

Anomaly Detection in Large Graphs

Christos Faloutsos

CMU

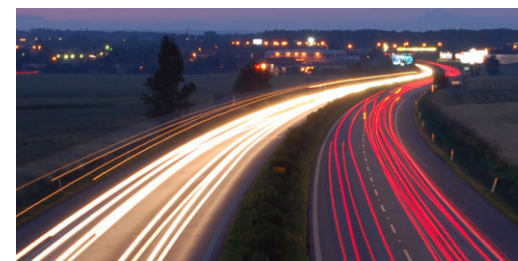
Thank you!

Prof. David Brumley

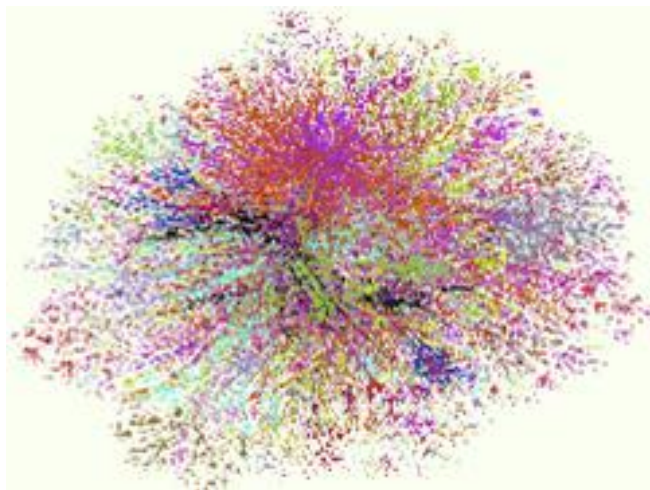
Megan Kearns

Roadmap

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



Graphs - why should we care?



Internet Map
[lumeta.com]

computer network security:

- Email traffic
- IP traffic (src, dst, dst-port, t)

Malware propagation

- (machine-id, infected-file-id)

Graphs - why should we care?

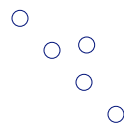


>\$10B; ~1B users



Motivating problems

- P1: patterns? Fraud detection?



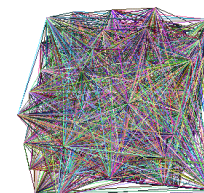
- P2: patterns in time-evolving graphs / tensors

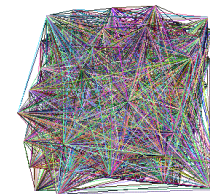
destination



source

time

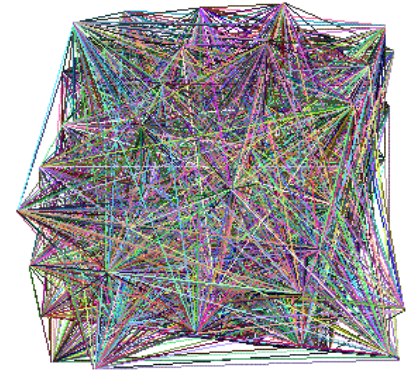




Part 1: Patterns, & fraud detection

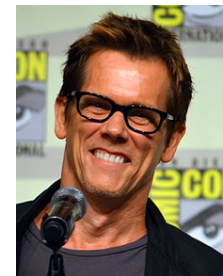
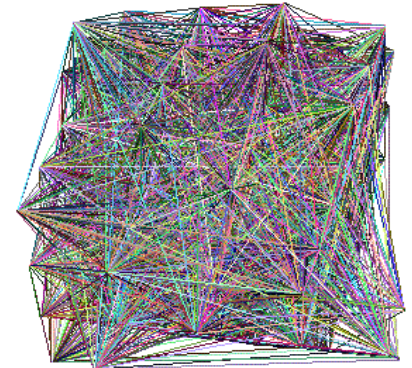
Laws and patterns

- Q1: Are real graphs random?

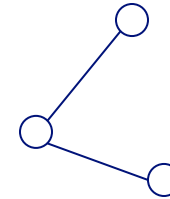


Laws and patterns

- 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

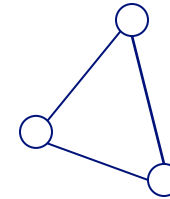


Solution# S.3: Triangle ‘Laws’

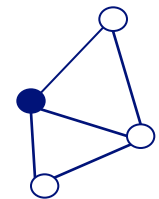


- Real social networks have a lot of triangles

Solution# S.3: Triangle ‘Laws’



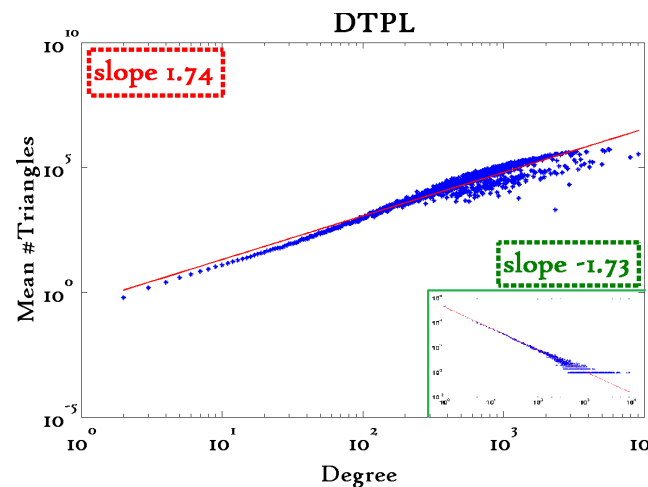
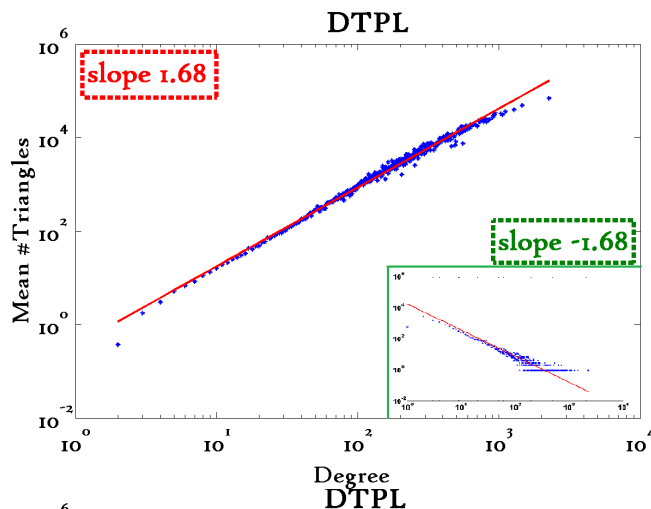
- Real social networks have a lot of triangles
 - Friends of friends are friends
- Any patterns?
 - 2x the friends, 2x the triangles ?



Triangle Law: #S.3

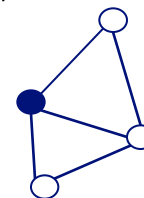
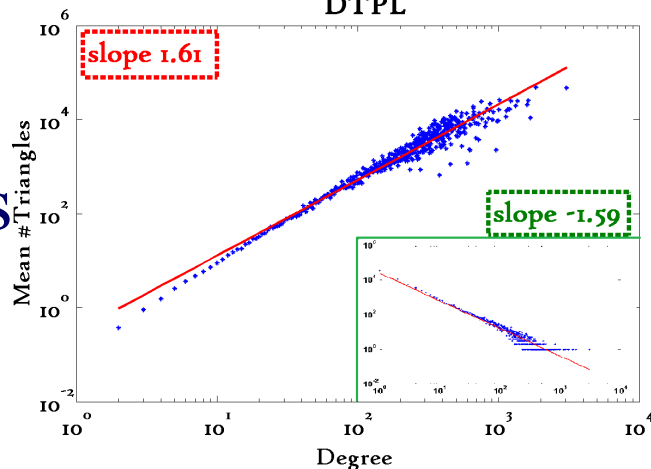
[Tsourakakis ICDM 2008]

Reuters



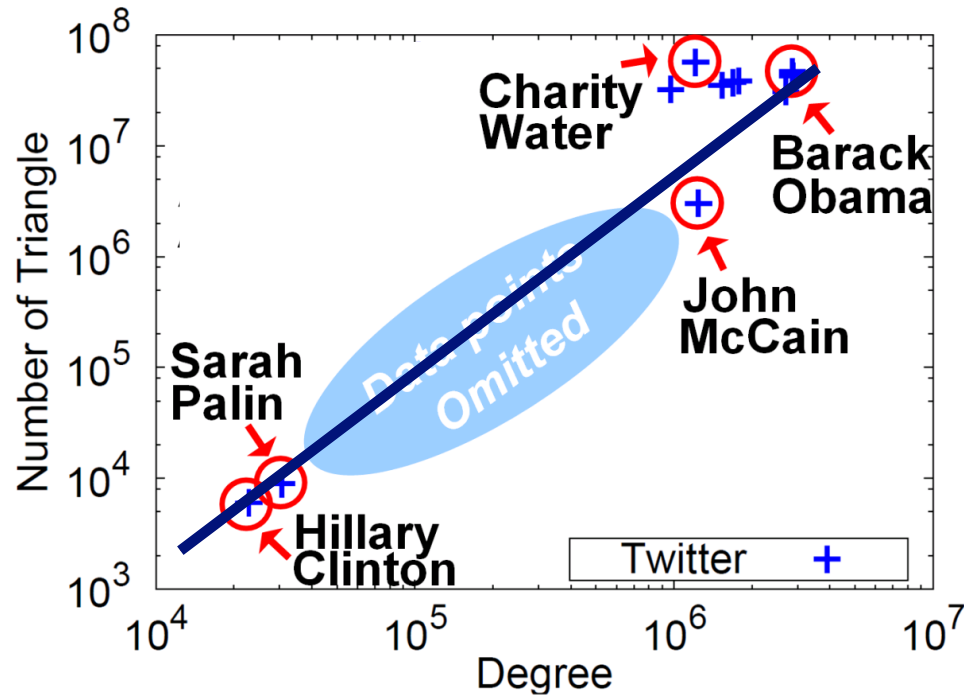
SN

Epinions



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

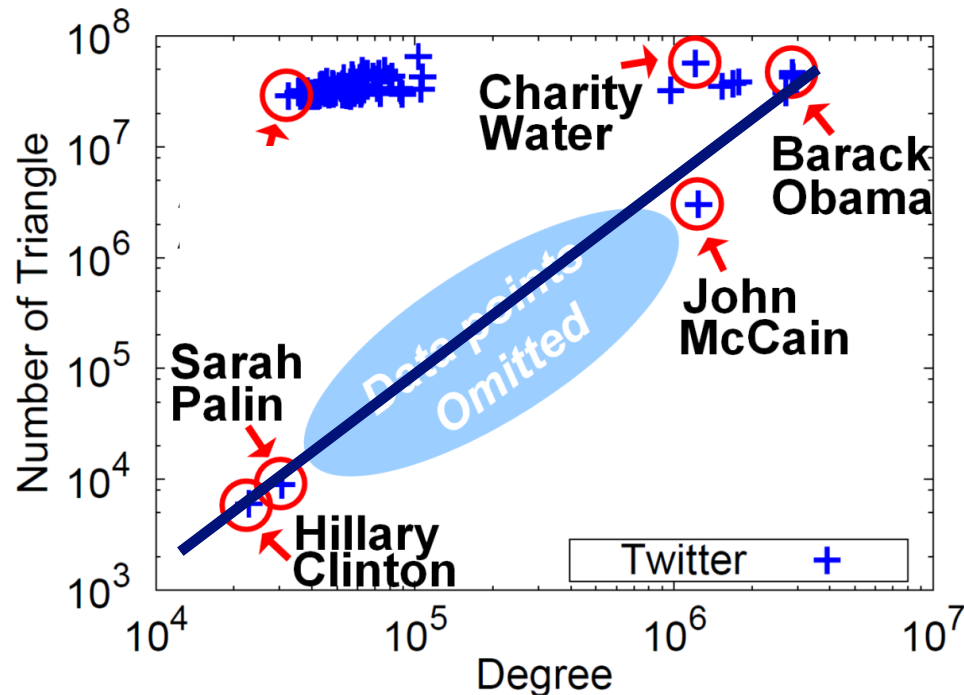
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

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

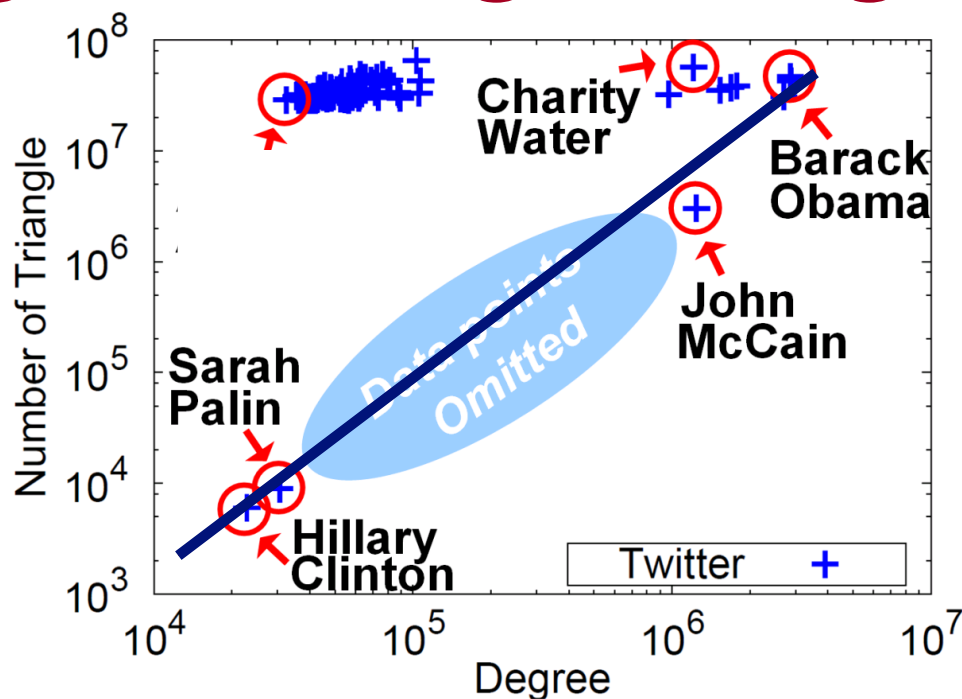
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

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

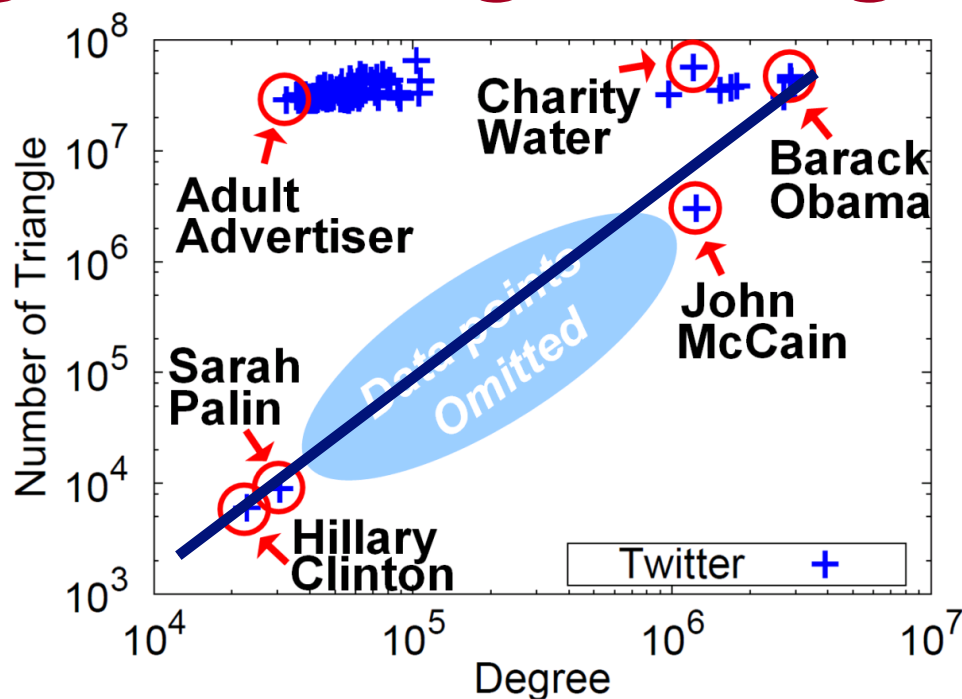
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

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

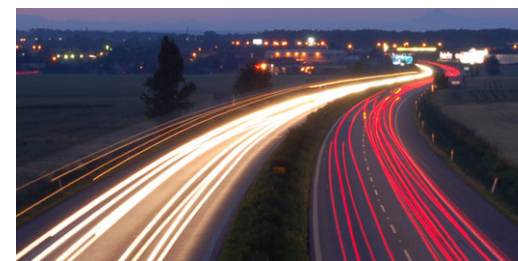
Triangle counting for large graphs?



Anomalous nodes in Twitter (~ 3 billion edges)

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

Roadmap

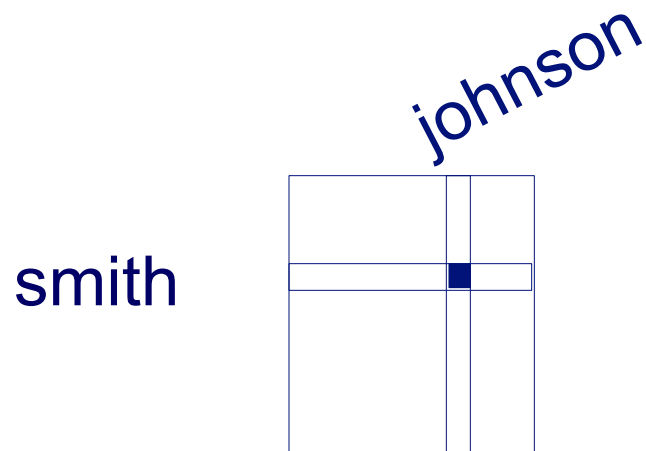


- Introduction – Motivation
- Part#1: Patterns in graphs
- ➔ • Part#2: time-evolving graphs; tensors
- Conclusions

Part 2: Time evolving graphs; tensors

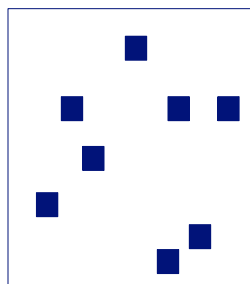
Graphs over time -> tensors!

- Problem #2:
 - Given who calls whom, and when
 - Find patterns / anomalies



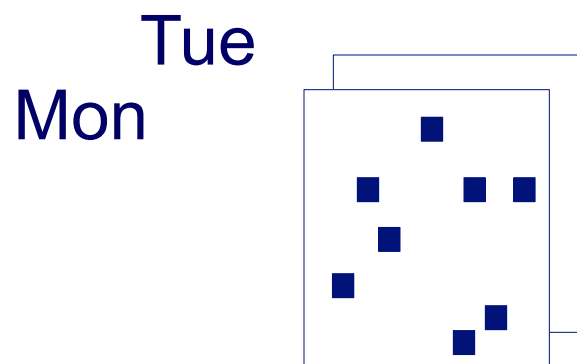
Graphs over time -> tensors!

- Problem #2:
 - Given who calls whom, and when
 - Find patterns / anomalies



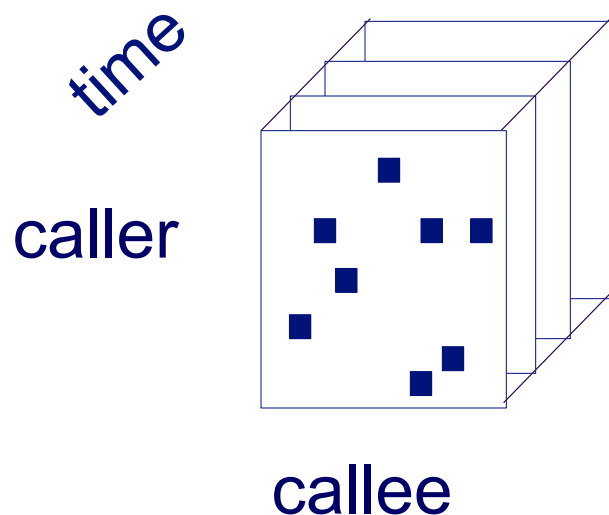
Graphs over time -> tensors!

- Problem #2:
 - Given who calls whom, and when
 - Find patterns / anomalies



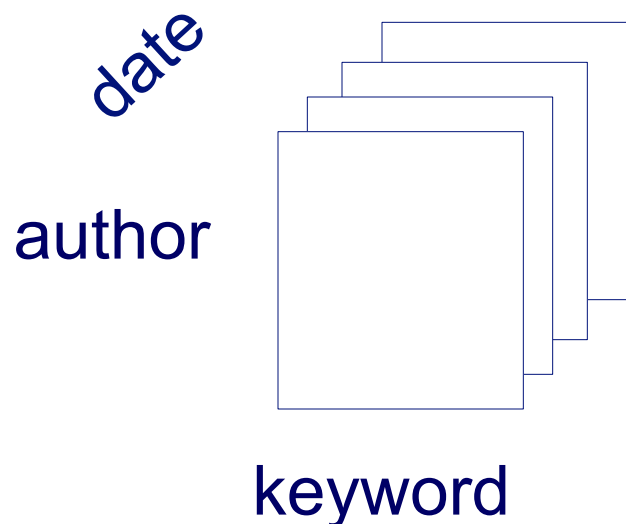
Graphs over time -> tensors!

- Problem #2:
 - Given who calls whom, and when
 - Find patterns / anomalies



Graphs over time -> tensors!

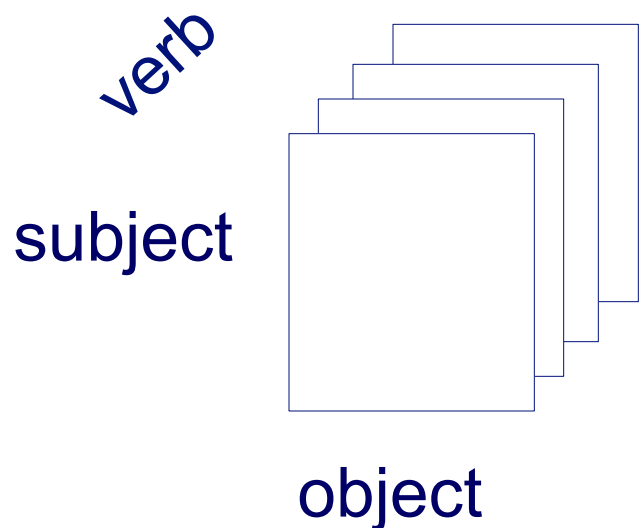
- Problem #2':
 - Given author-keyword-date
 - Find patterns / anomalies



MANY more settings,
with >2 'modes'

Graphs over time -> tensors!

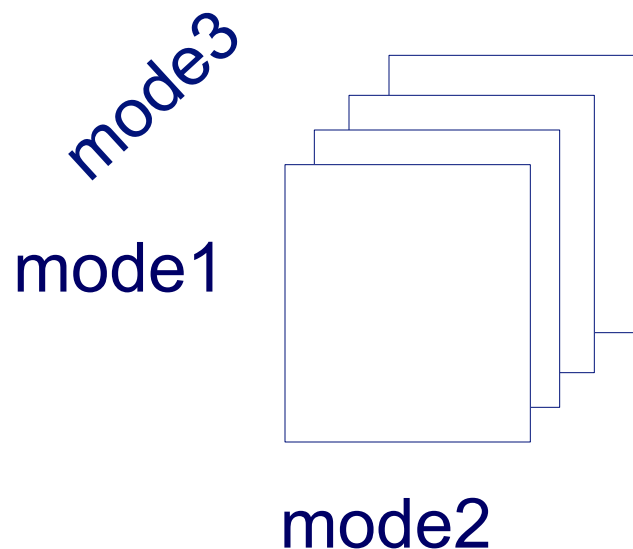
- Problem #2’’:
 - Given subject – verb – object facts
 - Find patterns / anomalies



MANY more settings,
with >2 ‘modes’

Graphs over time -> tensors!

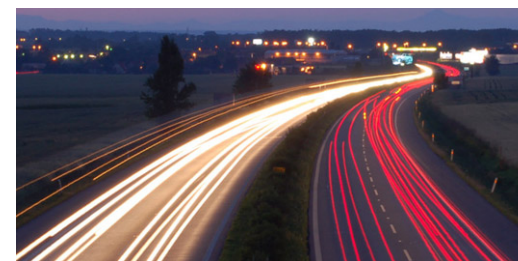
- Problem #2''':
 - Given <triplets>
 - Find patterns / anomalies



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

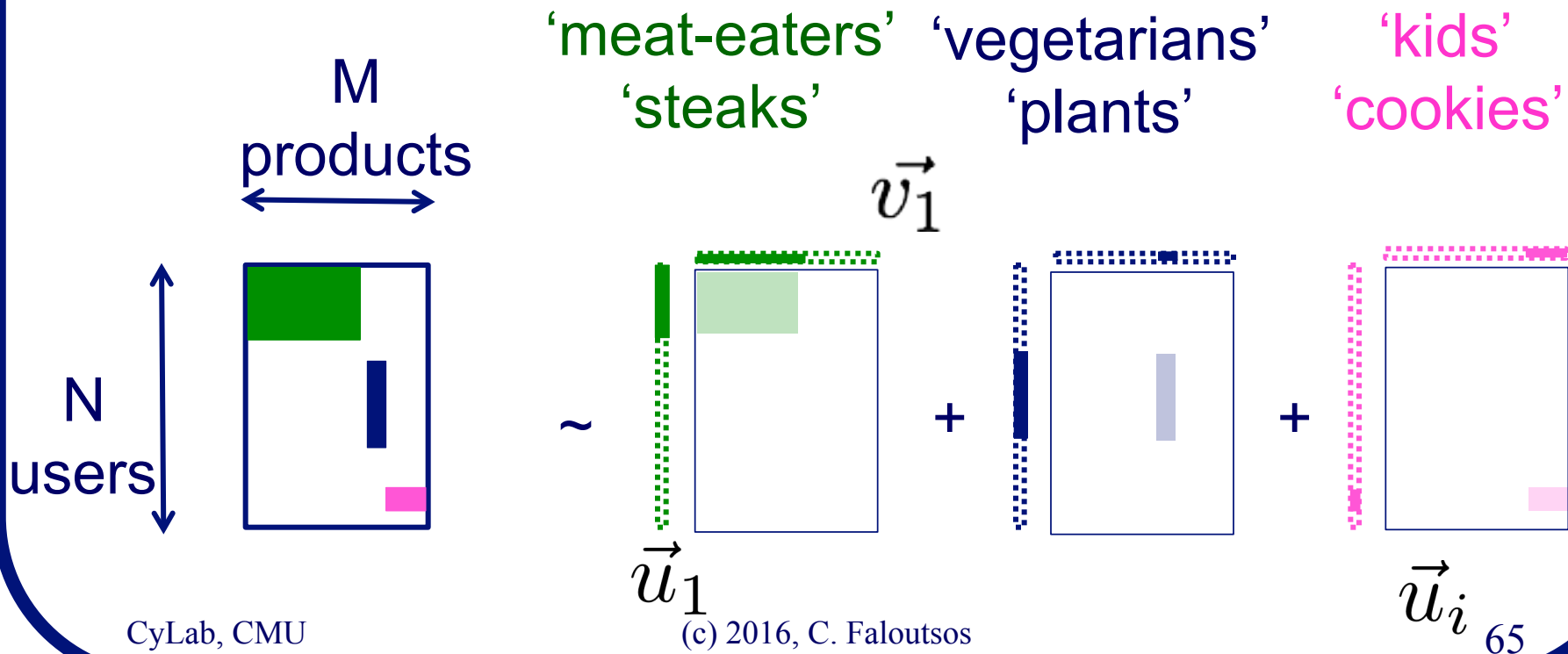
Roadmap

- Introduction – Motivation
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
 - ➔ – Intro to tensors
 - Results
 - Speed
- Conclusions



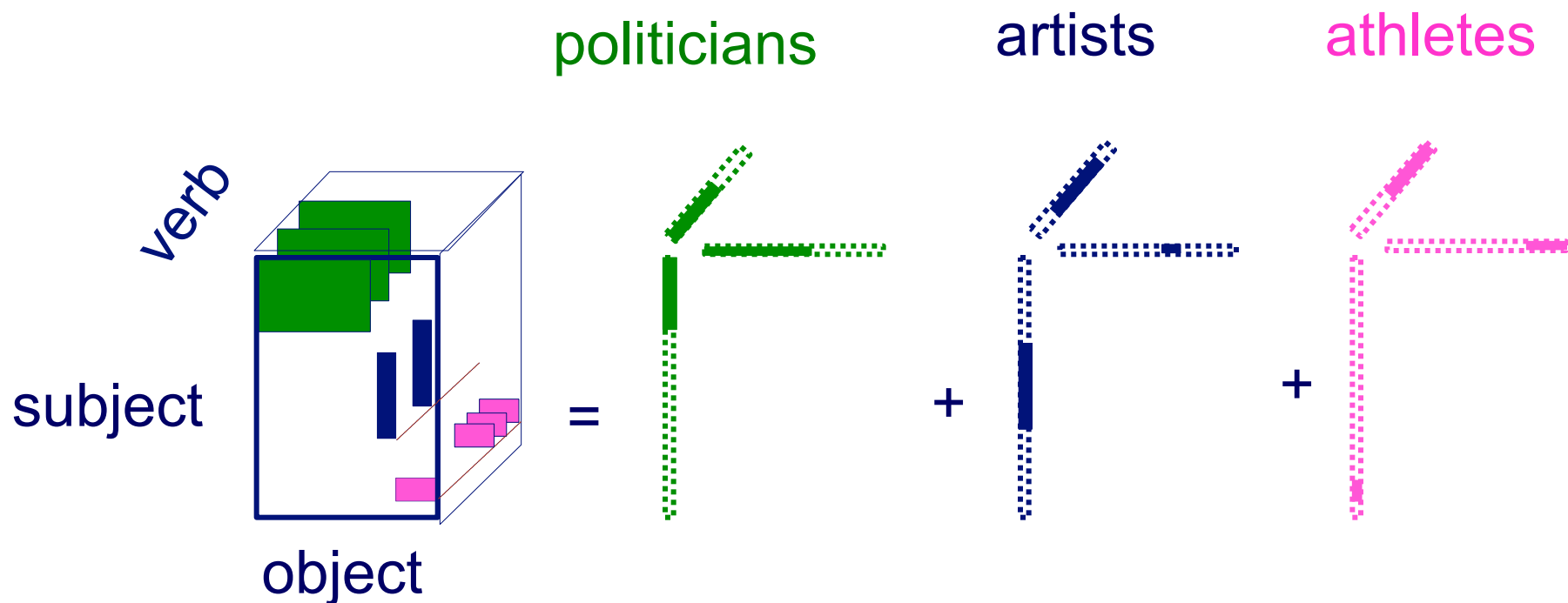
Answer to both: tensor factorization

- Recall: (SVD) matrix factorization: finds blocks



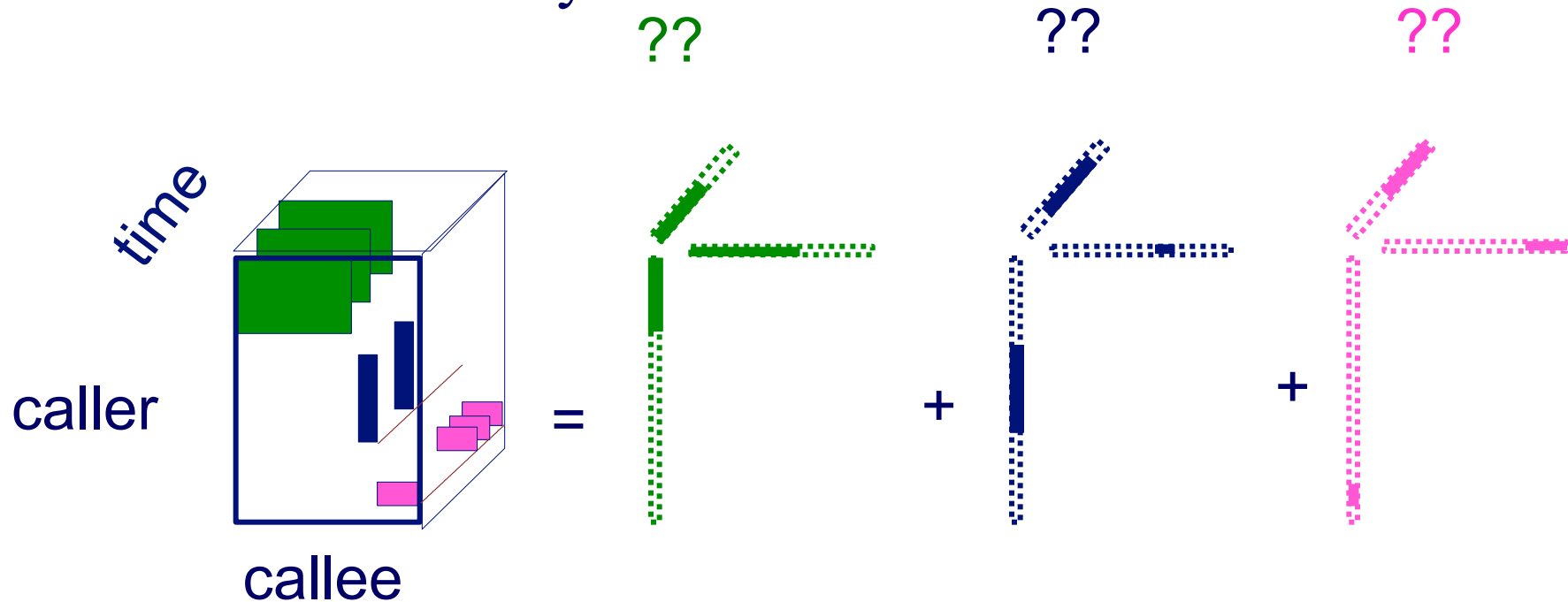
Answer to both: tensor factorization

- PARAFAC decomposition

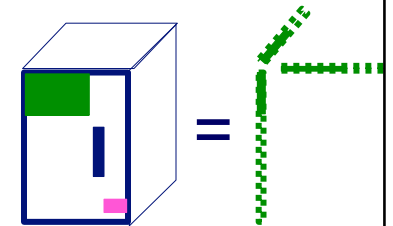


Answer: tensor factorization

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

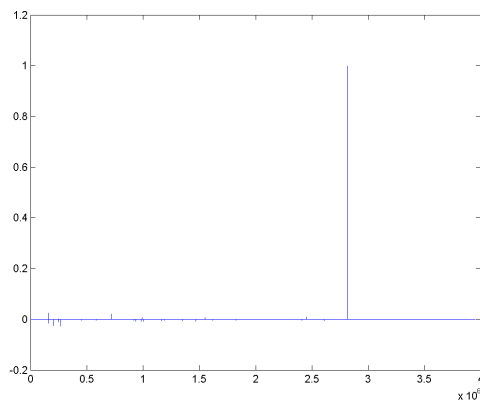


Anomaly detection in time-evolving graphs

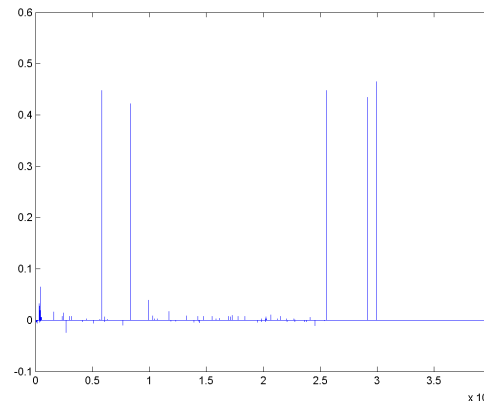


- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks

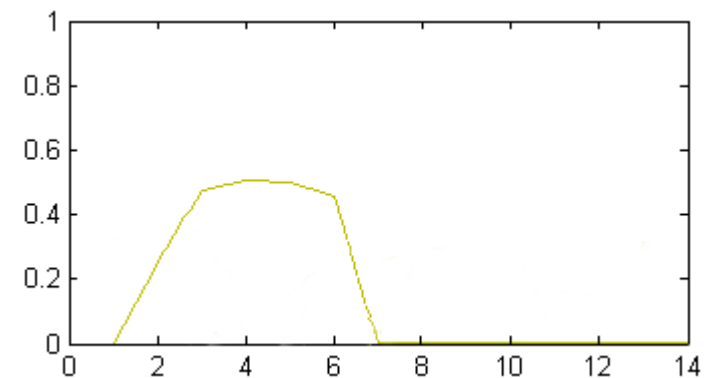
1 caller



5 receivers

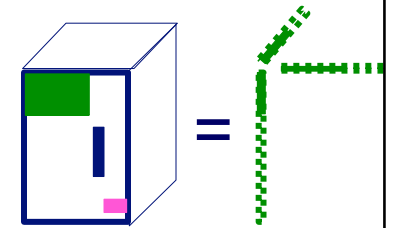


4 days of activity



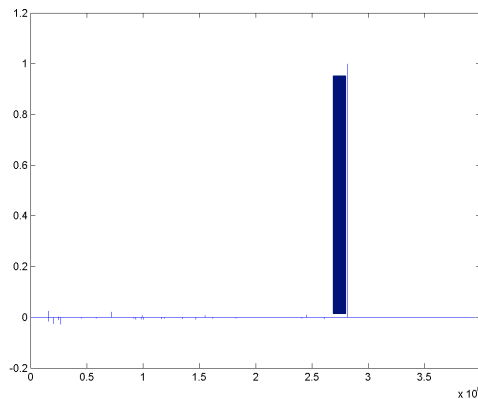
~200 calls to EACH receiver on EACH day!

Anomaly detection in time-evolving graphs

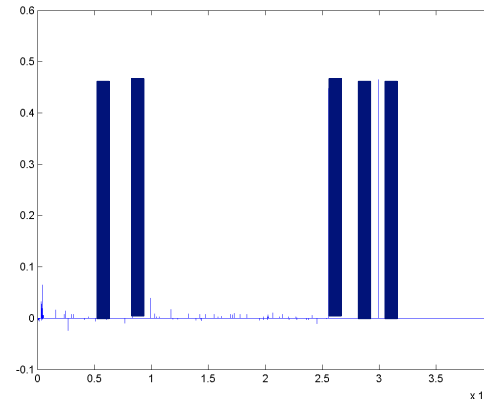


- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks

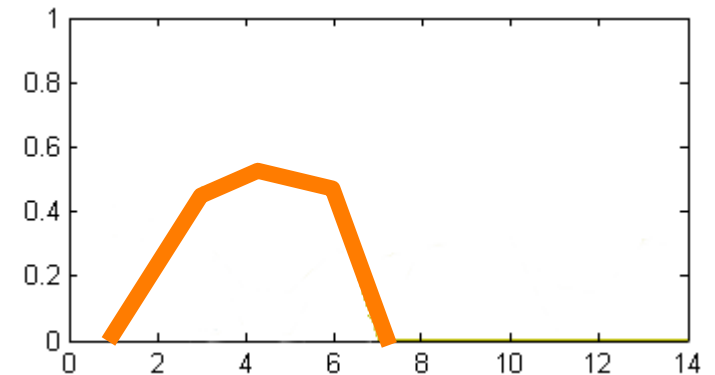
1 caller



5 receivers

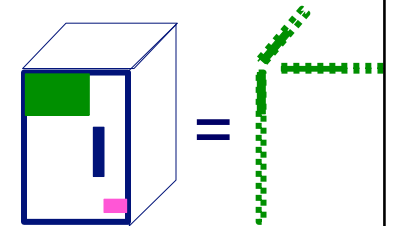


4 days of activity

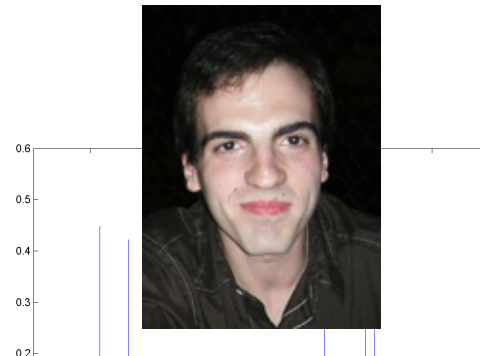
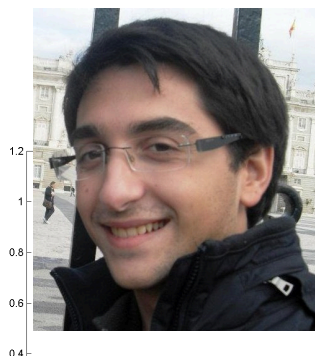


~200 calls to EACH receiver on EACH day!

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.

Roadmap



- Introduction – Motivation
 - Why study (big) graphs?
- Part#1: Patterns in graphs
- Part#2: time-evolving graphs; tensors
 - – Inter-arrival time patterns
- Acknowledgements and Conclusions

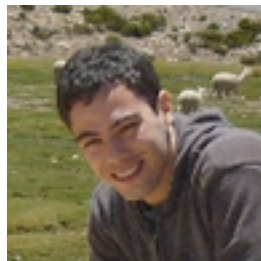
Carnegie Mellon



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

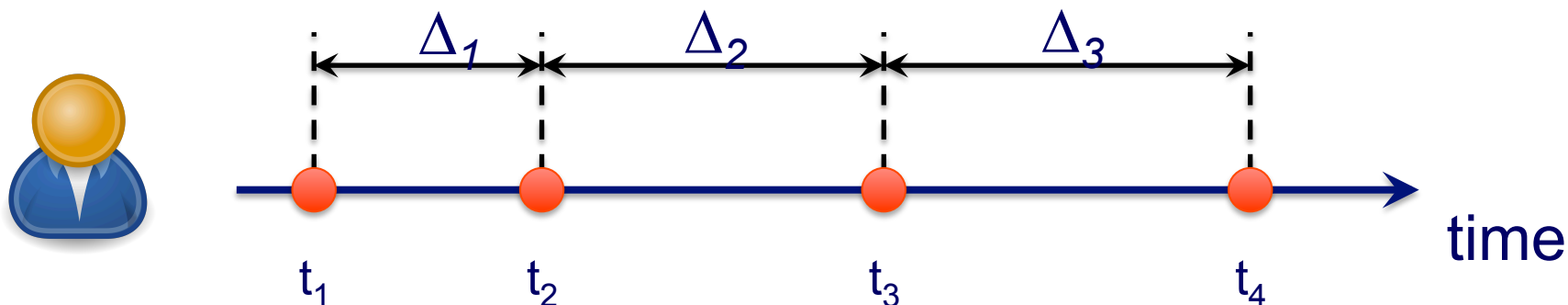
Twitter Dataset

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: $(\Delta_1, \Delta_2, \Delta_3, \dots)$

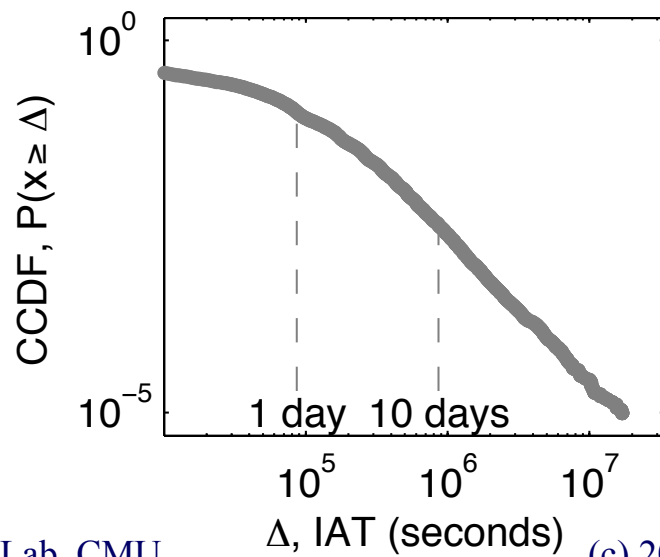


Pattern Mining

Pattern 1: Distribution of IAT is heavy-tailed

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

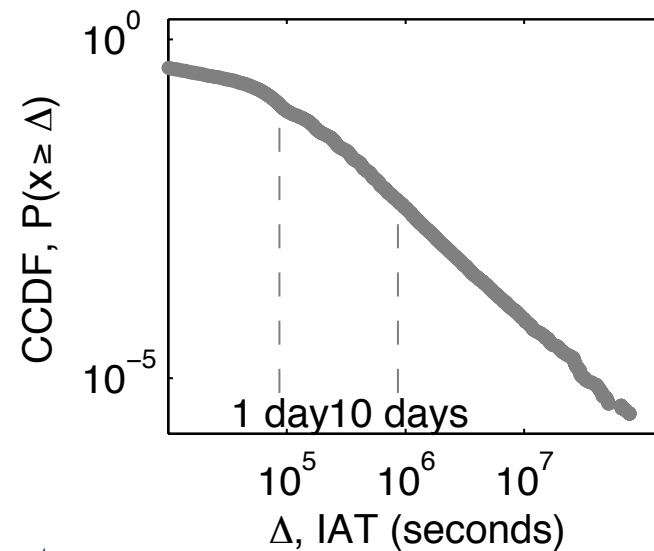
IAT Complementary Cumulative Distribution Function (CCDF)
(log-log axis)



CyLab, CMU

Reddit Users

(c) 2016, C. Faloutsos



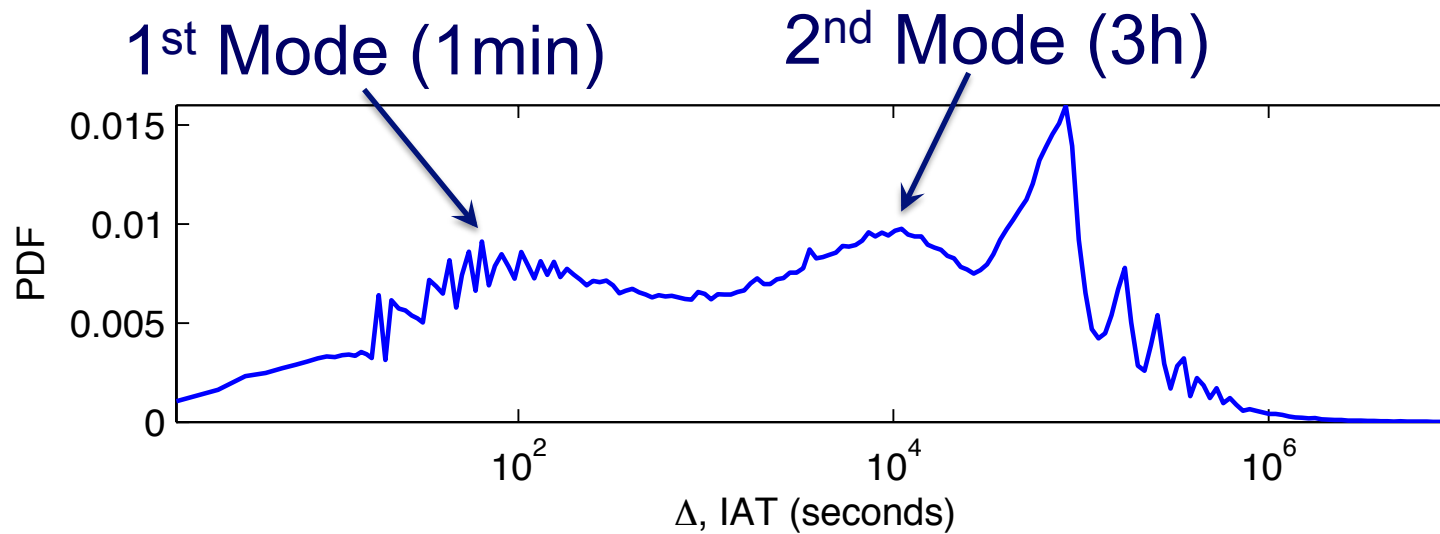
Twitter Users

Pattern Mining

Pattern 2: Bimodal IAT distribution

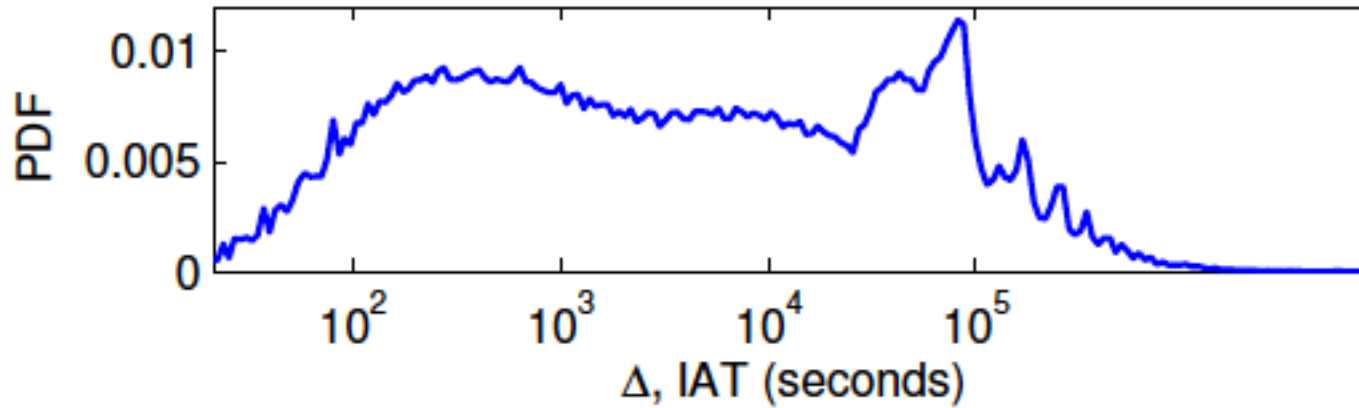
Users have highly active sections and resting periods

Log-binned histogram of postings IAT

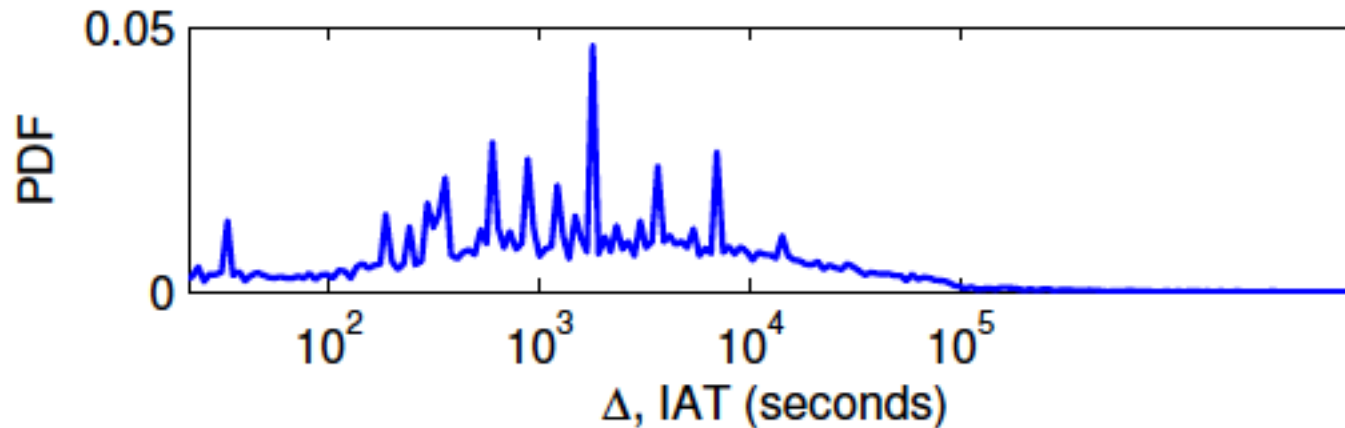


Human? Robots?

linear



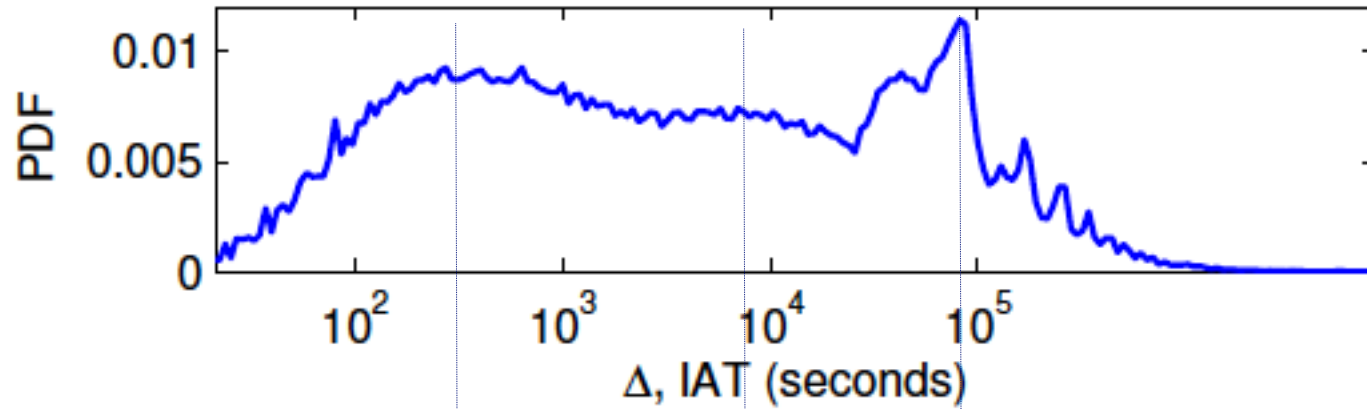
log



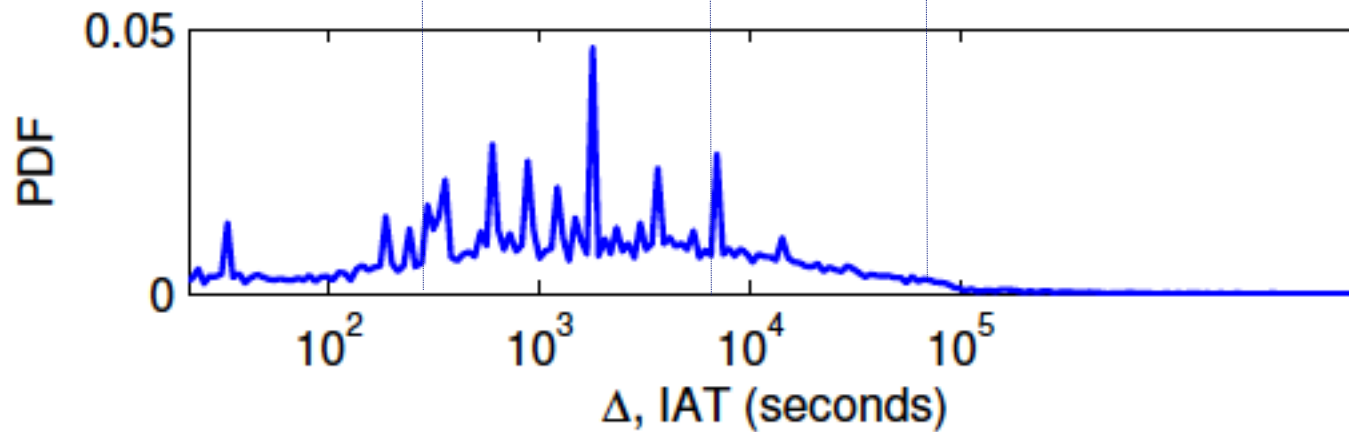
Human? Robots?

2' 3h 1day

linear



log

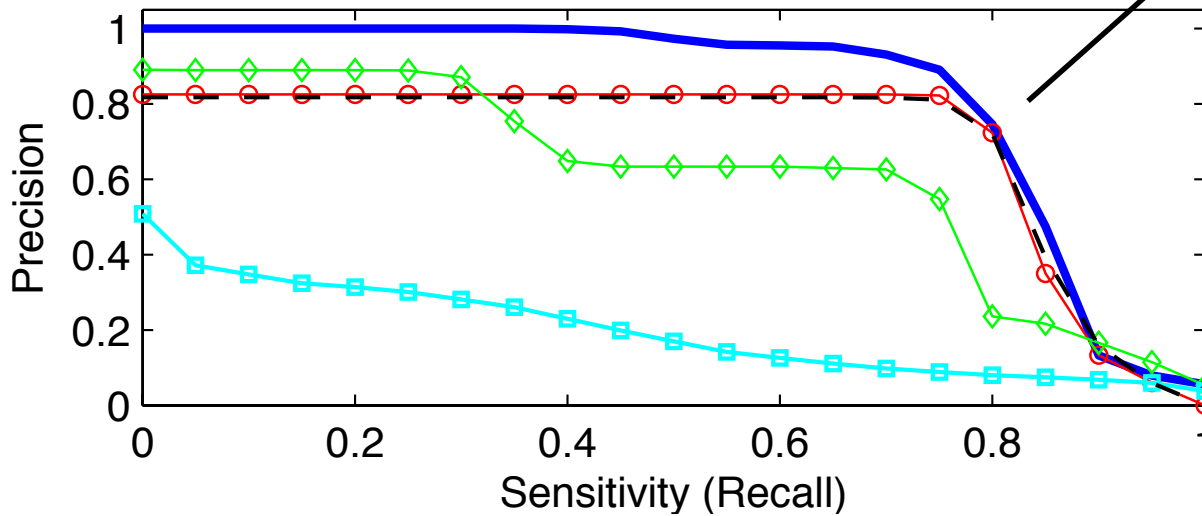


Experiments: Can RSC-Spotter Detect Bots?

Precision vs. Sensitivity Curves

Good performance: curve close to the top

Twitter



Precision > 94%
Sensitivity > 70%

With strongly imbalanced datasets

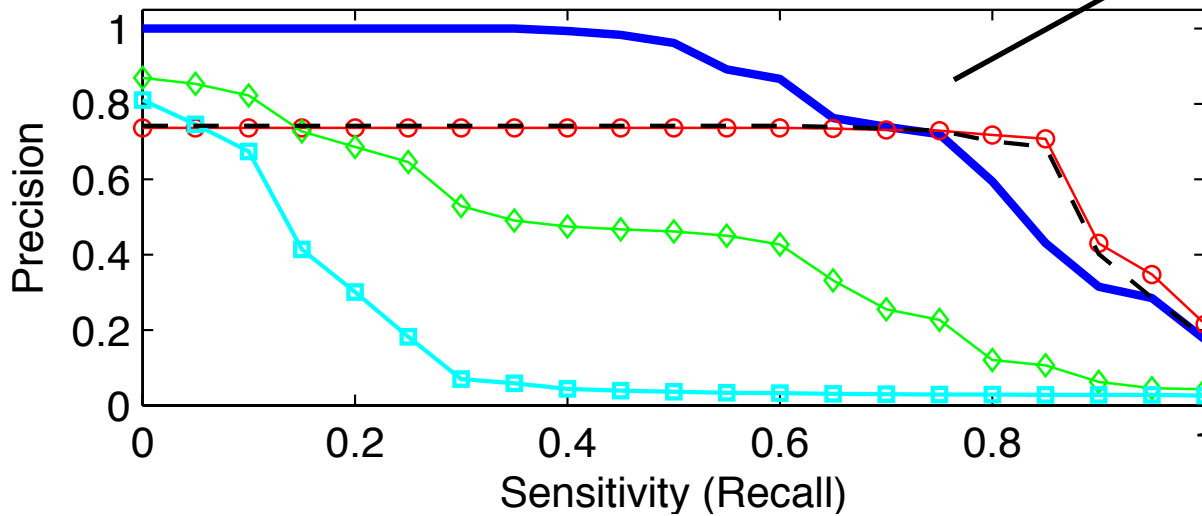
- RSC-Spotter
 - ◇— Entropy [6]
 - ◇— All Features
 - IAT Hist.
 - Weekday Hist.
- # humans >> # bots

Experiments: Can RSC-Spotter Detect Bots?

Precision vs. Sensitivity Curves

Good performance: curve close to the top

Reddit



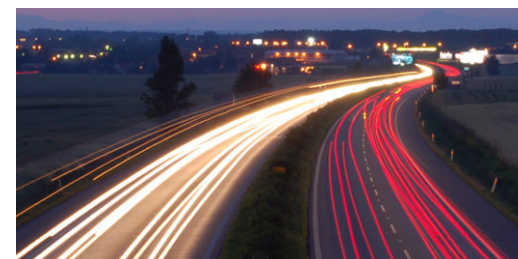
Precision > 96%
Sensitivity > 47%

With strongly imbalanced datasets

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

humans >> # bots

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



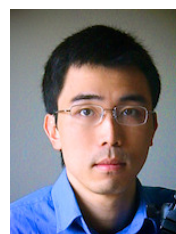
Akoglu,
Leman



Araujo,
Miguel



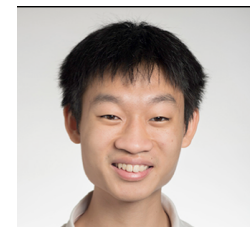
Beutel,
Alex



Chau,
Polo



Eswaran,
Dhivya



Hooi,
Bryan



Kang, U



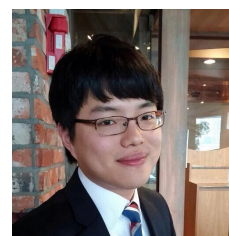
Koutra,
Danai



Papalexakis,
Vagelis



Shah,
Neil




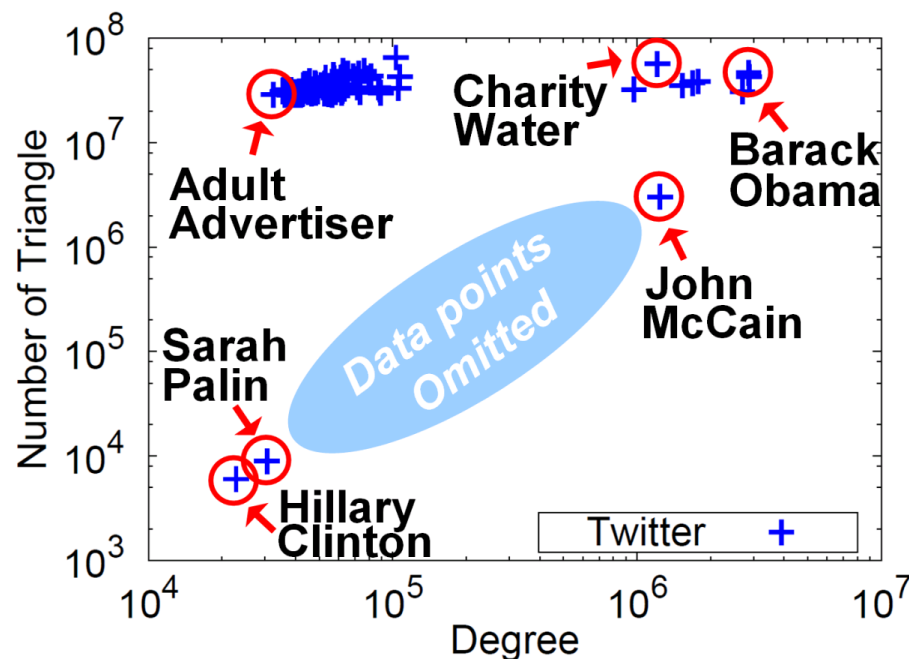
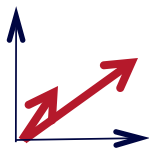
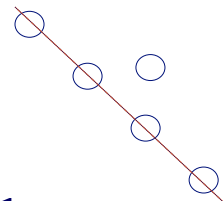
Shin,
Kijung



Song,
Hyun Ah

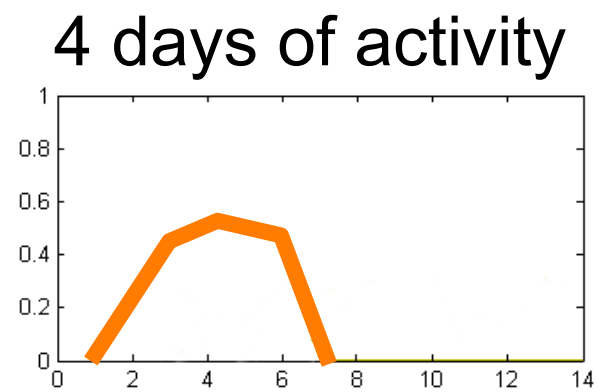
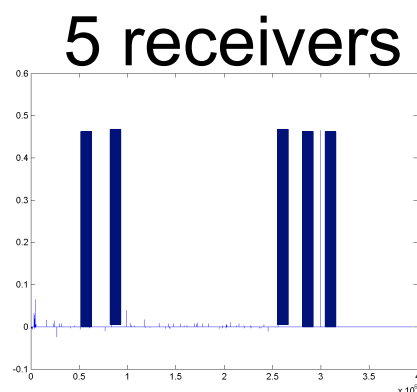
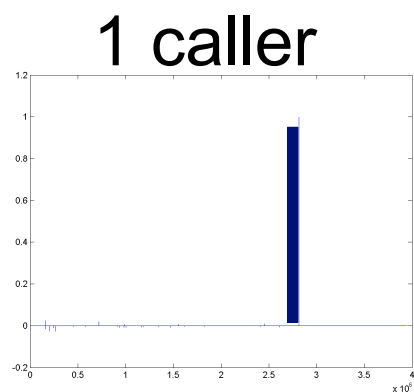
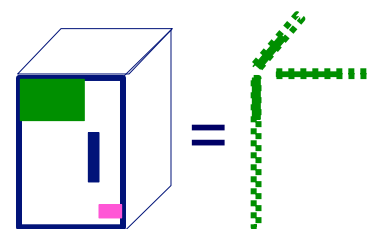
CONCLUSION#1 – Big data

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



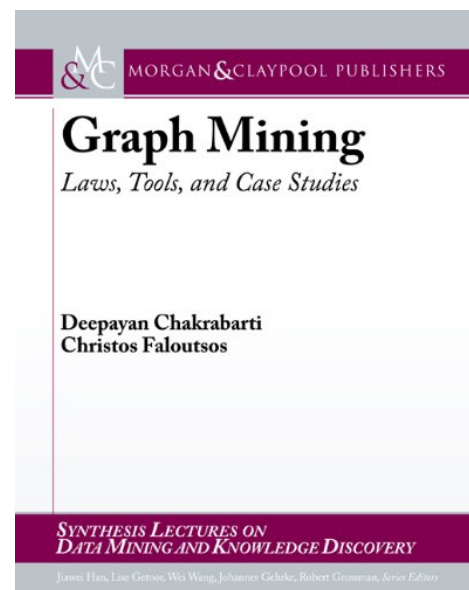
CONCLUSION#2 – tensors

- powerful tool



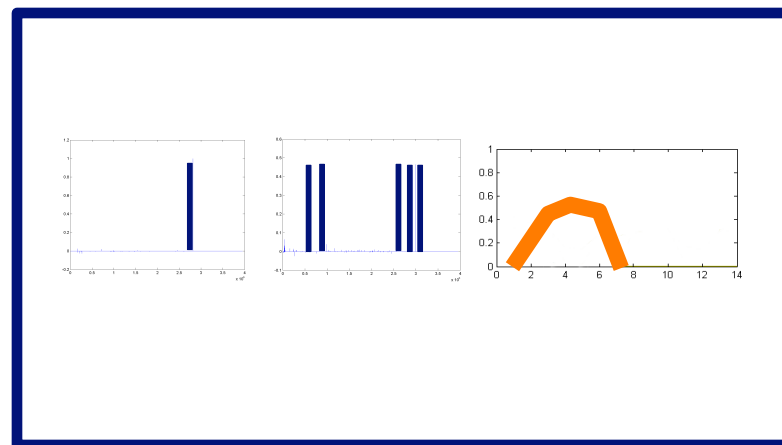
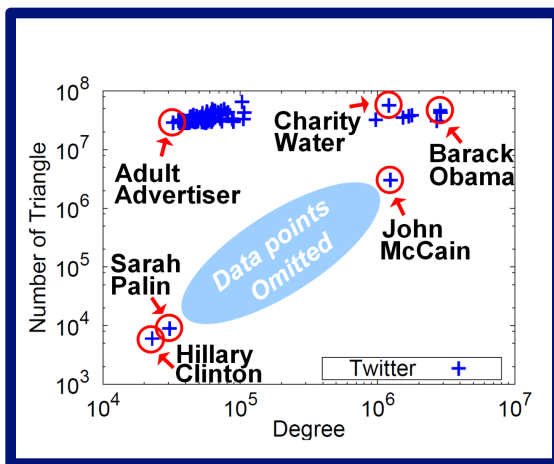
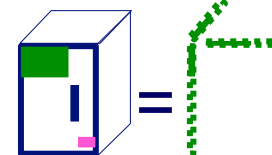
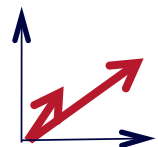
References

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



TAKE HOME MESSAGE:

Cross-disciplinarity



Thank you!

Cross-disciplinarity

