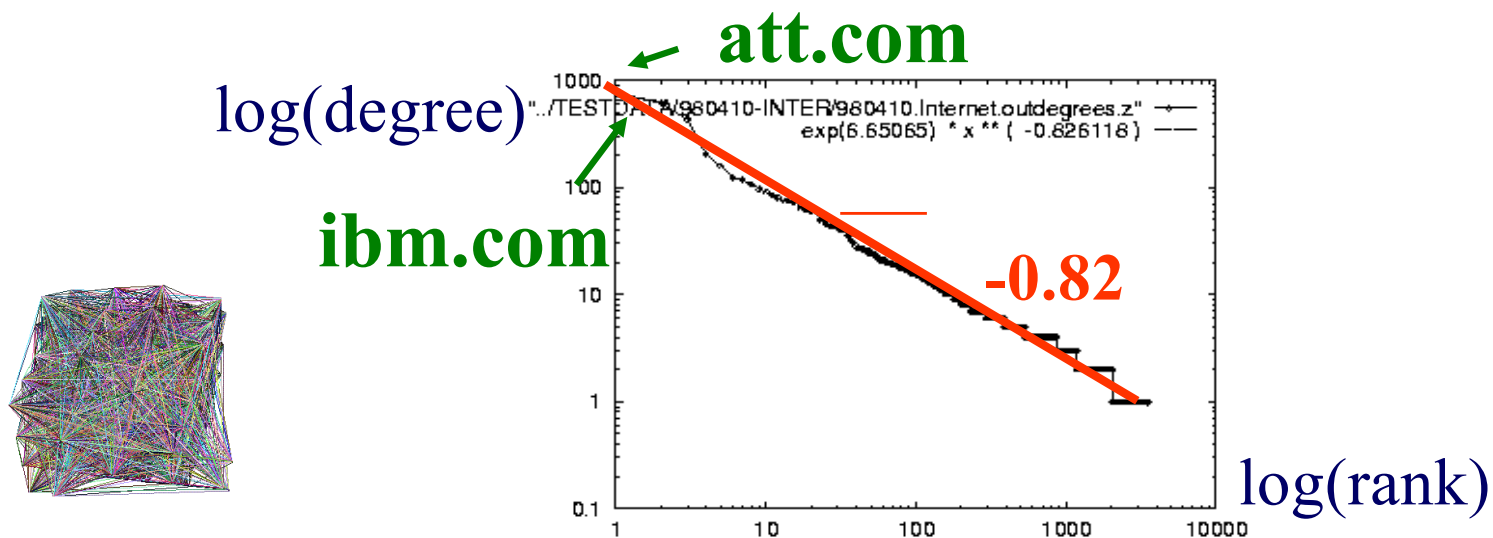
internet domains



Solution# S.1

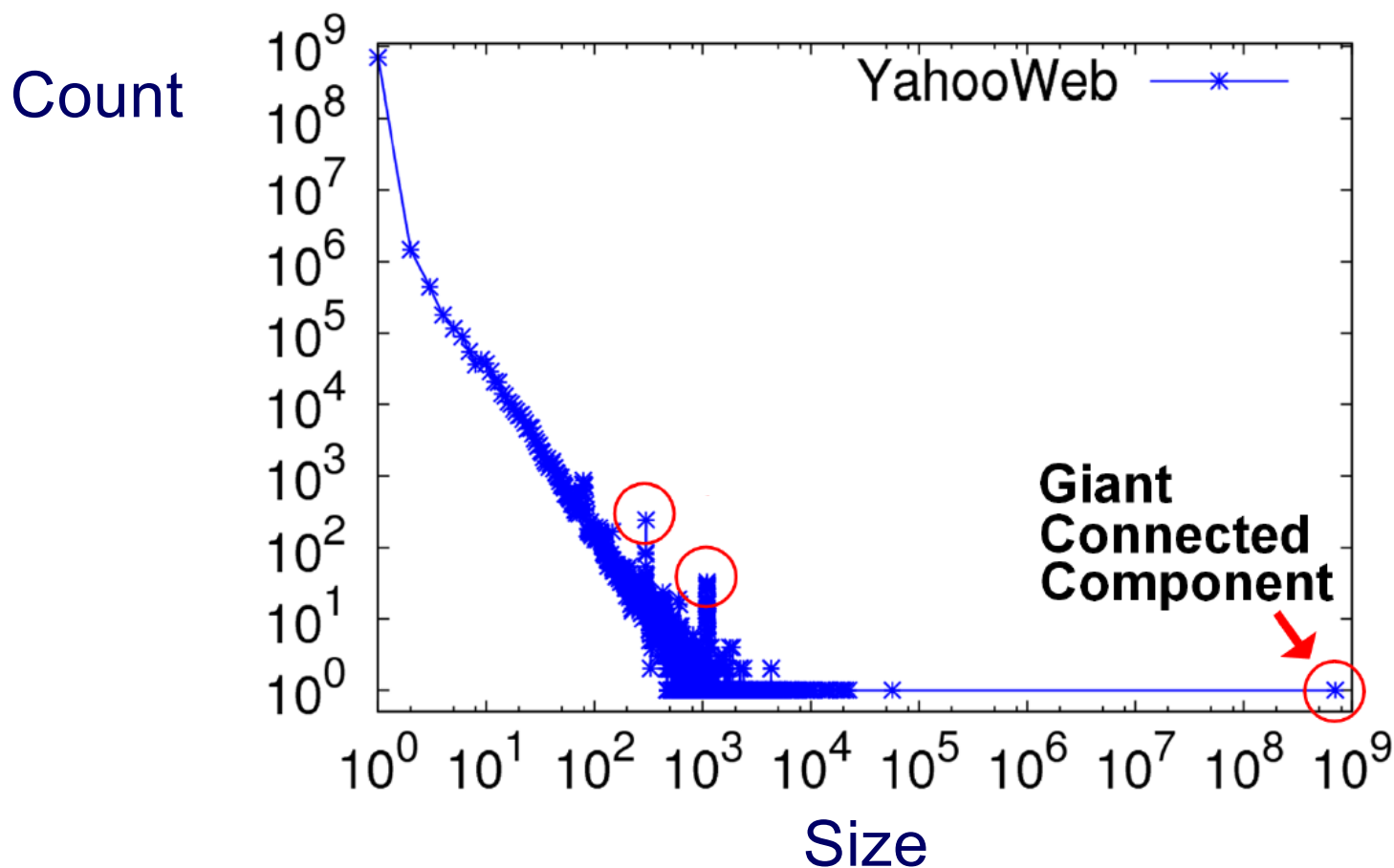
- Power law in the degree distribution [Faloutsos x 3 SIGCOMM99]

internet domains



S2: connected component sizes

- Connected Components – 4 observations:

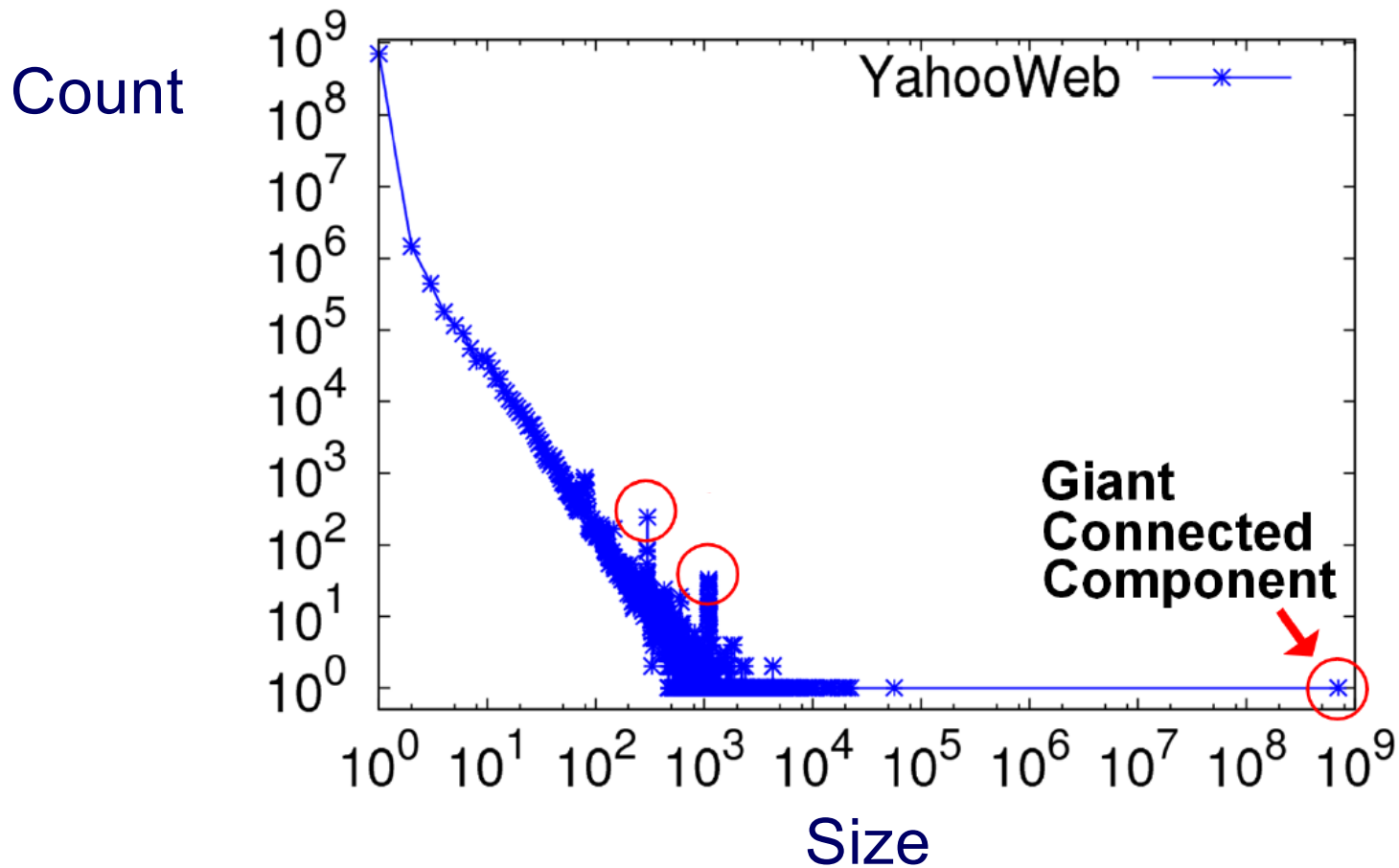


1.4B nodes
6B edges

S2: connected component sizes



- Connected Components

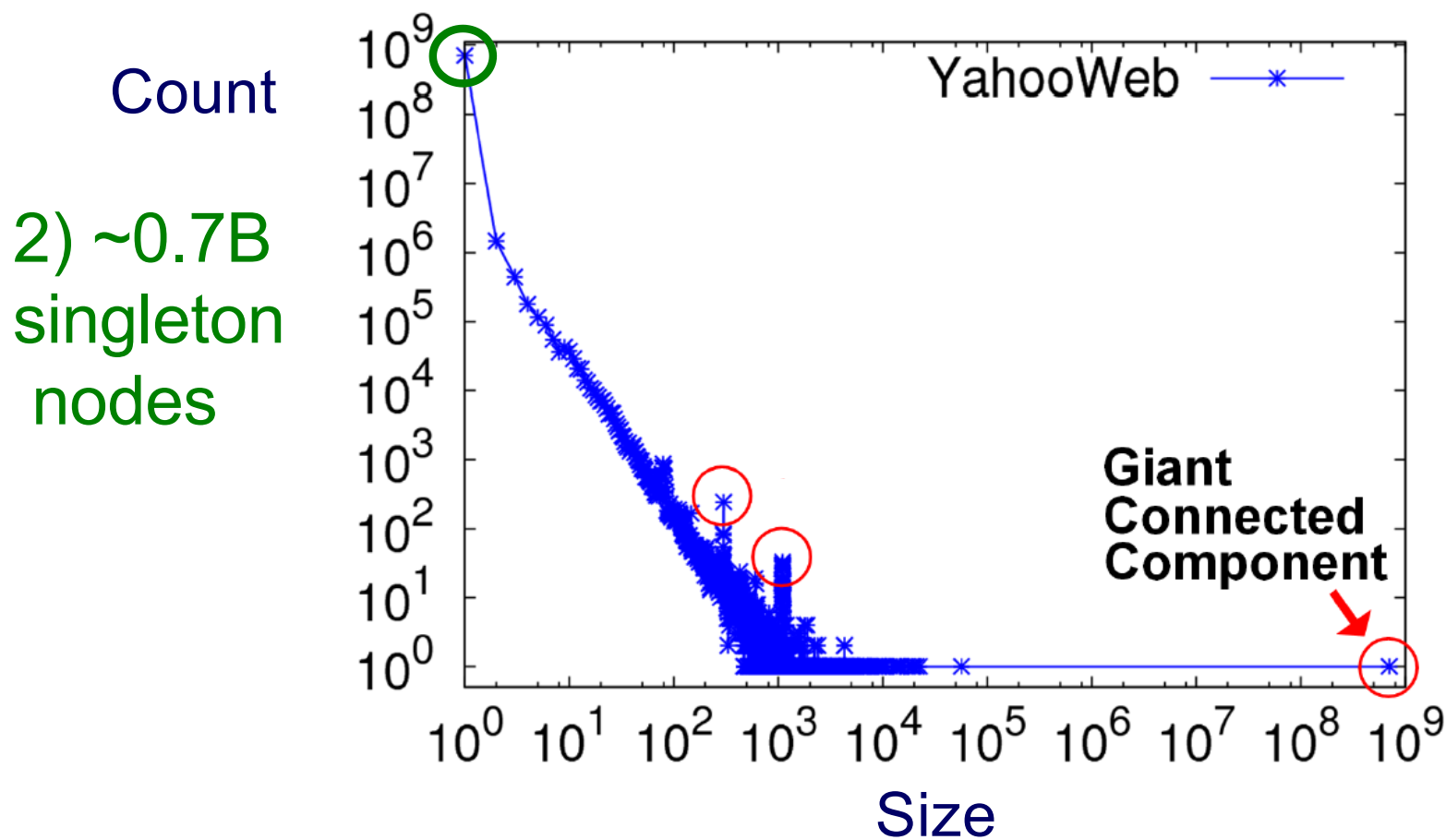


1) 10K x
larger
than next

S2: connected component sizes



- Connected Components

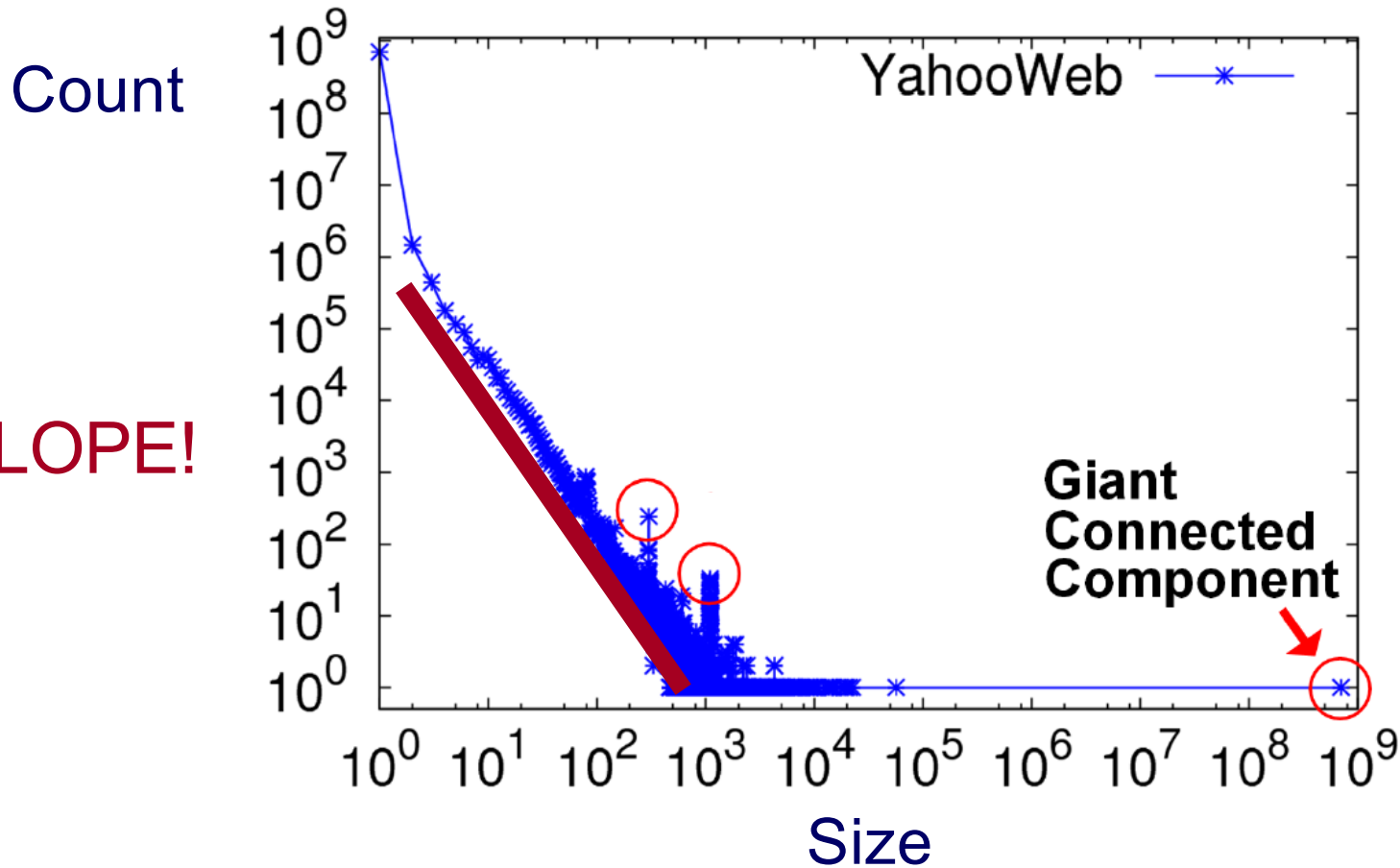


S2: connected component sizes



- Connected Components

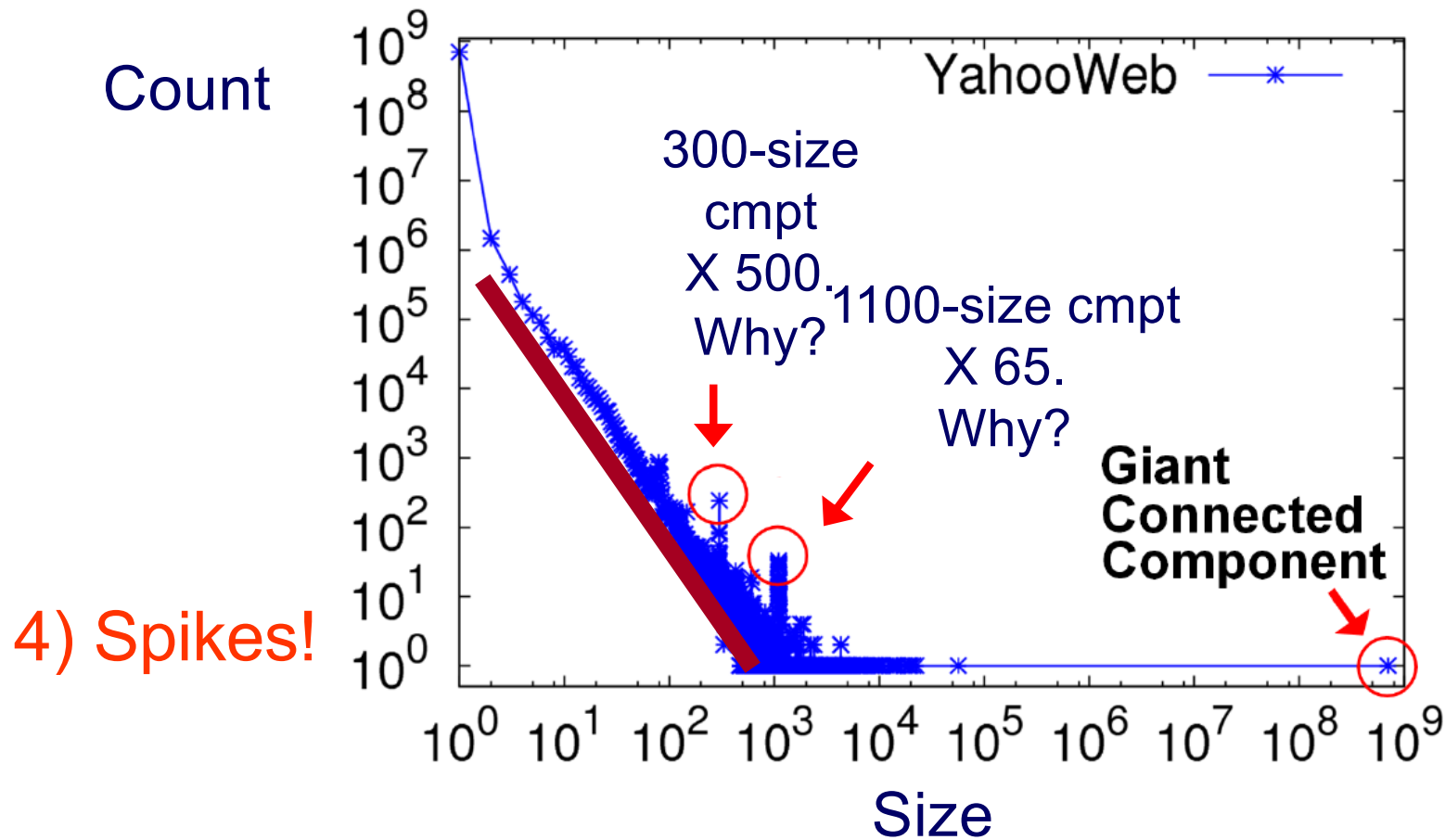
3) SLOPE!



S2: connected component sizes



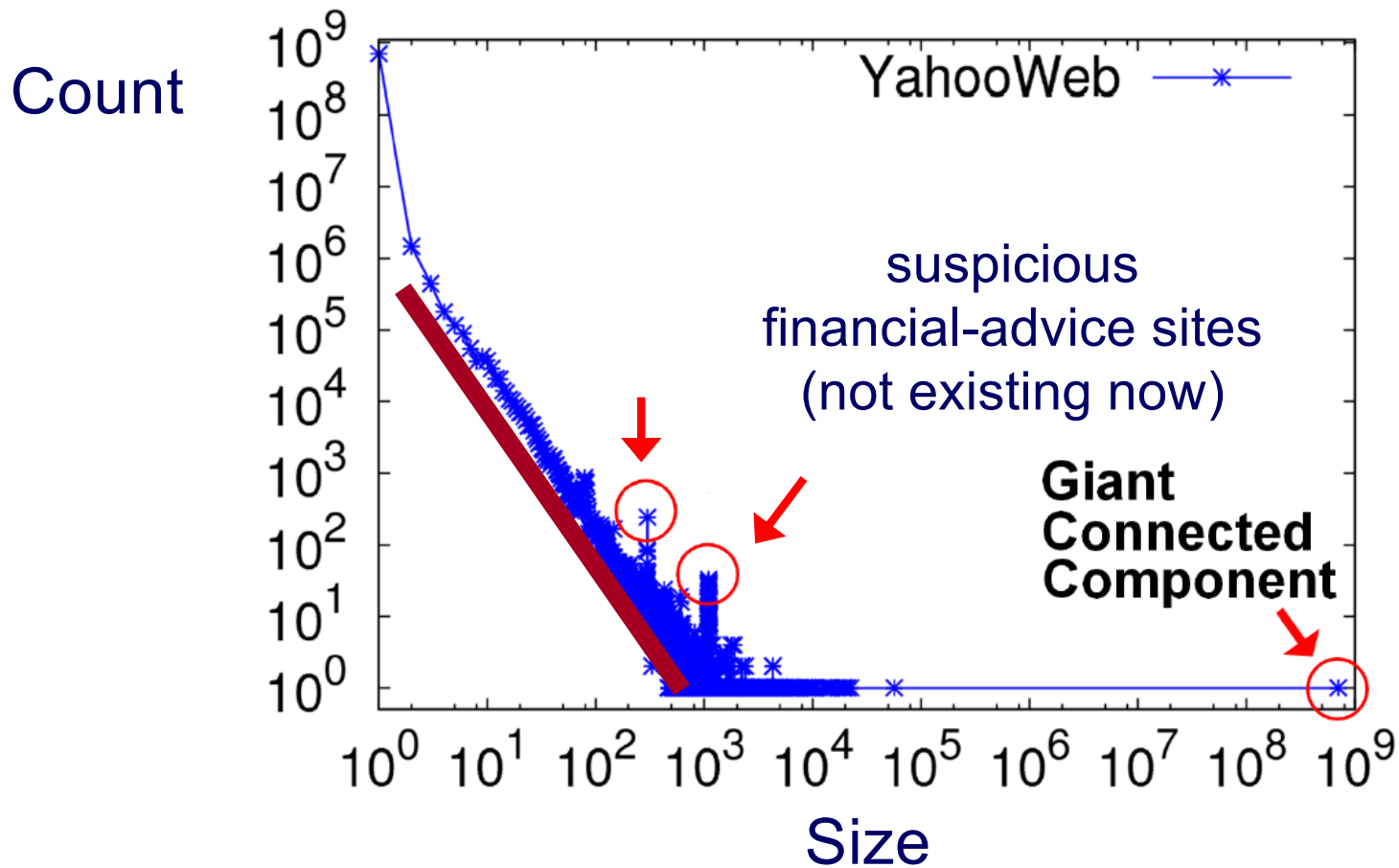
- Connected Components



S2: connected component sizes



- Connected Components

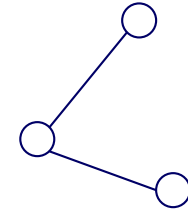


Roadmap



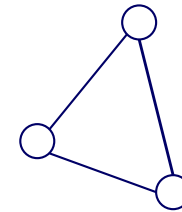
- Introduction – Motivation
- Part#1: Patterns in graphs
 - ➔ – P1.1: Patterns: Degree; Triangles
 - P1.2: Anomaly/fraud detection
- Part#2: time-evolving graphs; tensors
- Conclusions

Solution# S.3: Triangle ‘Laws’

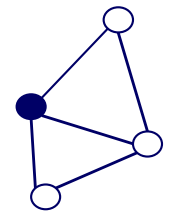


- Real social networks have a lot of triangles

Solution# S.3: Triangle ‘Laws’



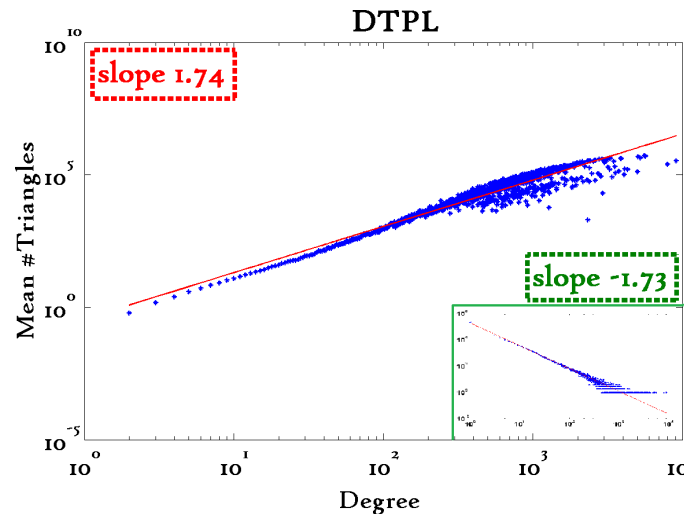
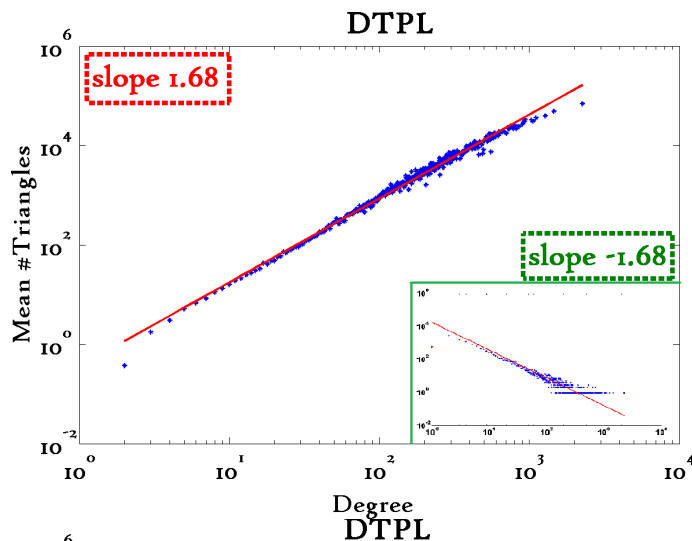
- Real social networks have a lot of triangles
 - Friends of friends are friends
- Any patterns?
 - 2x the friends, 2x the triangles ?



Triangle Law: #S.3

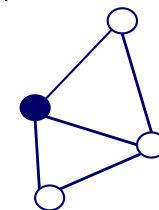
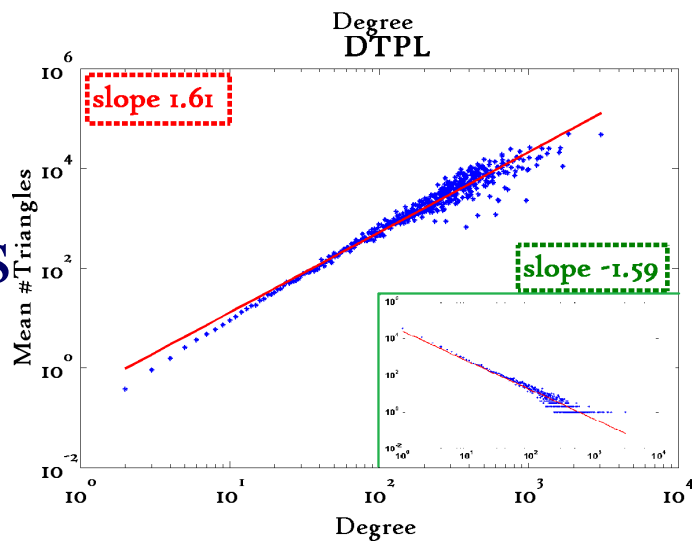
[Tsourakakis ICDM 2008]

Reuters



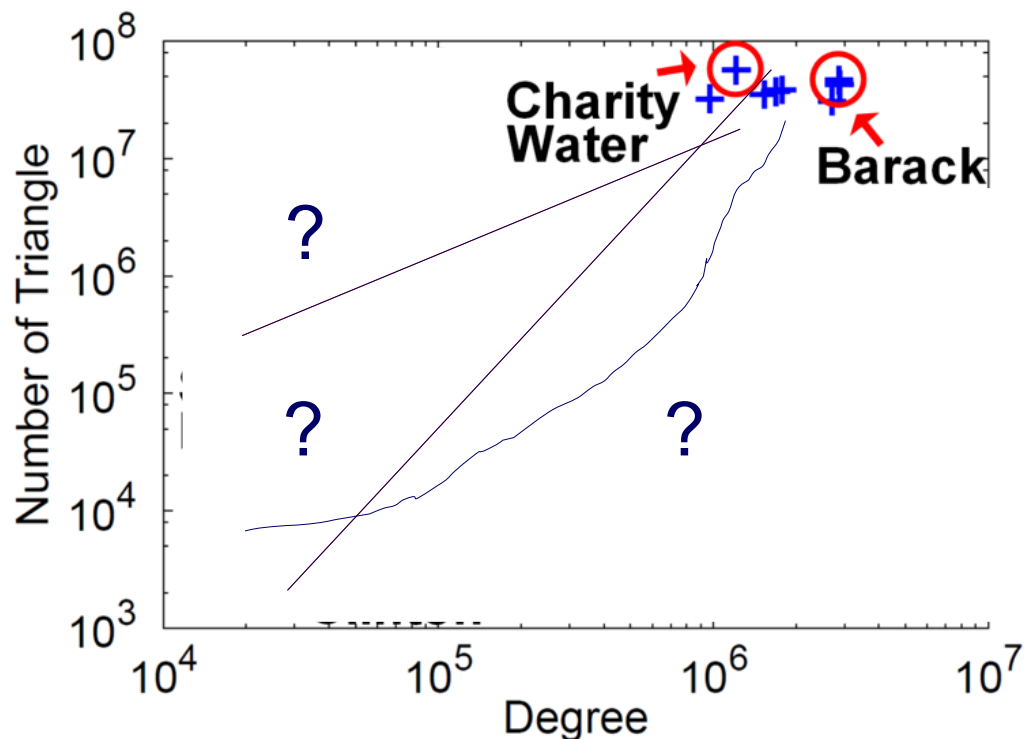
SN

Epinions



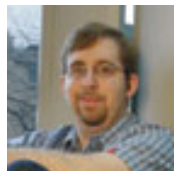
X-axis: degree
 Y-axis: mean # triangles
 n friends $\rightarrow \sim n^{1.6}$ triangles

Triangle counting for large graphs?

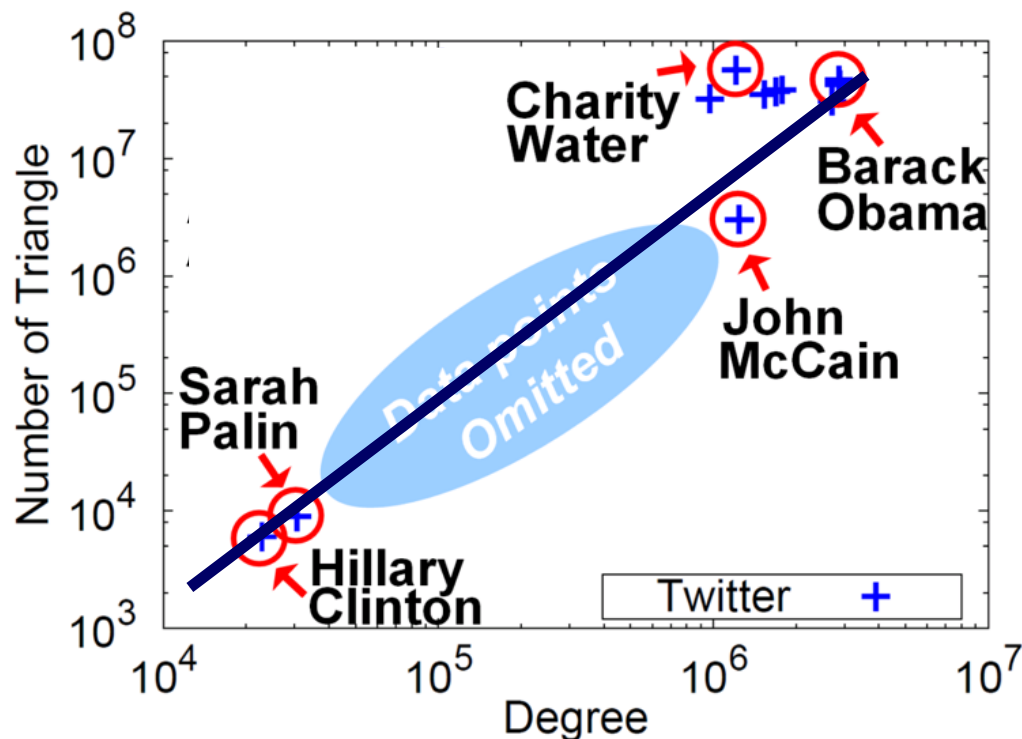


Anomalous nodes in Twitter (~ 3 billion edges)

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



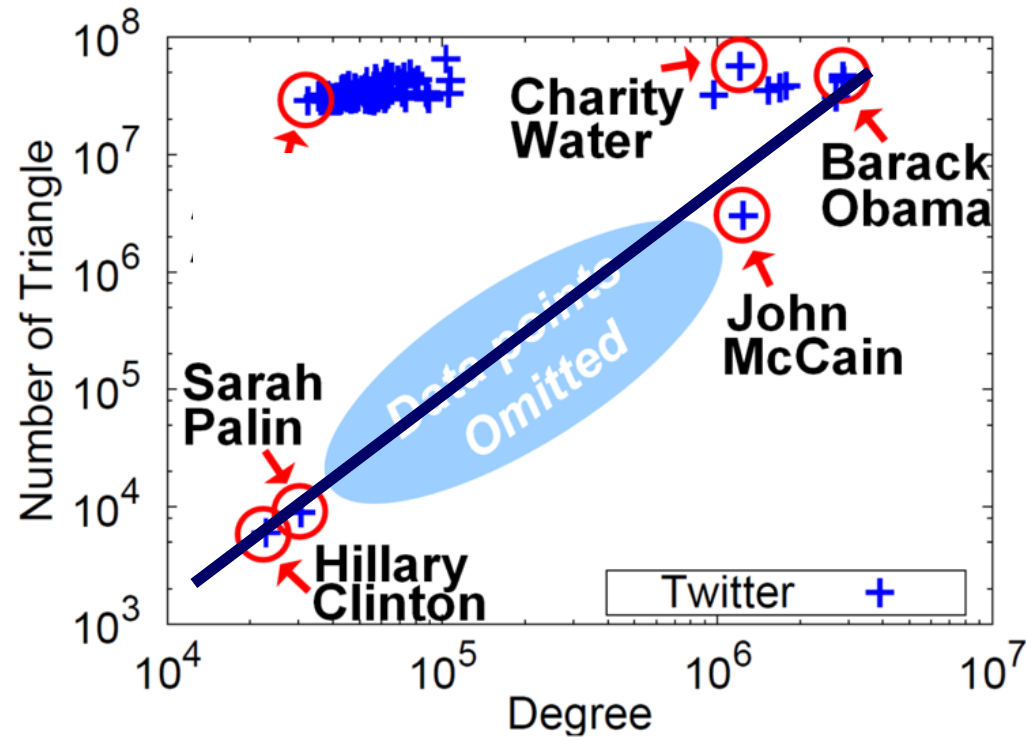
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

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

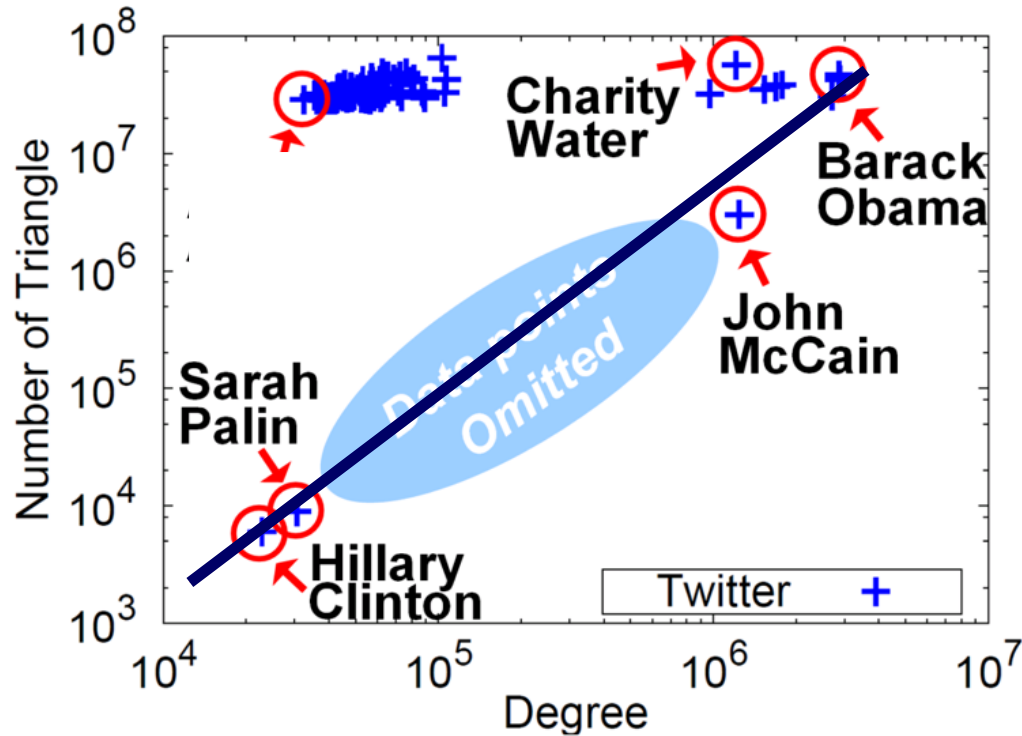
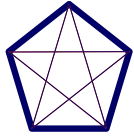
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

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

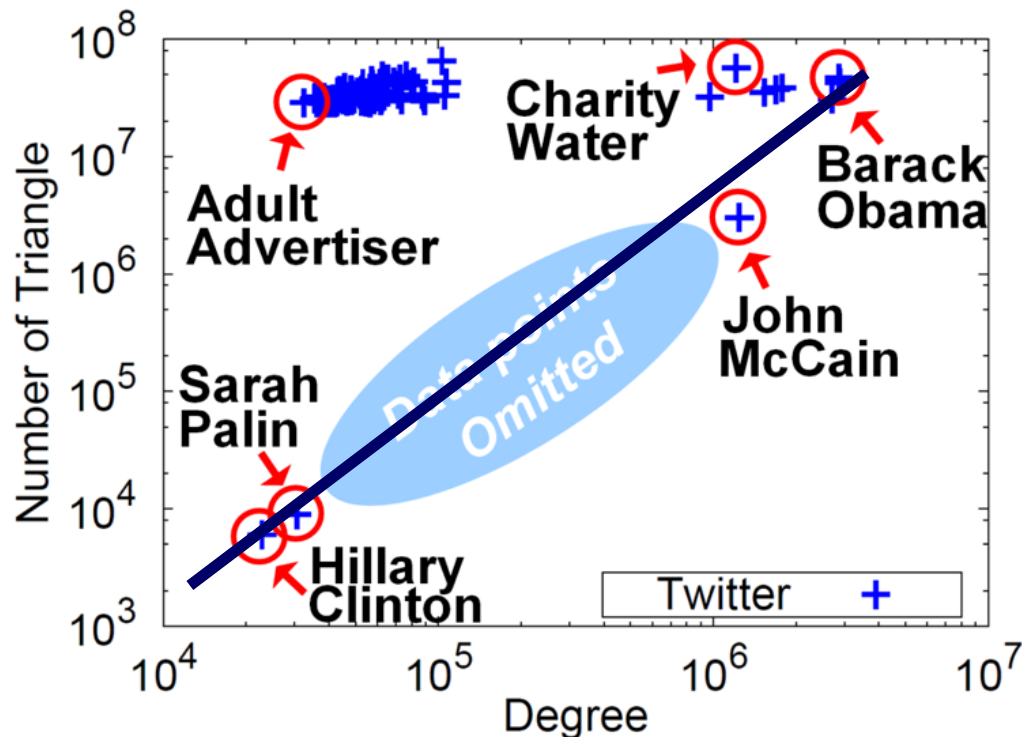
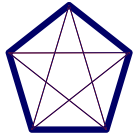
Triangle counting for large graphs?



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

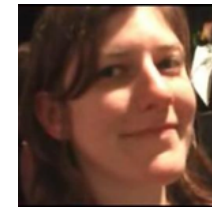
	Unweighted	Weighted
Static	<p>L01. Power-law degree distribution [Faloutsos et al. '99, Kleinberg et al. '99, Chakrabarti et al. '04, Newman '04]</p> <p>L02. Triangle Power Law (TPL) [Tsourakakis '08]</p> <p>L03. Eigenvalue Power Law (EPL) [Siganos et al. '03]</p> <p>L04. Community structure [Flake et 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
 - P1.1: Patterns
 - ➔ – P1.2: Anomaly / fraud detection
 - No labels – spectral
 - With labels: Belief Propagation
- Part#2: time-evolving graphs; tensors
- Conclusions



Patterns



anomalies

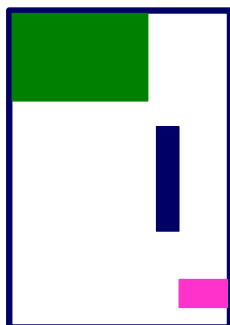
How to find ‘suspicious’ groups?

- ‘blocks’ are normal, right?



idols

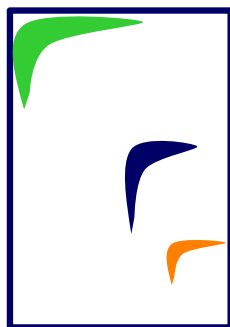
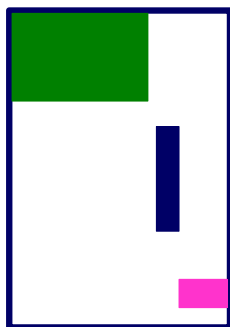
fans



Except that:



- ‘blocks’ are normal, ~~right?~~
- ‘hyperbolic’ communities are more realistic [Araujo+, PKDD’14]

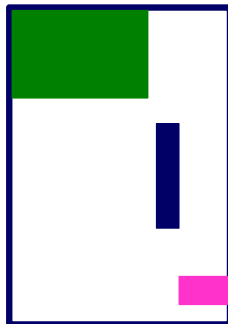


Except that:



- ‘blocks’ are usually **suspicious**
- ‘hyperbolic’ communities are more realistic
[Araujo+, PKDD’14]

Q: Can we spot blocks, easily?



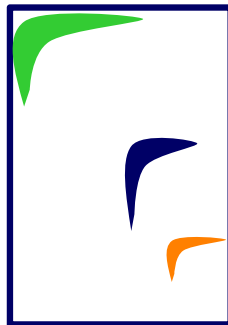
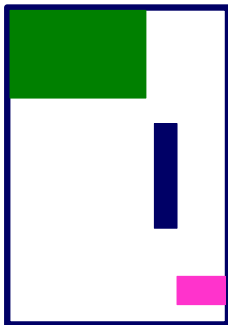
Except that:



- ‘blocks’ are usually **suspicious**
- ‘hyperbolic’ communities are more realistic
[Araujo+, PKDD’14]

Q: Can we spot blocks, easily?

A: Silver bullet: SVD!



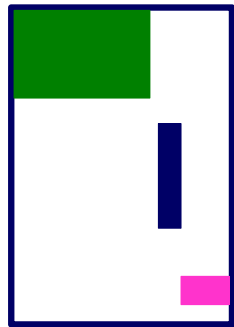
Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks



M
idols

N
fans

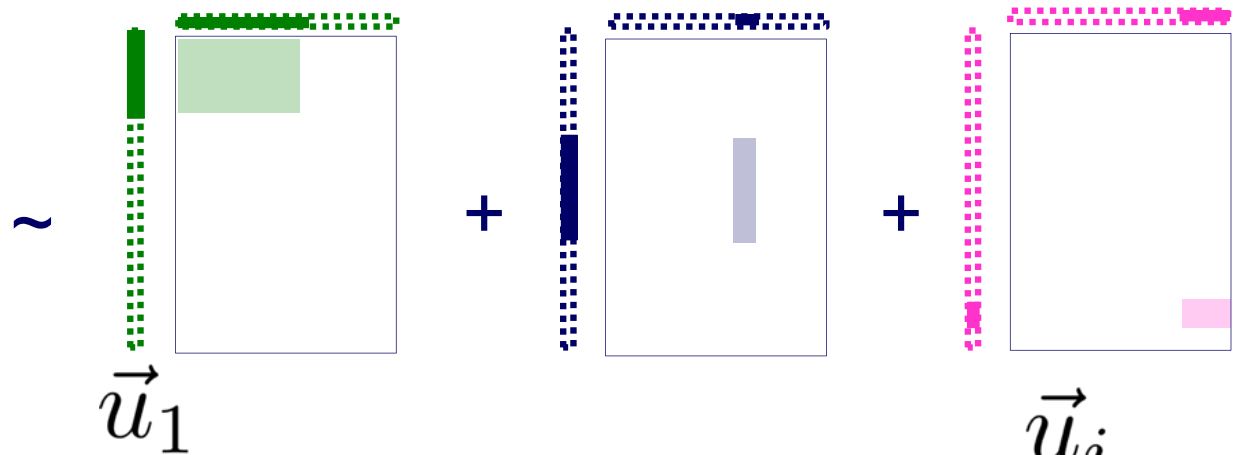


'music lovers'
'singers'

'sports lovers'
'athletes'

'citizens'
'politicians'

\vec{v}_1



\vec{u}_i 37

Christos Faloutsos

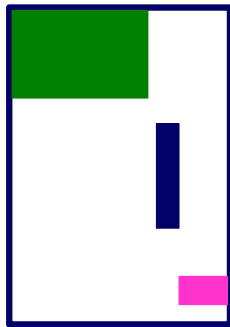
Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks



M
products

N
users

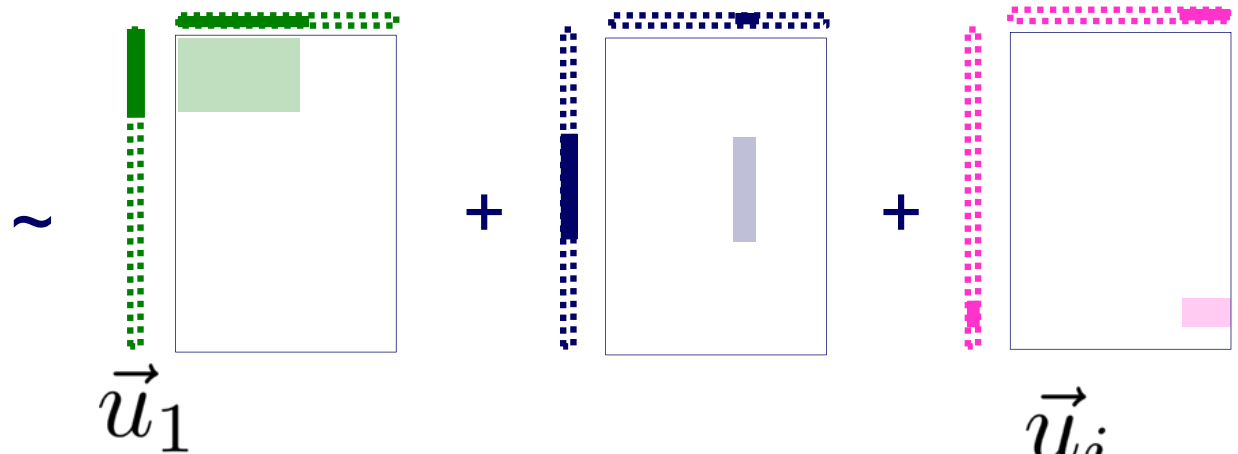


'meat-eaters'
'steaks'

\vec{v}_1

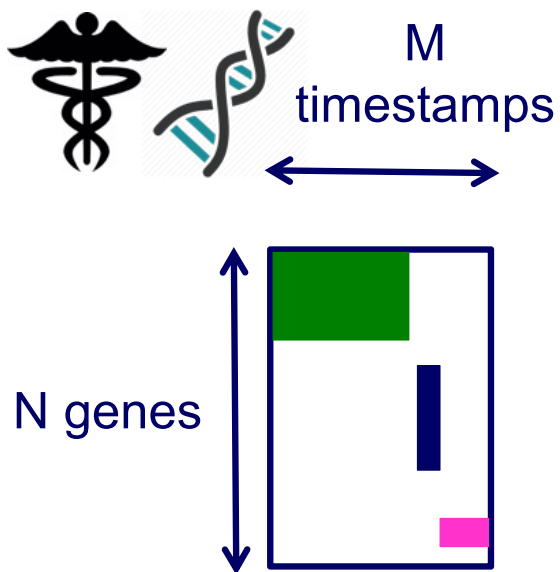
'vegetarians'
'plants'

'kids'
'cookies'

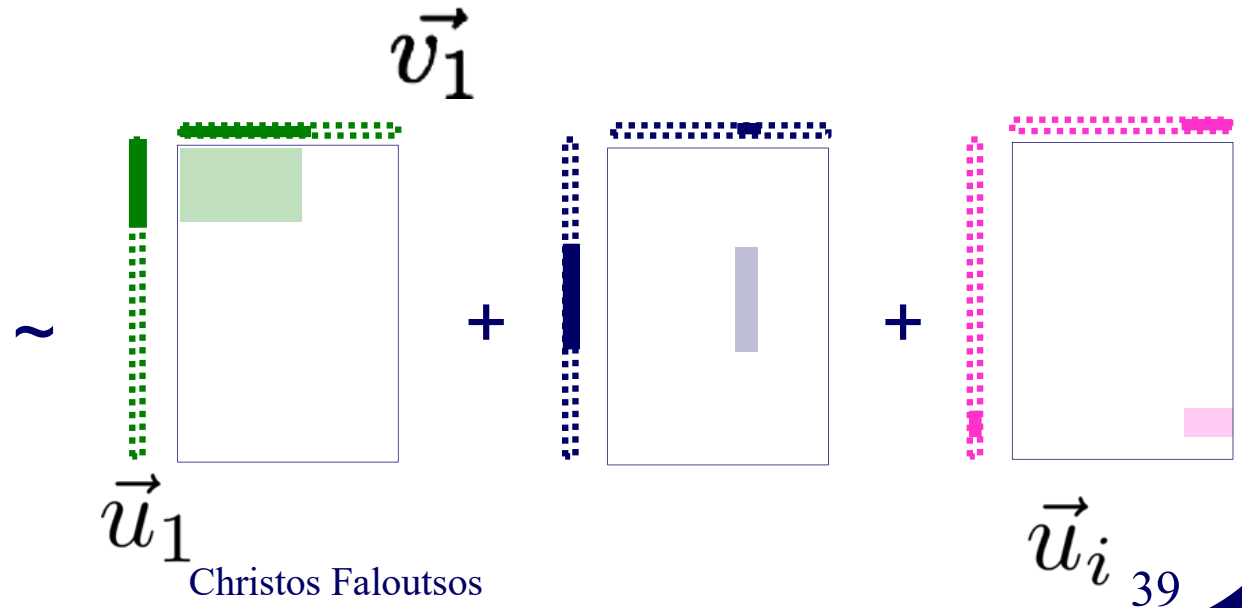


Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks

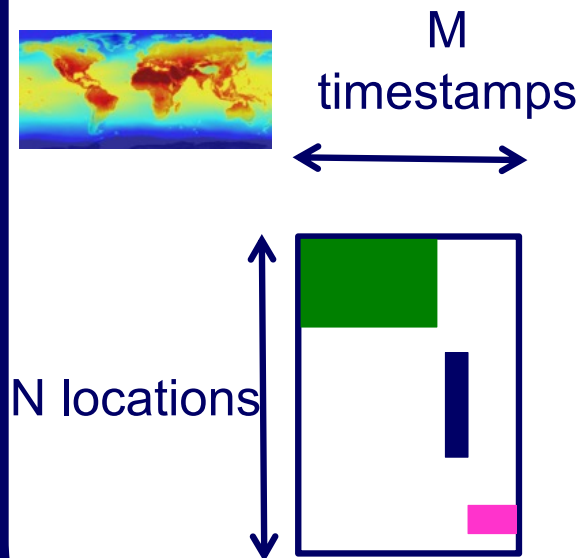


'cancer' 'alzheimer' 'Parkinson'



Crash intro to SVD

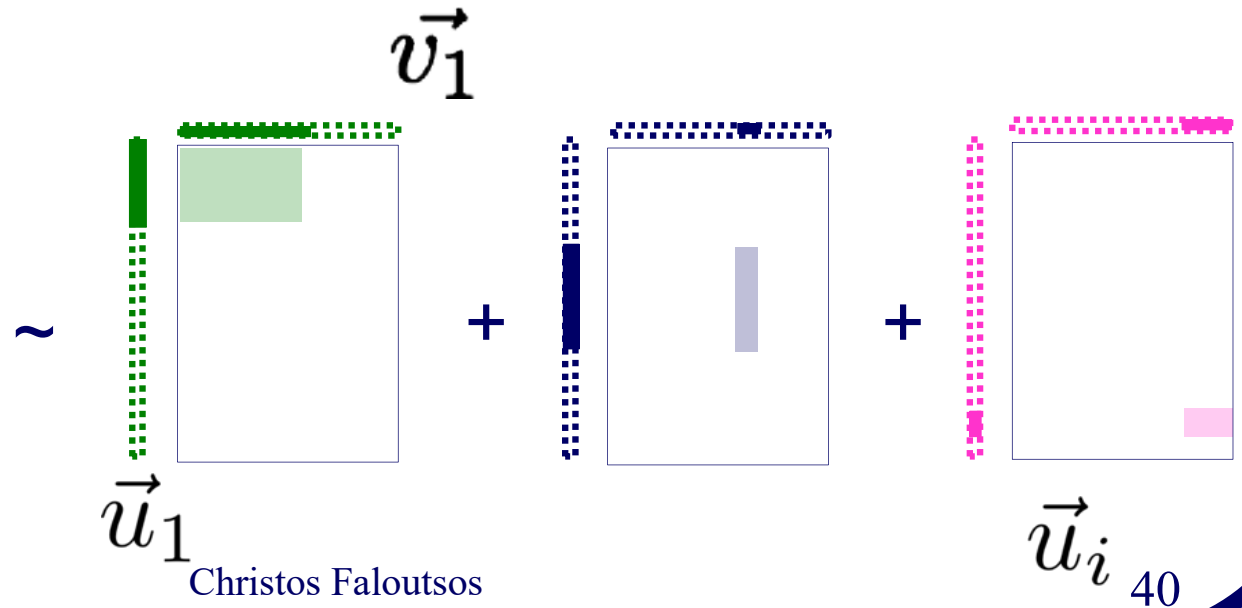
- Recall: (SVD) matrix factorization: finds blocks



'hurricane'

'cold-spell'

'heat-wave'



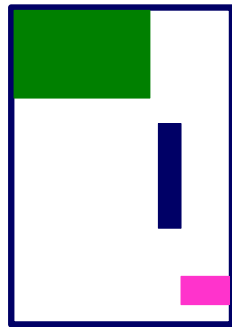
Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks

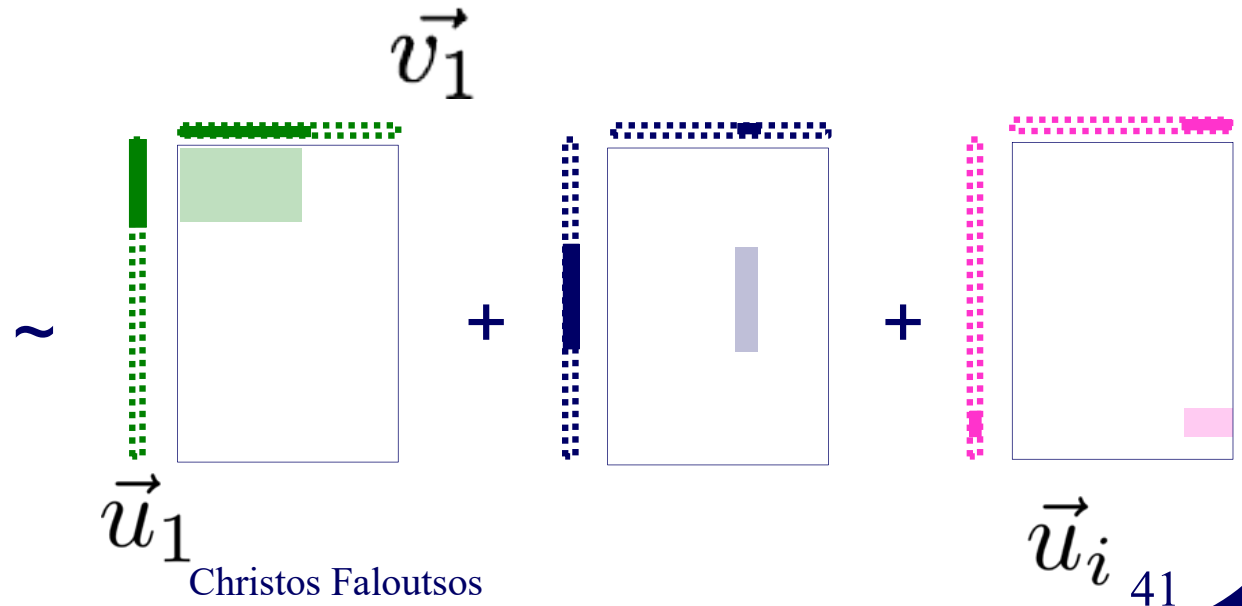


M
idols

N
fans



'music lovers' 'singers'
'sports lovers' 'athletes'
'citizens' 'politicians'



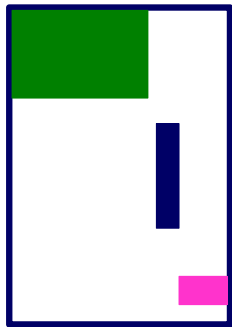
Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks



M
idols

N
fans



'music lovers'
'singers'

'sports lovers'
'athletes'

'citizens'
'politicians'

$$\sim \vec{u}_1 + \vec{v}_1 + \vec{u}_i$$

Christos Faloutsos

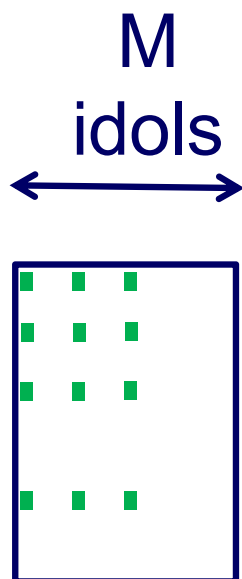
42

Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks **Even if shuffled!**



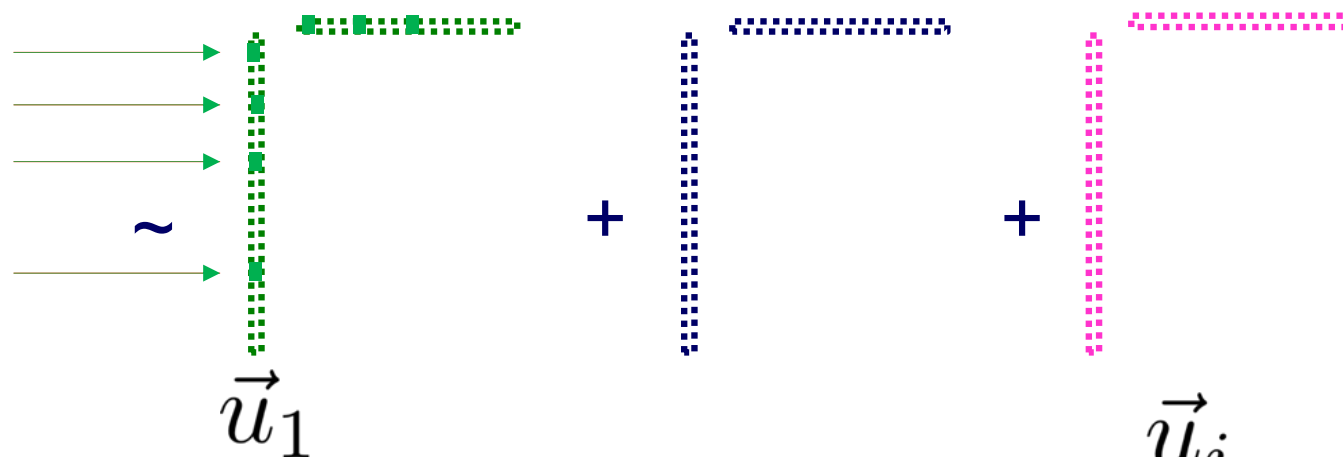
N
fans



'music lovers'
'singers'

'sports lovers'
'athletes'

'citizens'
'politicians'

 \vec{v}_1


Inferring Strange Behavior from Connectivity Pattern in Social Networks


PAKDD'14



Meng Jiang, Peng Cui, Shiqiang Yang (Tsinghua)
Alex Beutel, Christos Faloutsos (CMU)



Dataset

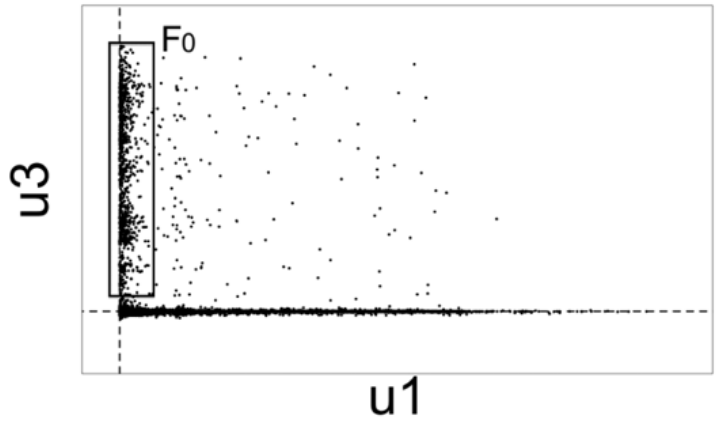
- Tencent Weibo 
- 117 million nodes (with profile and UGC data)
- 3.33 billion directed edges



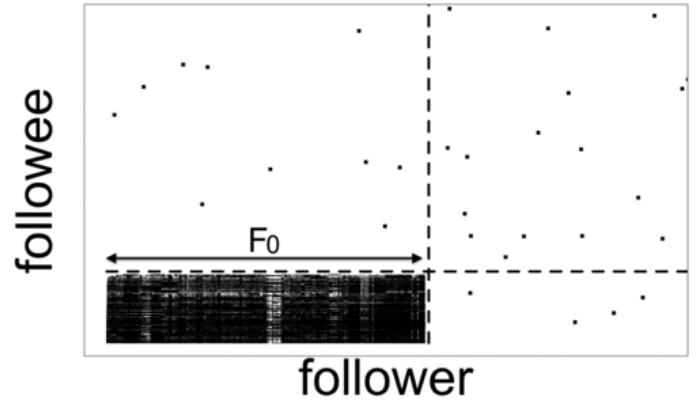
Real Data



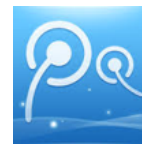
“Rays”



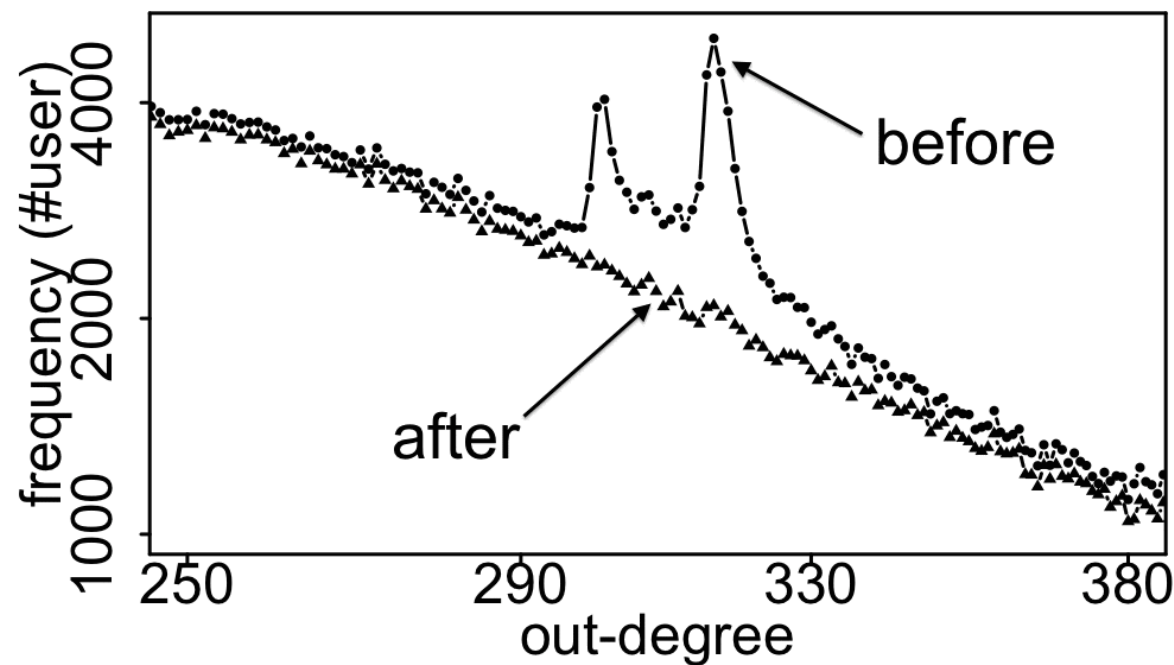
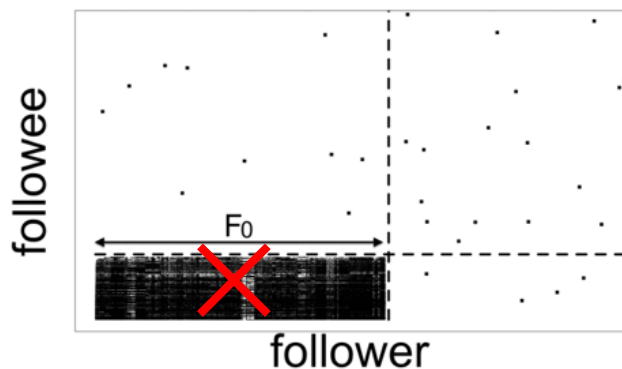
“Block”



Real Data



- Spikes on the out-degree distribution



Roadmap

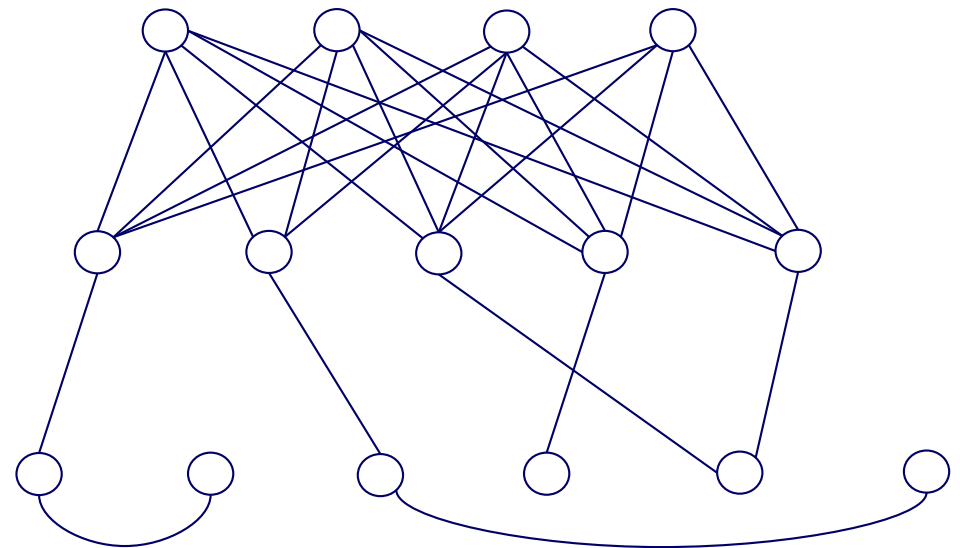
- Introduction – Motivation
- Part#1: Patterns in graphs
 - P1.1: Patterns
 - P1.2: Anomaly / fraud detection
 - No labels – spectral methods
 - With labels: Belief Propagation
- ➔ • Part#2: time-evolving graphs; tensors
- Conclusions



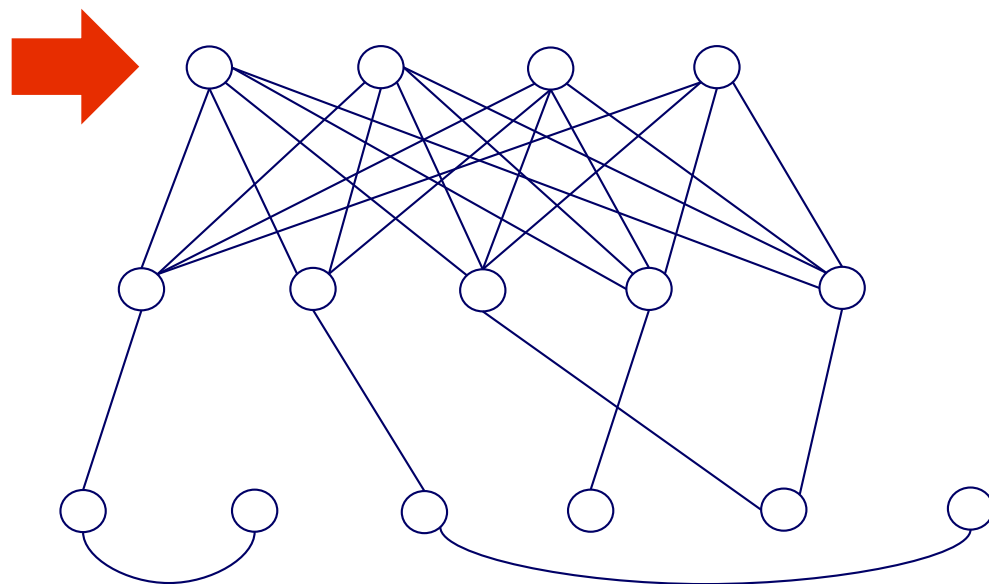
E-bay Fraud detection



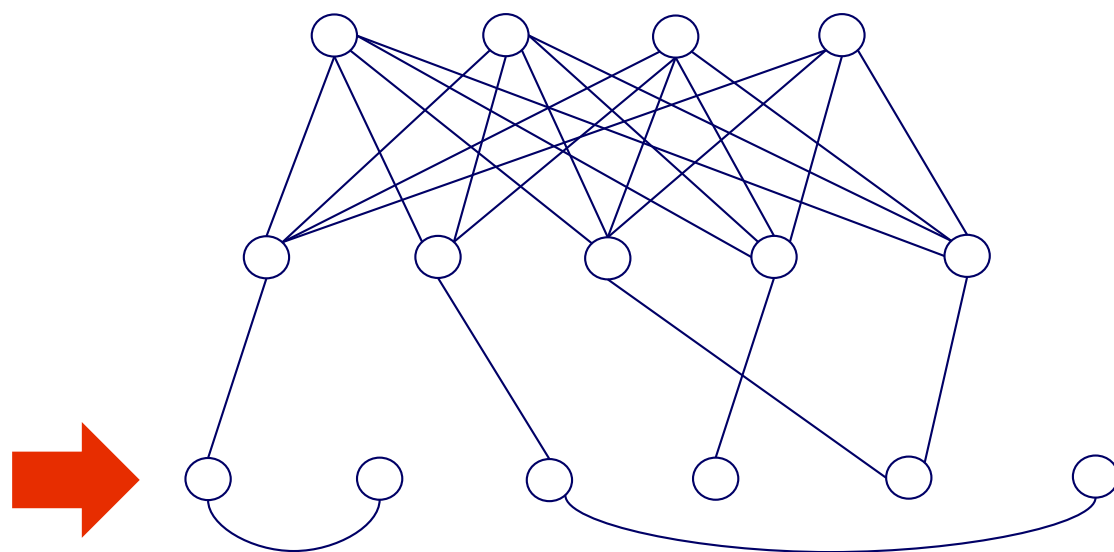
w/ Polo Chau &
Shashank Pandit, CMU
[www'07]



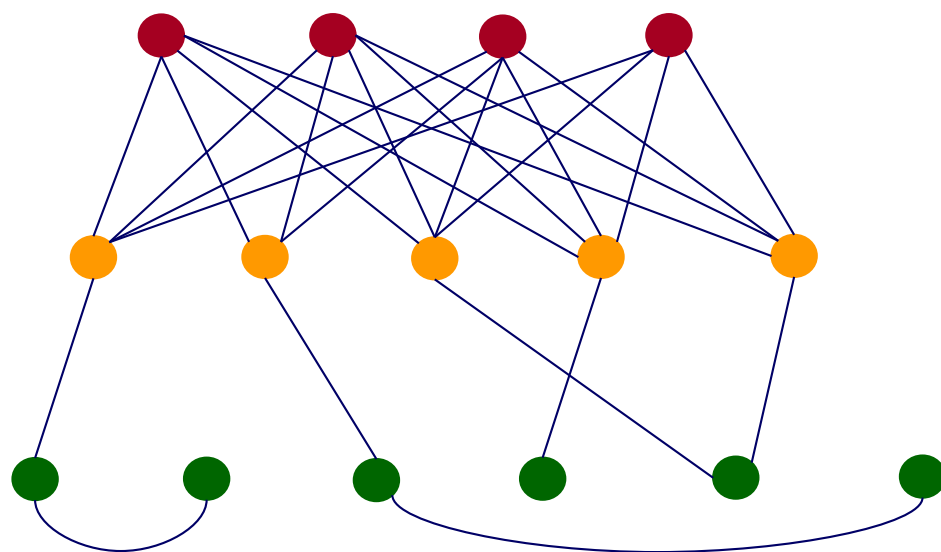
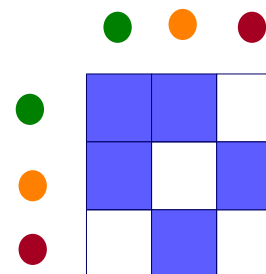
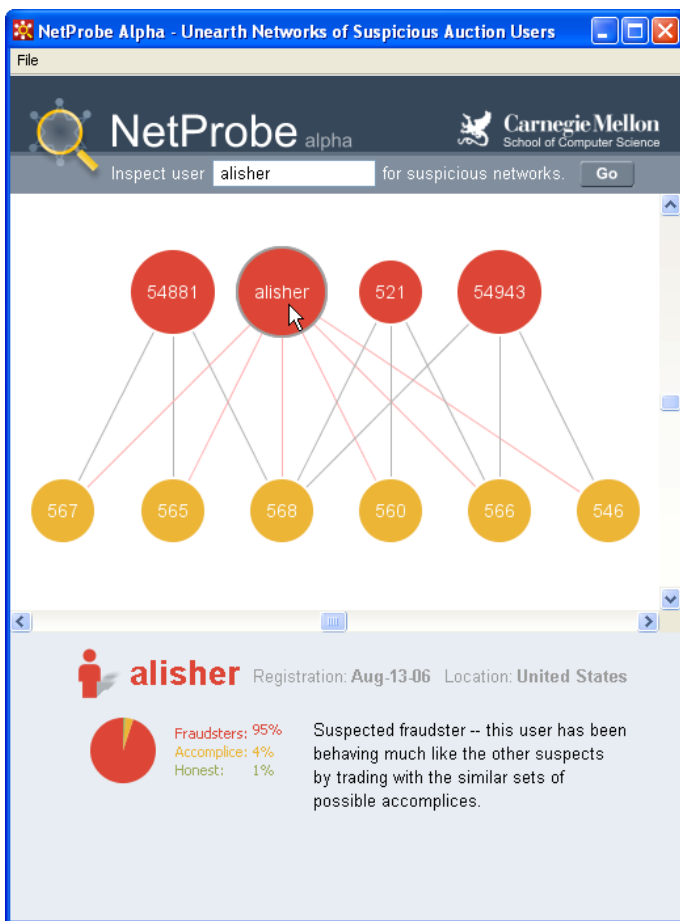
E-bay Fraud detection



E-bay Fraud detection



E-bay Fraud detection - NetProbe



Popular press



The Washington Post

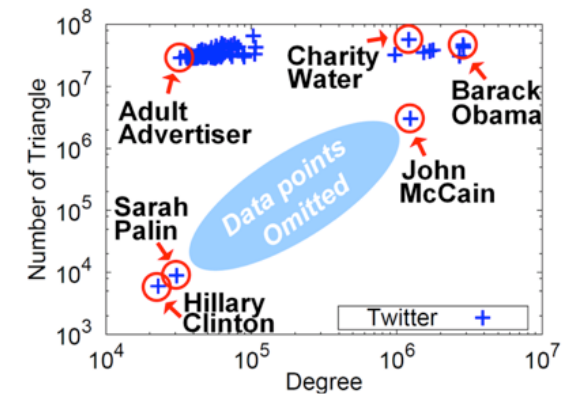
Los Angeles Times

And less desirable attention:

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

Summary of Part#1

- **many** patterns in real graphs
 - Power-laws everywhere
 - Long (and growing) list of tools for anomaly/fraud detection



Patterns



anomalies

Roadmap

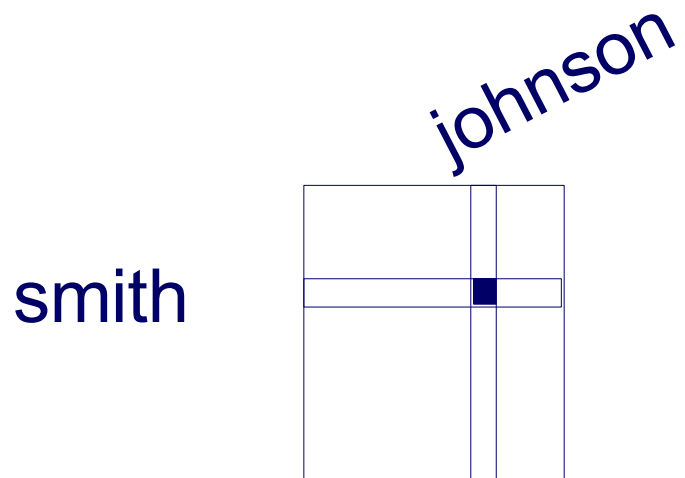
- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs
 - ➔ – P2.1: tools/tensors
 - P2.2: other patterns
- Conclusions



Part 2: Time evolving graphs; tensors

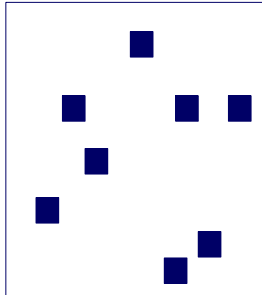
Graphs over time -> tensors!

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



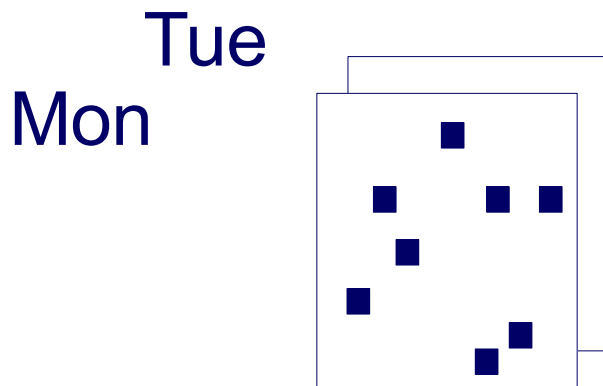
Graphs over time \rightarrow tensors!

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



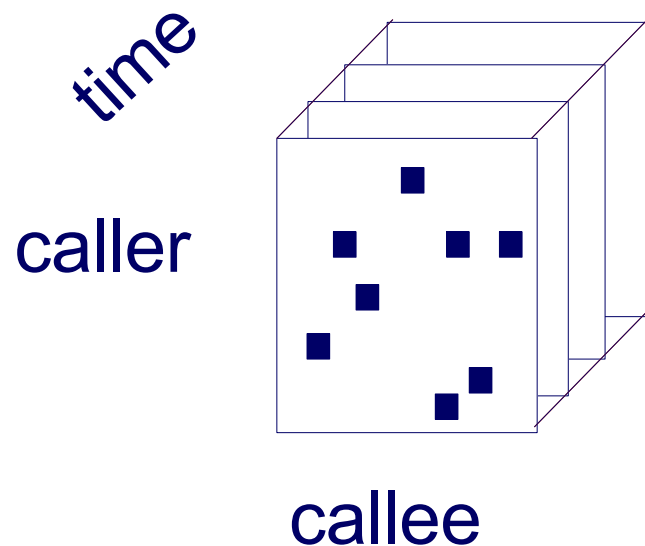
Graphs over time -> tensors!

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



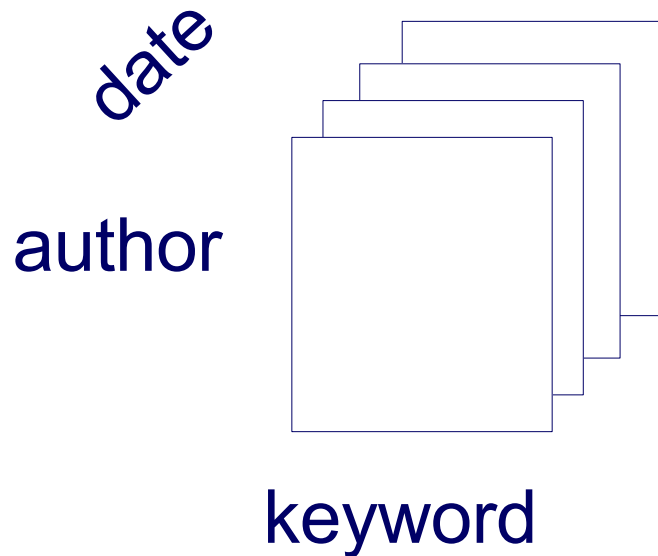
Graphs over time -> tensors!

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



Graphs over time -> tensors!

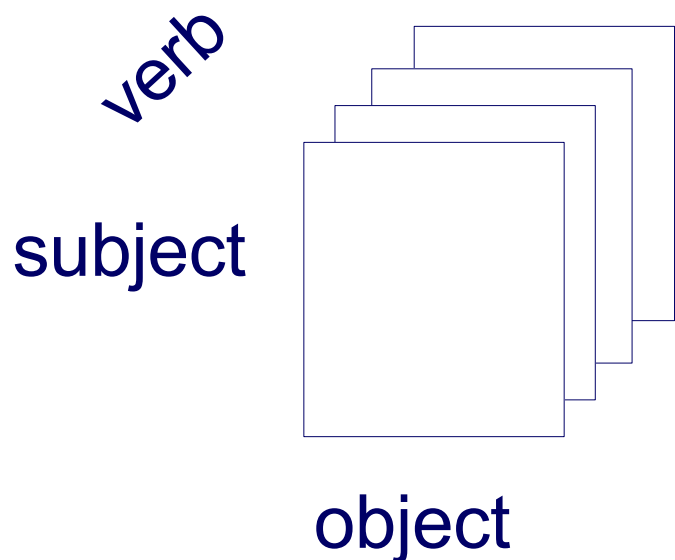
- Problem #2.1':
 - Given author-keyword-date
 - Find patterns / anomalies



MANY more settings,
with >2 'modes'

Graphs over time -> tensors!

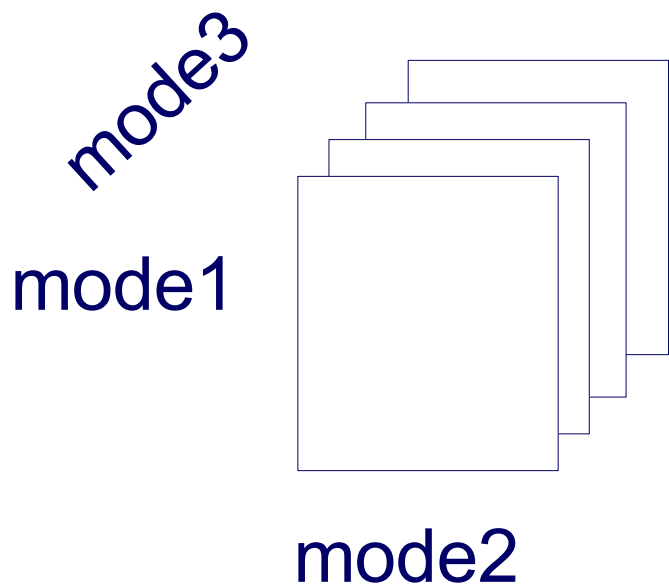
- Problem #2.1’’:
 - Given subject – verb – object facts
 - Find patterns / anomalies



MANY more settings,
with >2 ‘modes’

Graphs over time -> tensors!

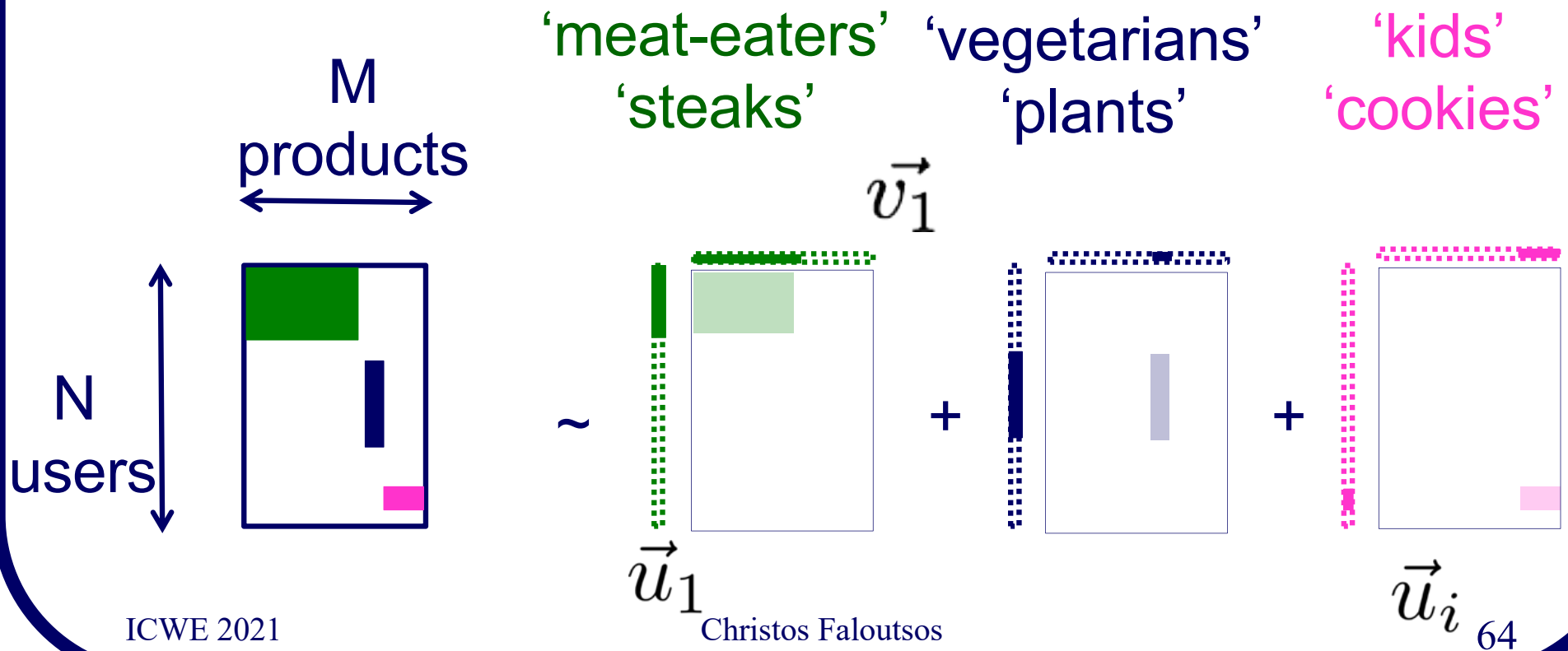
- Problem #2.1''':
 - Given <triplets>
 - Find patterns / anomalies



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

Answer : tensor factorization

- Recall: (SVD) matrix factorization: finds blocks



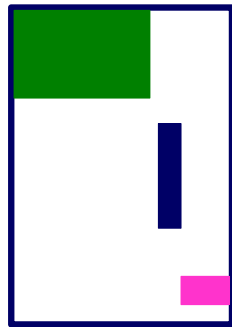
Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks



M
idols

N
fans

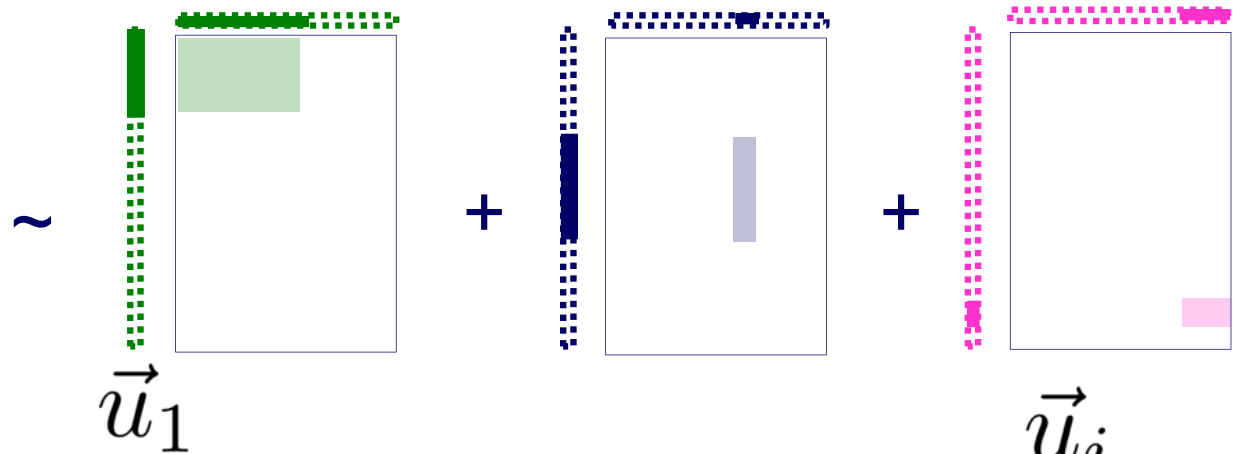


'music lovers'
'singers'

'sports lovers'
'athletes'

'citizens'
'politicians'

\vec{v}_1



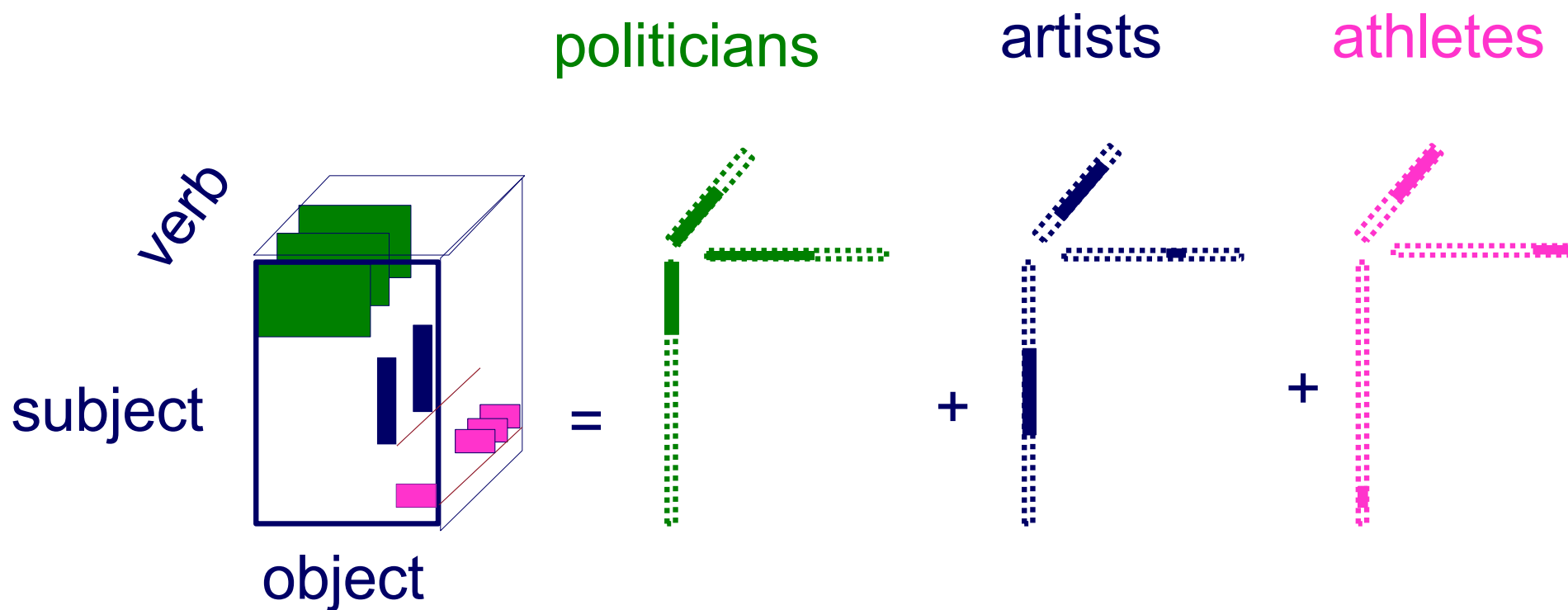
\vec{u}_1

Christos Faloutsos

\vec{u}_i 65

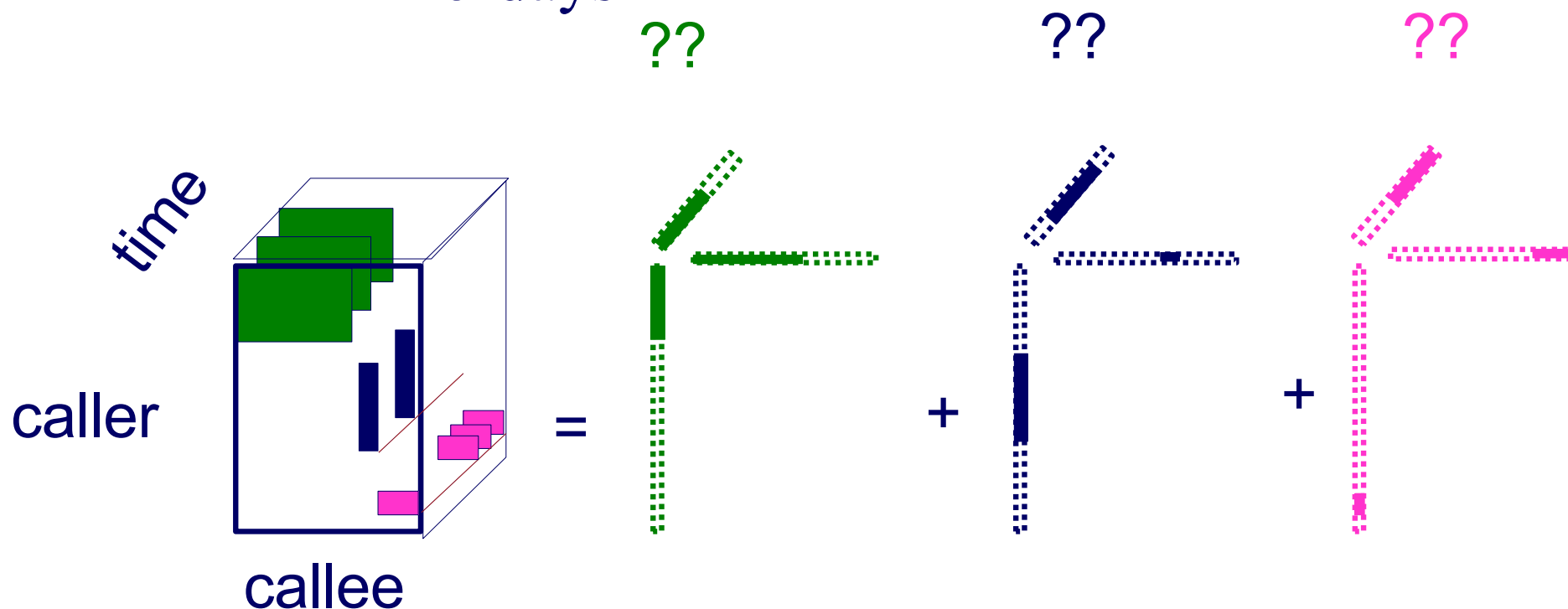
Answer: tensor factorization

- PARAFAC decomposition

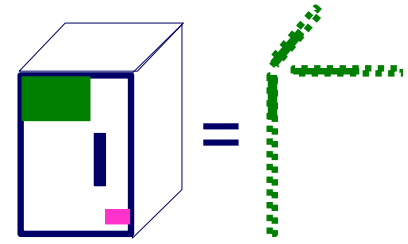


Answer: tensor factorization

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

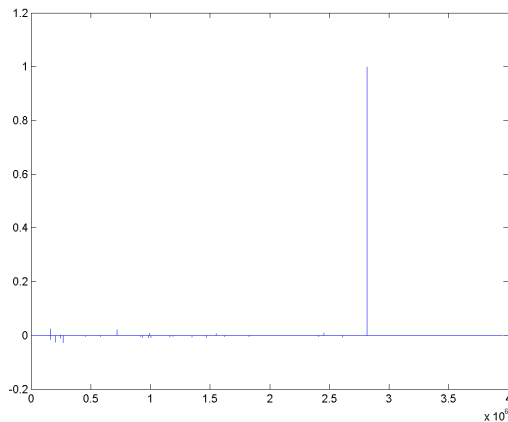


Anomaly detection in time-evolving graphs

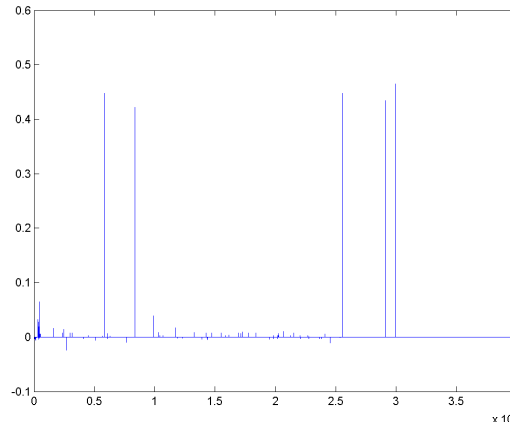


- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks

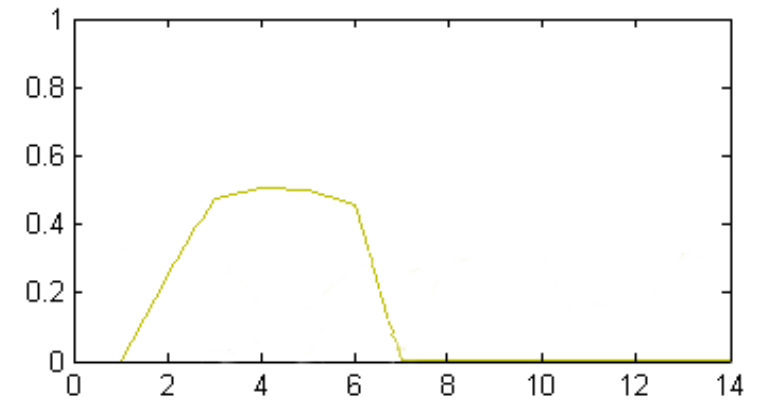
1 caller



5 receivers

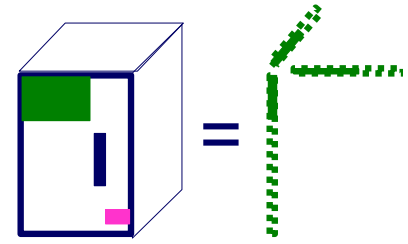


4 days of activity



~200 calls to EACH receiver on EACH day!

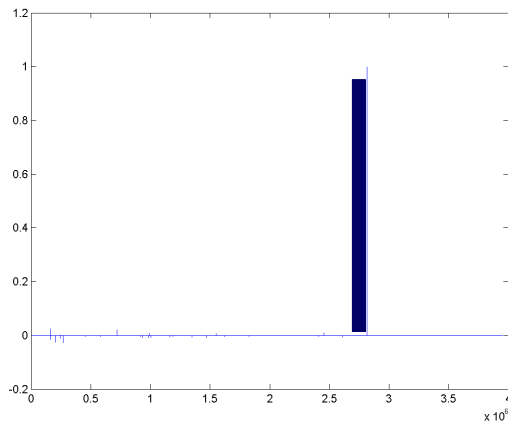
Anomaly detection in time-evolving graphs



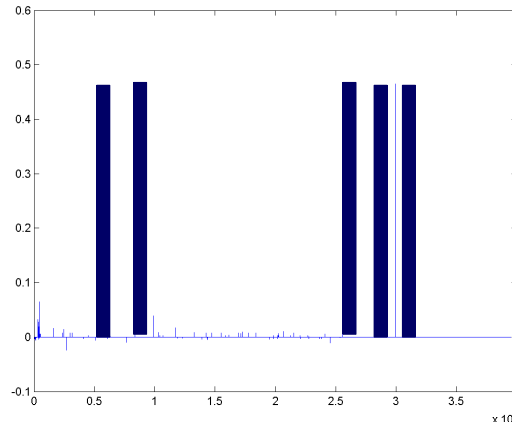
- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks



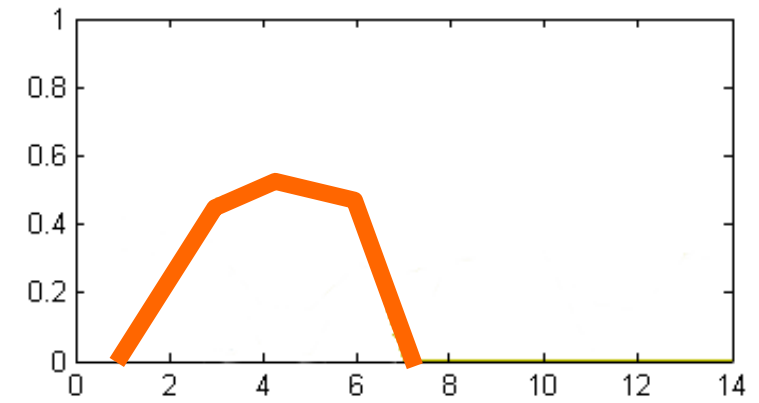
1 caller



5 receivers

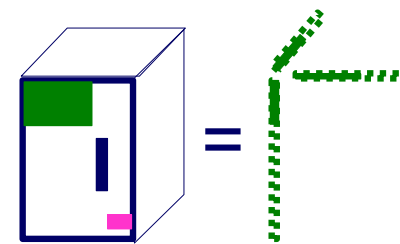


4 days of activity



~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.

Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs
 - P2.1: tools/tensors
 - ➔ – P2.2: other patterns – inter-arrival time
- Conclusions



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

Pattern Mining: Datasets

Reddit Dataset

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

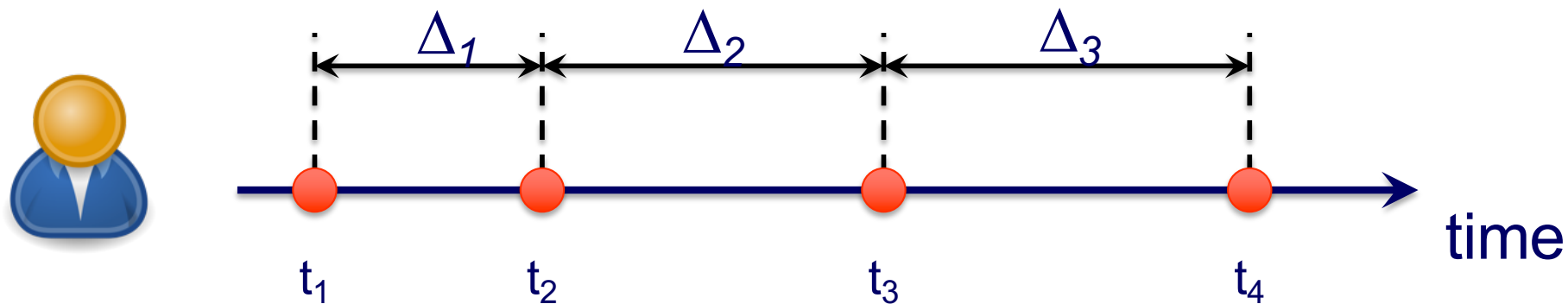
Twitter Dataset

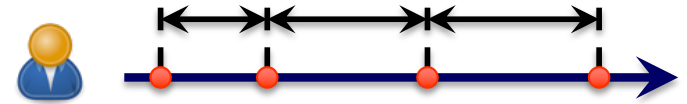
Time-stamp from tweets
6,790 users
16 Million time-stamps

For each user we have:

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

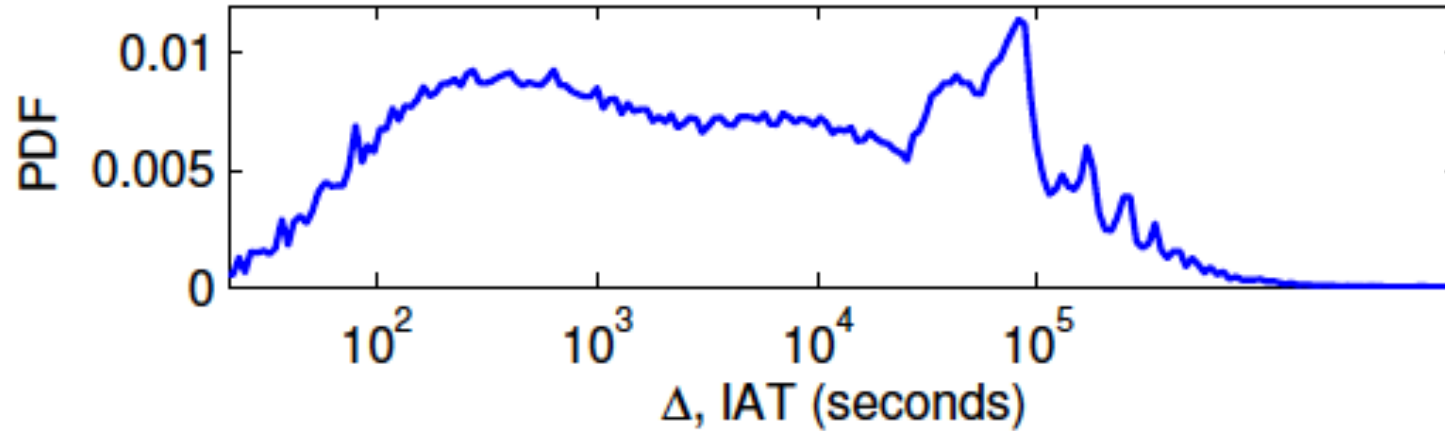
Inter-arrival times (IAT) of postings: $(\Delta_1, \Delta_2, \Delta_3, \dots)$



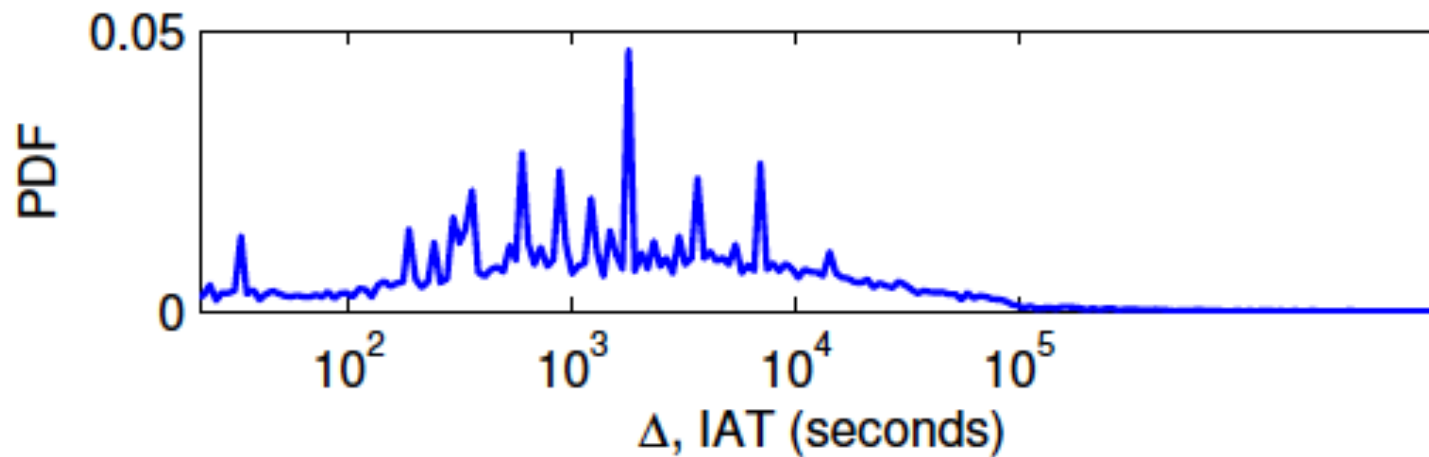


Human? Robots?

linear



log

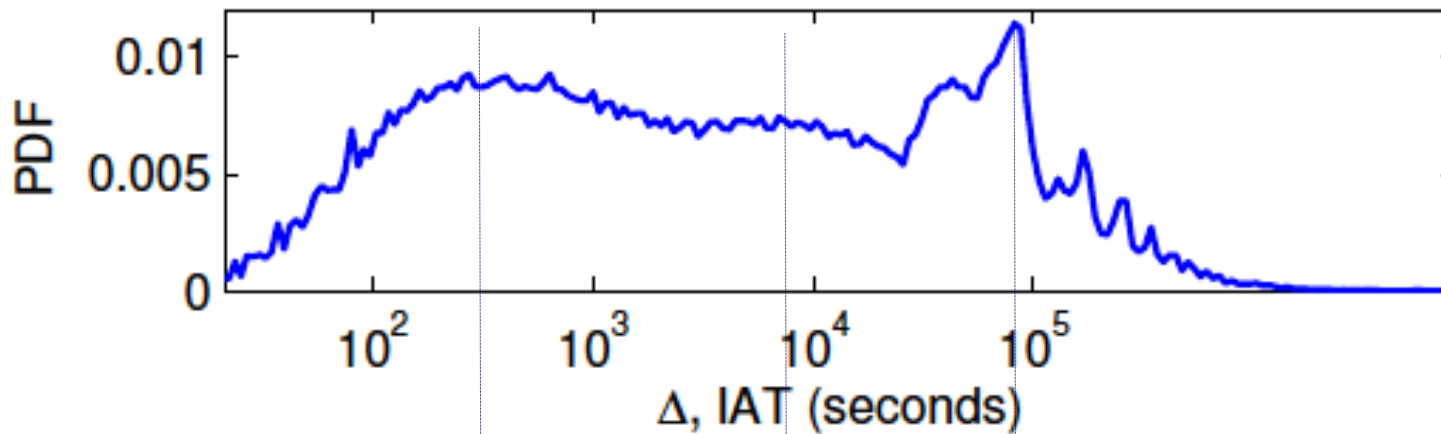




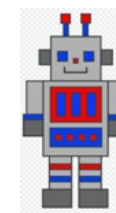
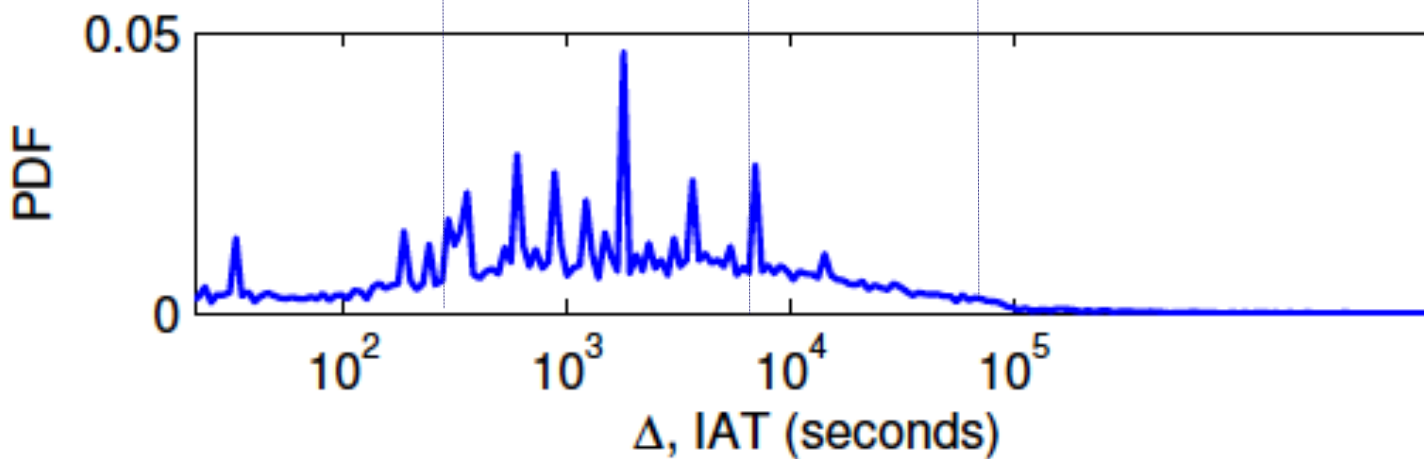
Human? Robots?

2' 3h 1day

linear



log

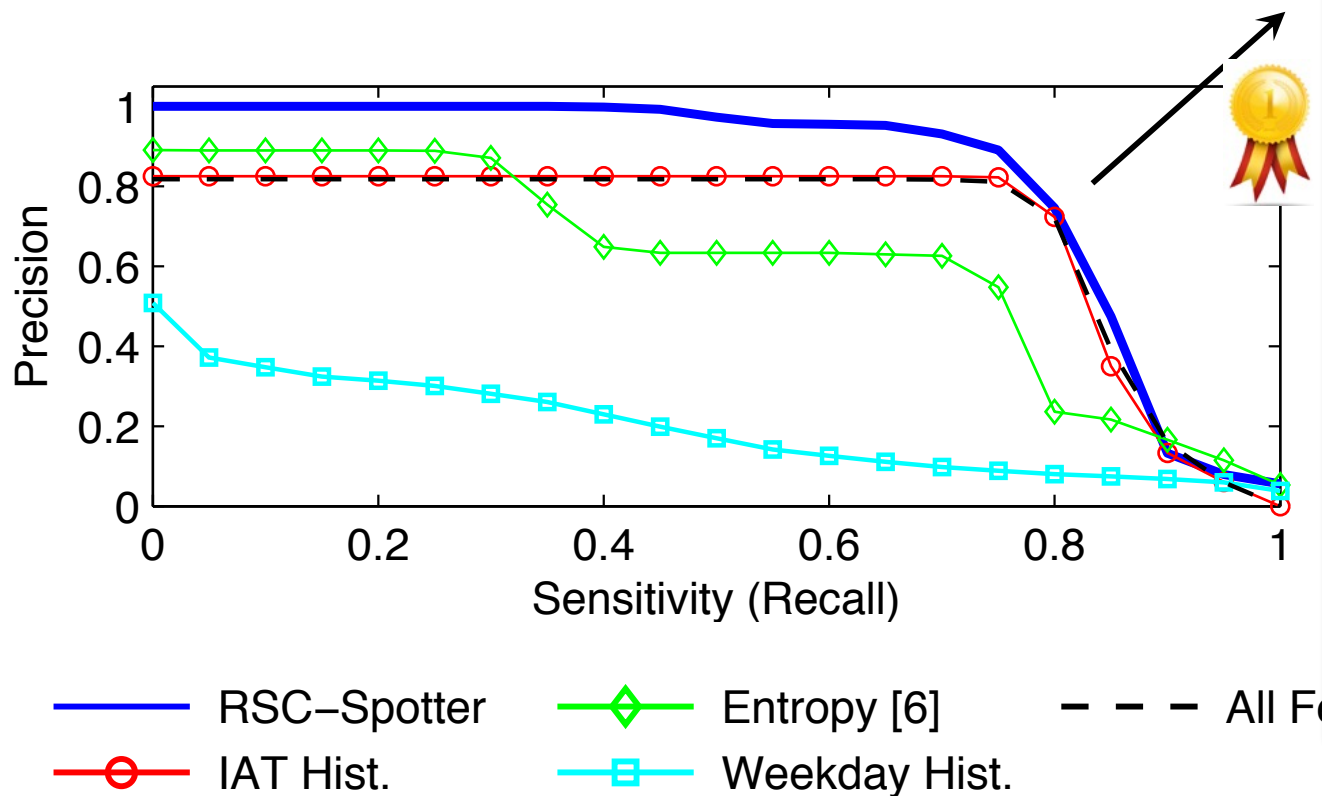


Experiments: Can RSC-Spotter Detect Bots?

Precision vs. Sensitivity Curves

Good performance: curve close to the top

Twitter



Precision > 94%
Sensitivity > 70%

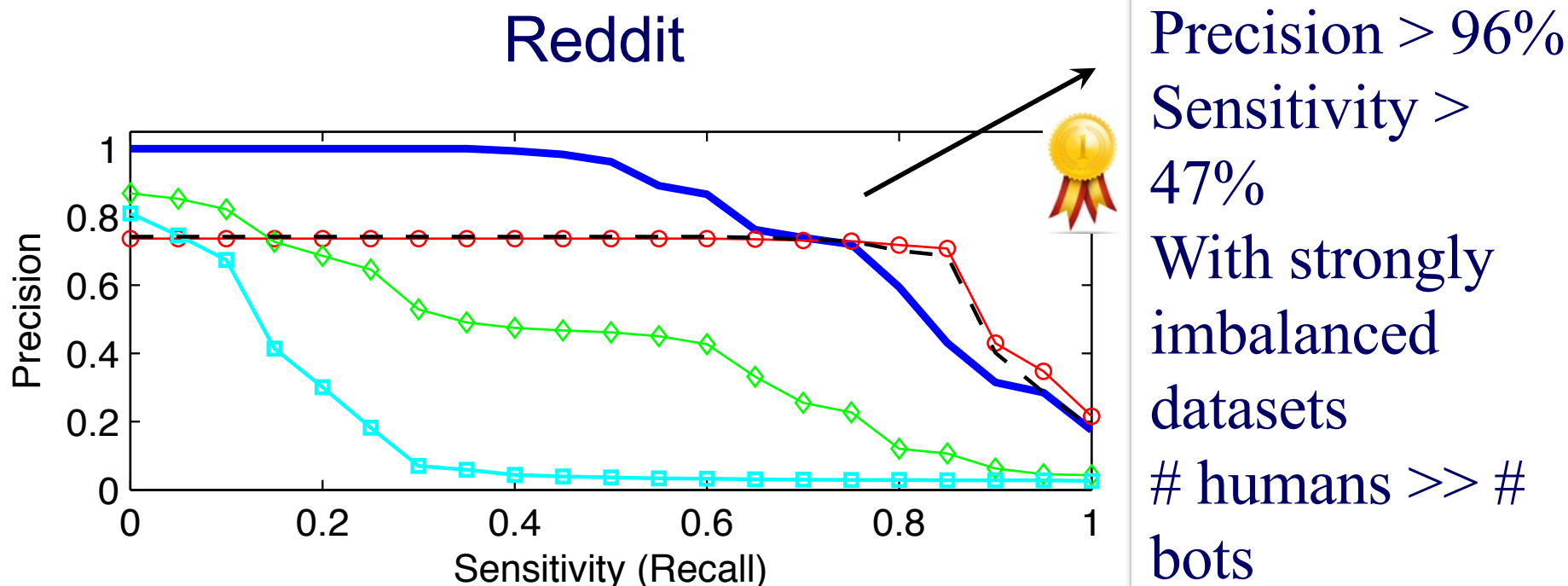
With strongly imbalanced datasets

humans >> # bots

Experiments: Can RSC-Spotter Detect Bots?

Precision vs. Sensitivity Curves

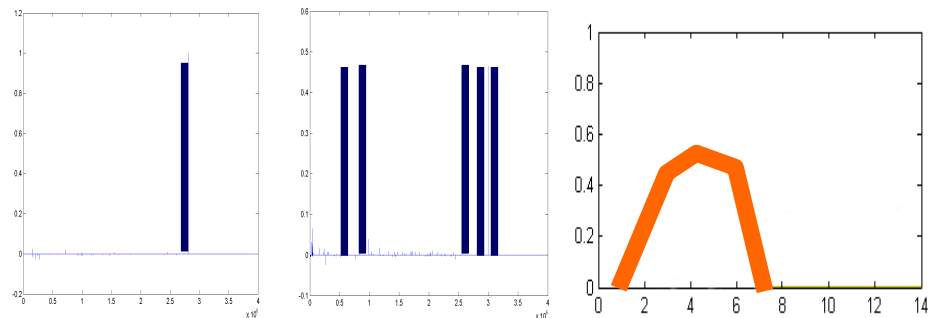
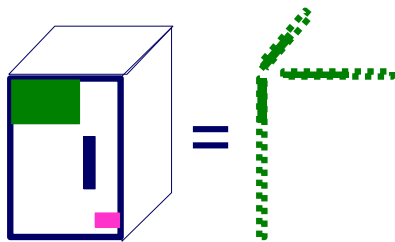
Good performance: curve close to the top



— RSC-Spotter —◇— Entropy [6] - - - All Features
—○— IAT Hist. —□— Weekday Hist.

Part 2: Conclusions

- Time-evolving / heterogeneous graphs \rightarrow tensors
- PARAFAC finds patterns
- Surprising temporal patterns



Roadmap

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



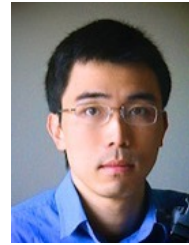
Akoglu,
Leman



Araujo,
Miguel



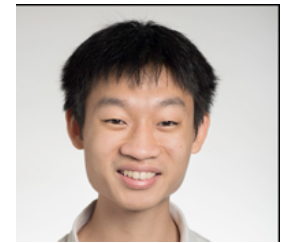
Beutel,
Alex



Chau,
Polo



Eswaran,
Dhivya



Hooi,
Bryan



Kang, U



Koutra,
Danai



Papalexakis,
Vagelis



Shah,
Neil




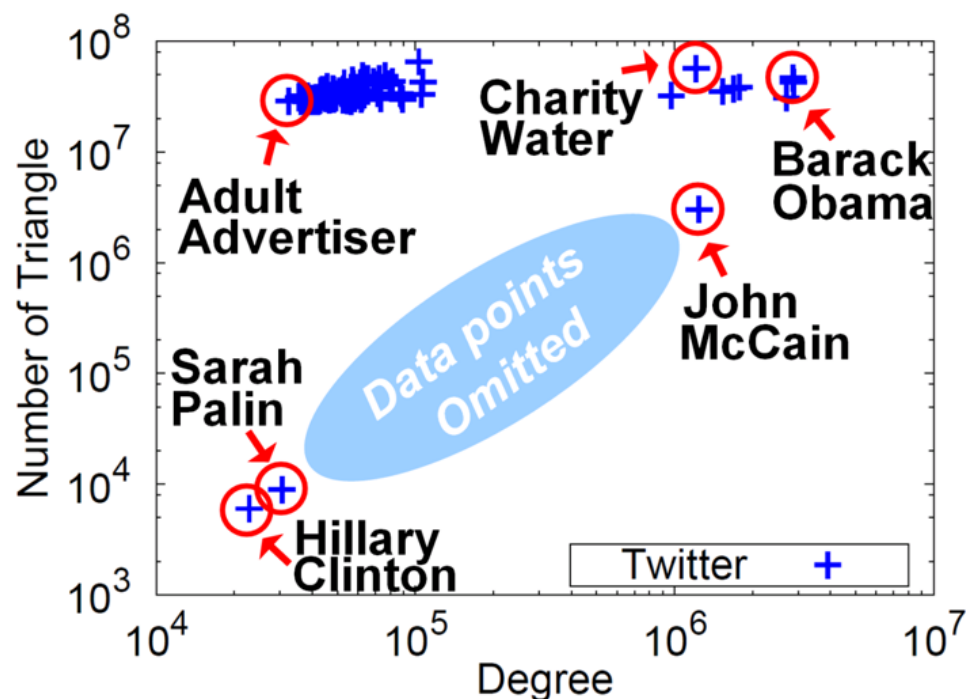
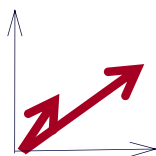
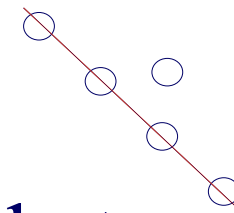
Shin,
Kijung



Song,
Hyun Ah

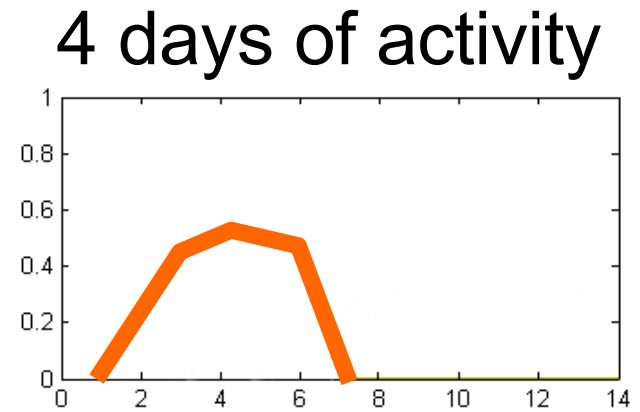
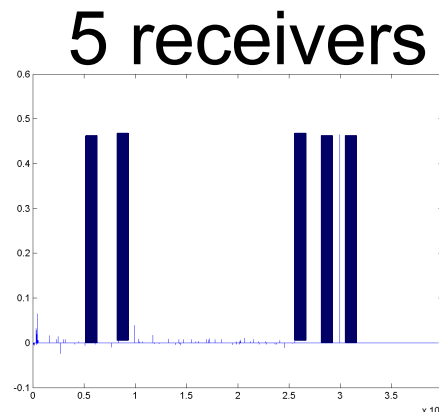
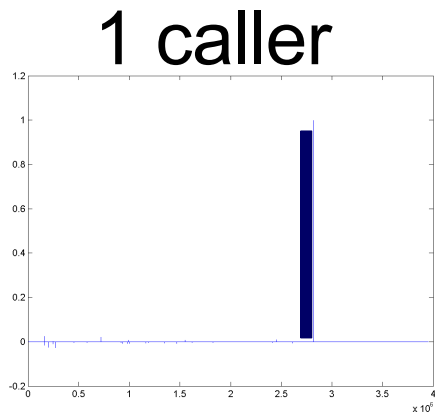
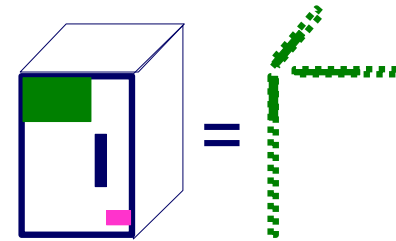
CONCLUSION#1 – Big data

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



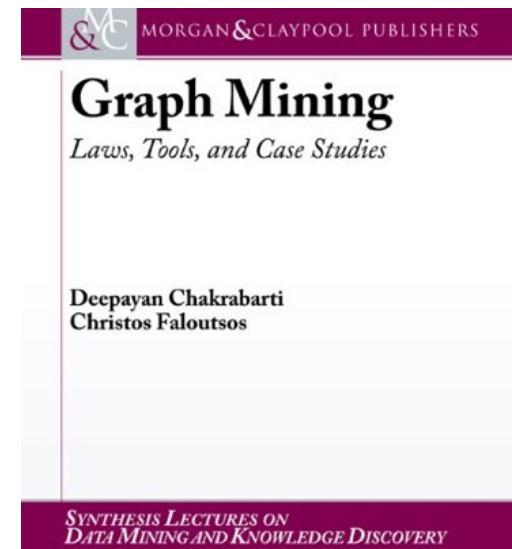
CONCLUSION#2 – tensors

- powerful tool



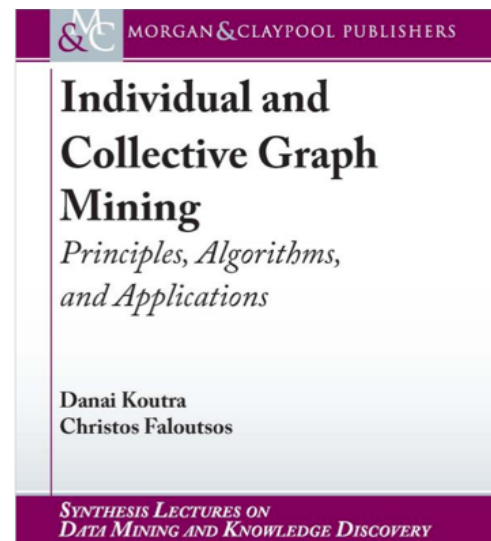
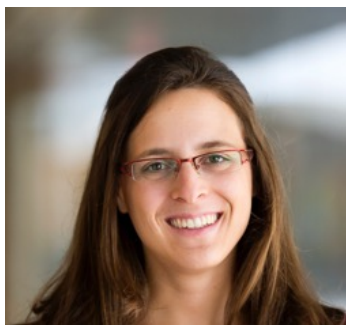
References

- D. Chakrabarti, C. Faloutsos: *Graph Mining – Laws, Tools and Case Studies*, Morgan Claypool 2012
- <http://www.morganclaypool.com/doi/abs/10.2200/S00449ED1V01Y201209DMK006>



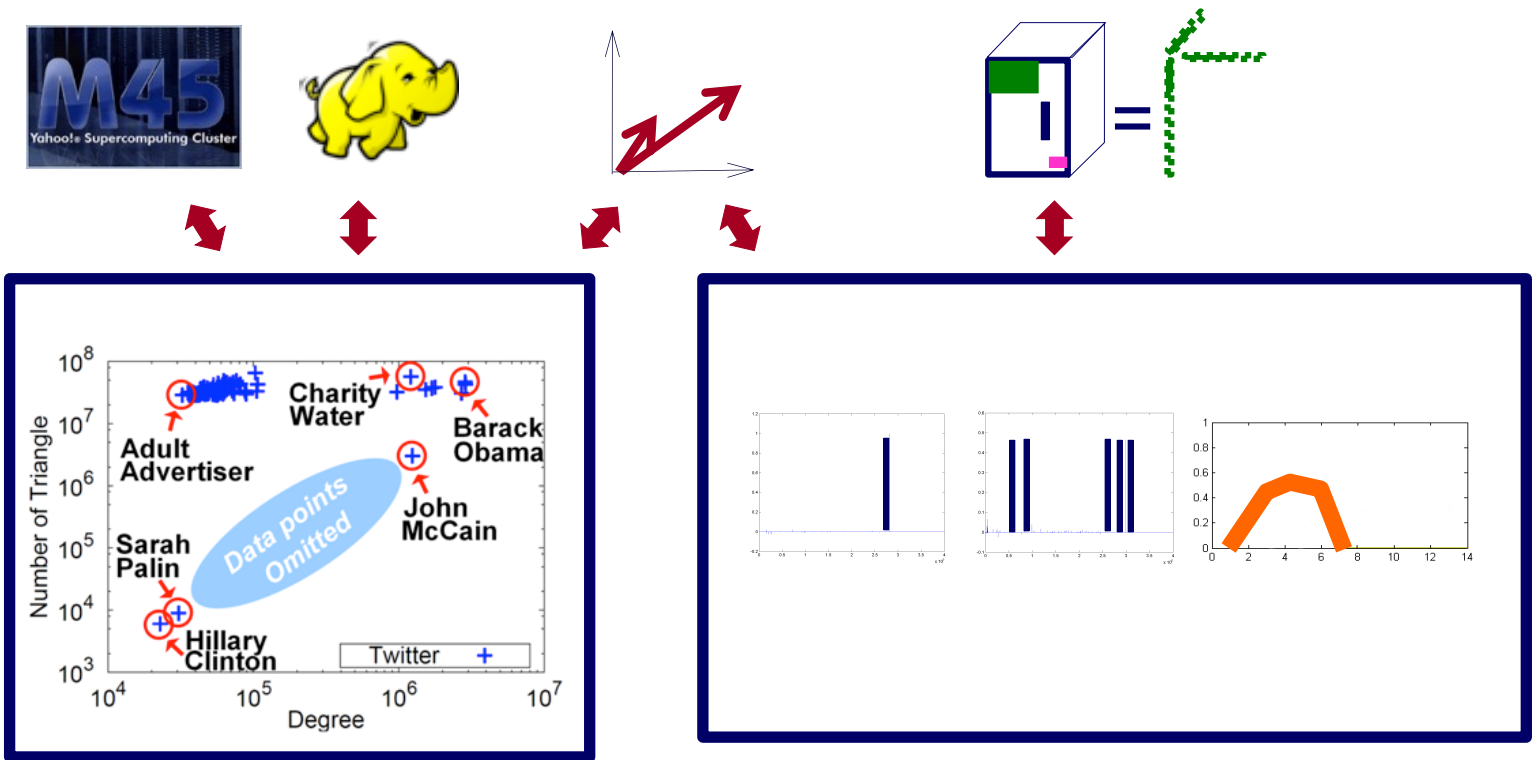
References

- Danai Koutra and Christos Faloutsos, *Individual and Collective Graph Mining: Principles, Algorithms, and Applications*, Morgan Claypool 2017
(<https://doi.org/10.2200/S00796ED1V01Y201708DMK014>)



TAKE HOME MESSAGE:

Cross-disciplinarity



Thank you!

Cross-disciplinarity

