# Anomaly detection in large graphs
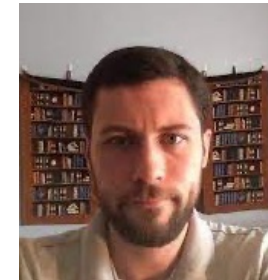
*Christos Faloutsos*

CMU

# Thank you!

- Badel Mbanga

- Clifton Denning

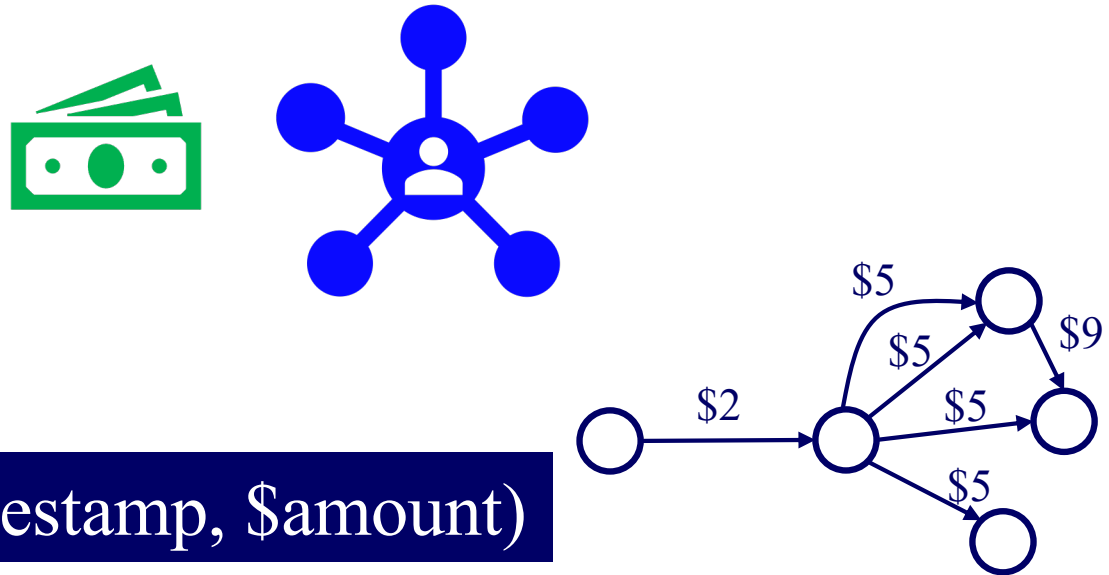- Matthew Berezo

Christos Faloutsos

# Roadmap

➡ • Introduction – Motivation
    – Why study (big) graphs?
• Part#1: Patterns in graphs
• Part#2: time-evolving graphs; tensors
• Conclusions
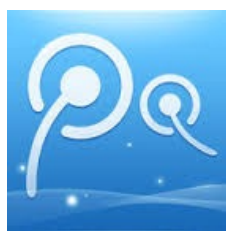
# Graphs – why should we care?
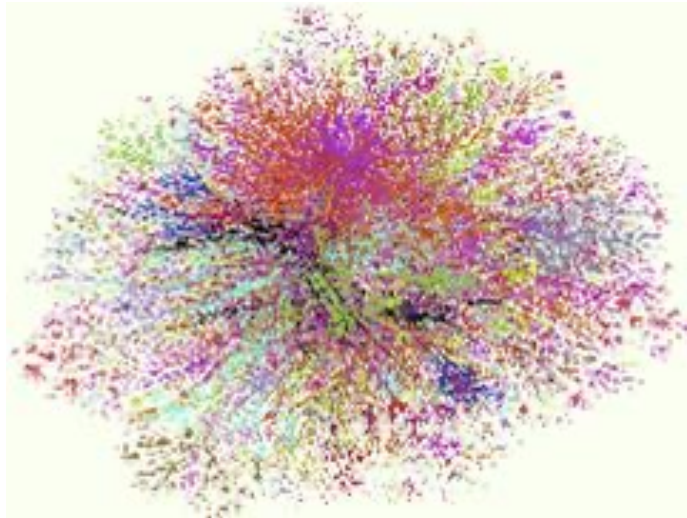
PNC

(source, destination, timestamp, $amount)

$5

$5
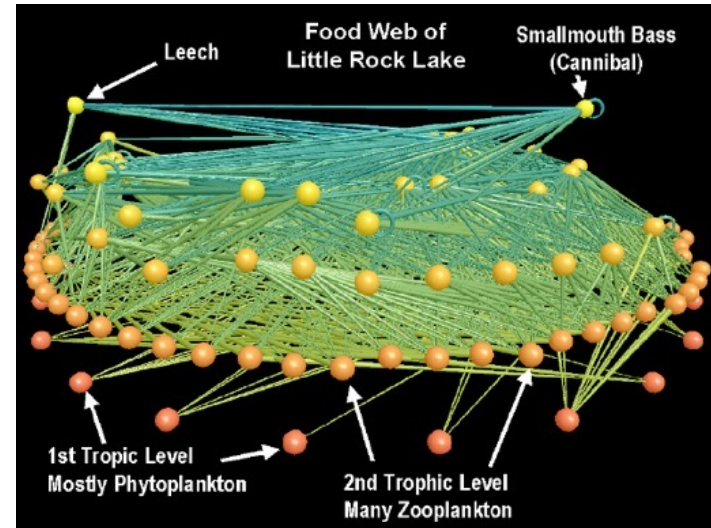
$9

$2

$5

$5

# Graphs - why should we care?



>$10B; ~1B users

# Graphs - why should we care?



## Internet Map
## [lumeta.com]



Food Web of Little Rock Lake

Leech

Smallmouth Bass (Cannibal)

1st Tropic Level Mostly Phytoplankton
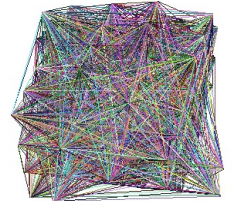
2nd Trophic Level Many Zooplankton

## 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?

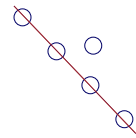- P2: patterns in time-evolving graphs / tensors

destination

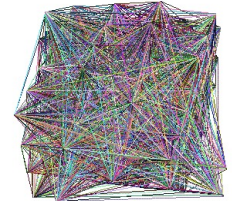source    time

# Motivating problems

- P1: patterns? Fraud detection?

Patterns    anomalies

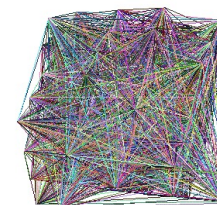- P2: patterns in time-evolving graphs / tensors

destination

source    time

# Roadmap



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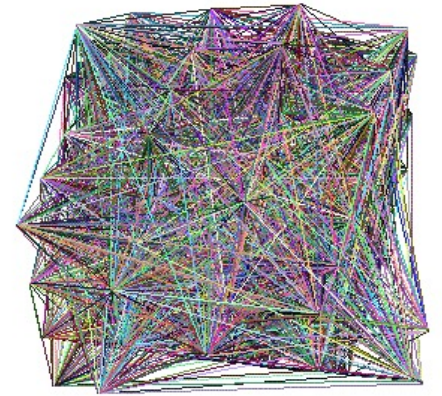- Part#2: time-evolving graphs; tensors
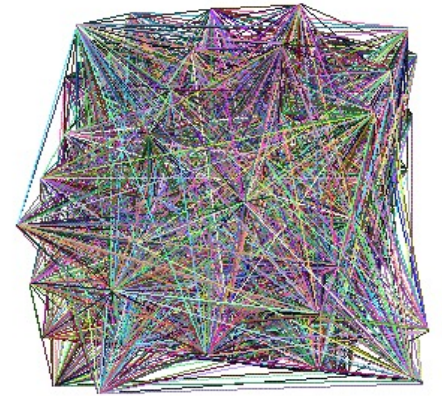- Conclusions

# Part 1: Patterns, & fraud detection

# Laws and patterns

- Q1: Are real graphs random?

Christos Faloutsos
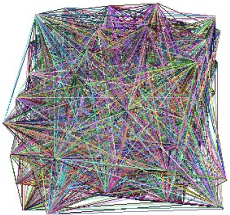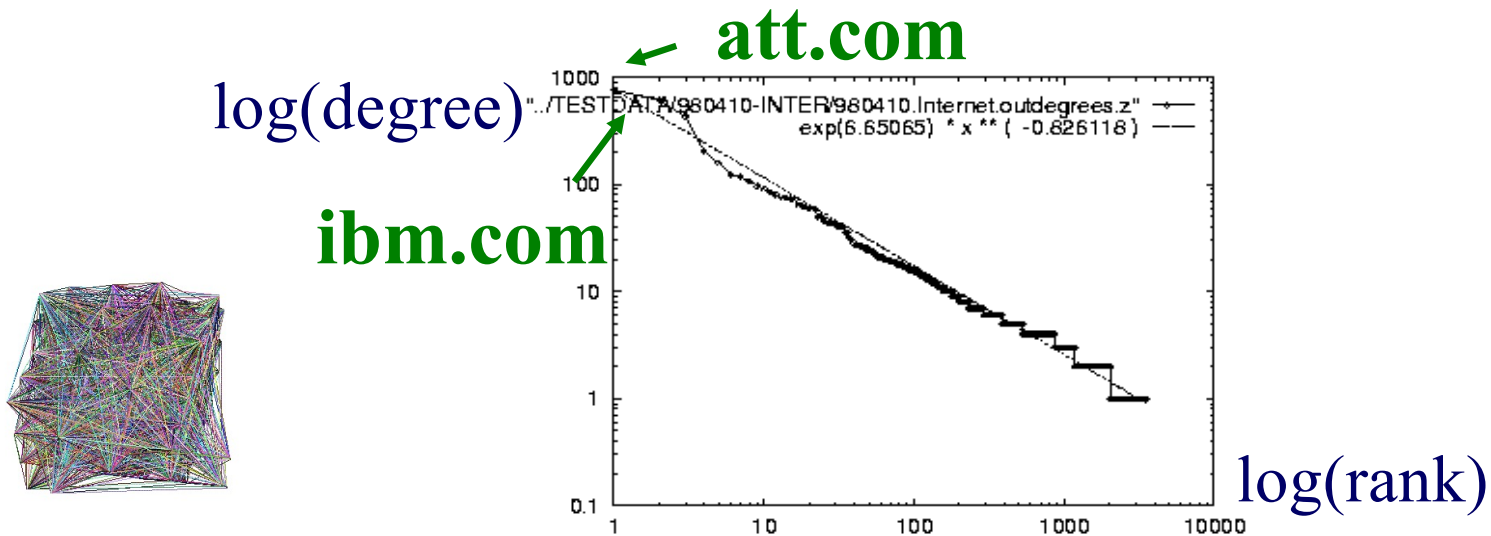
# Laws and patterns

- 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

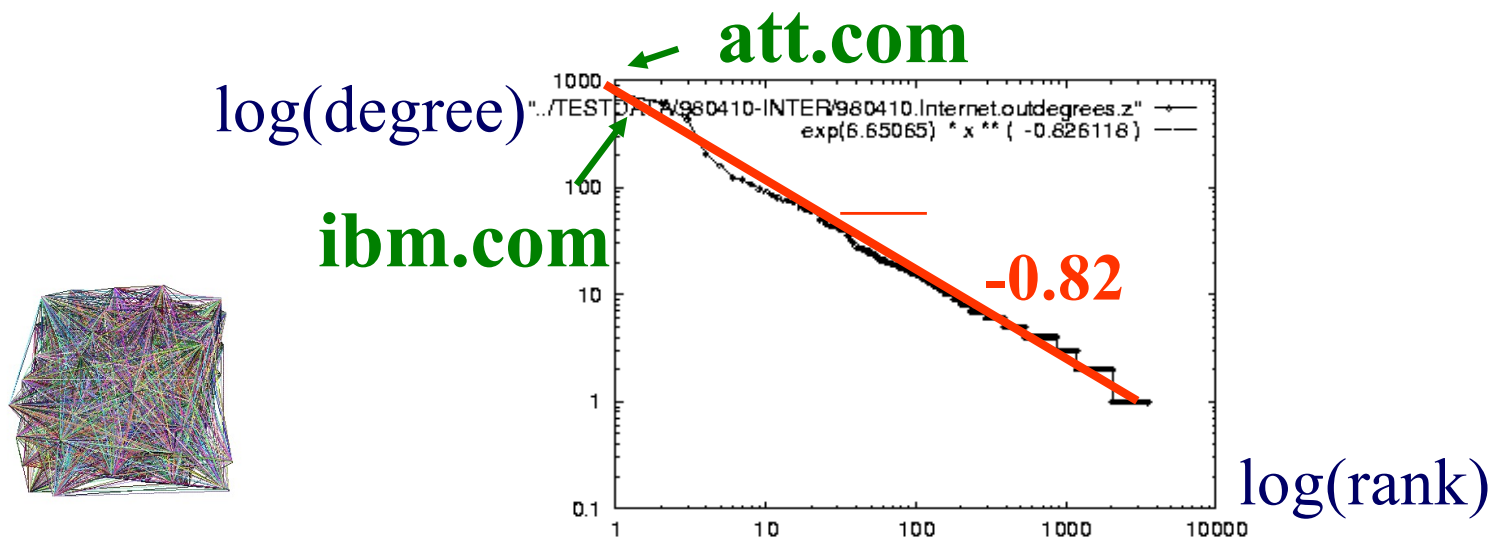- Power law in the degree distribution [Faloutsos x 3 SIGCOMM99]

**internet domains**

att.com

log(degree)
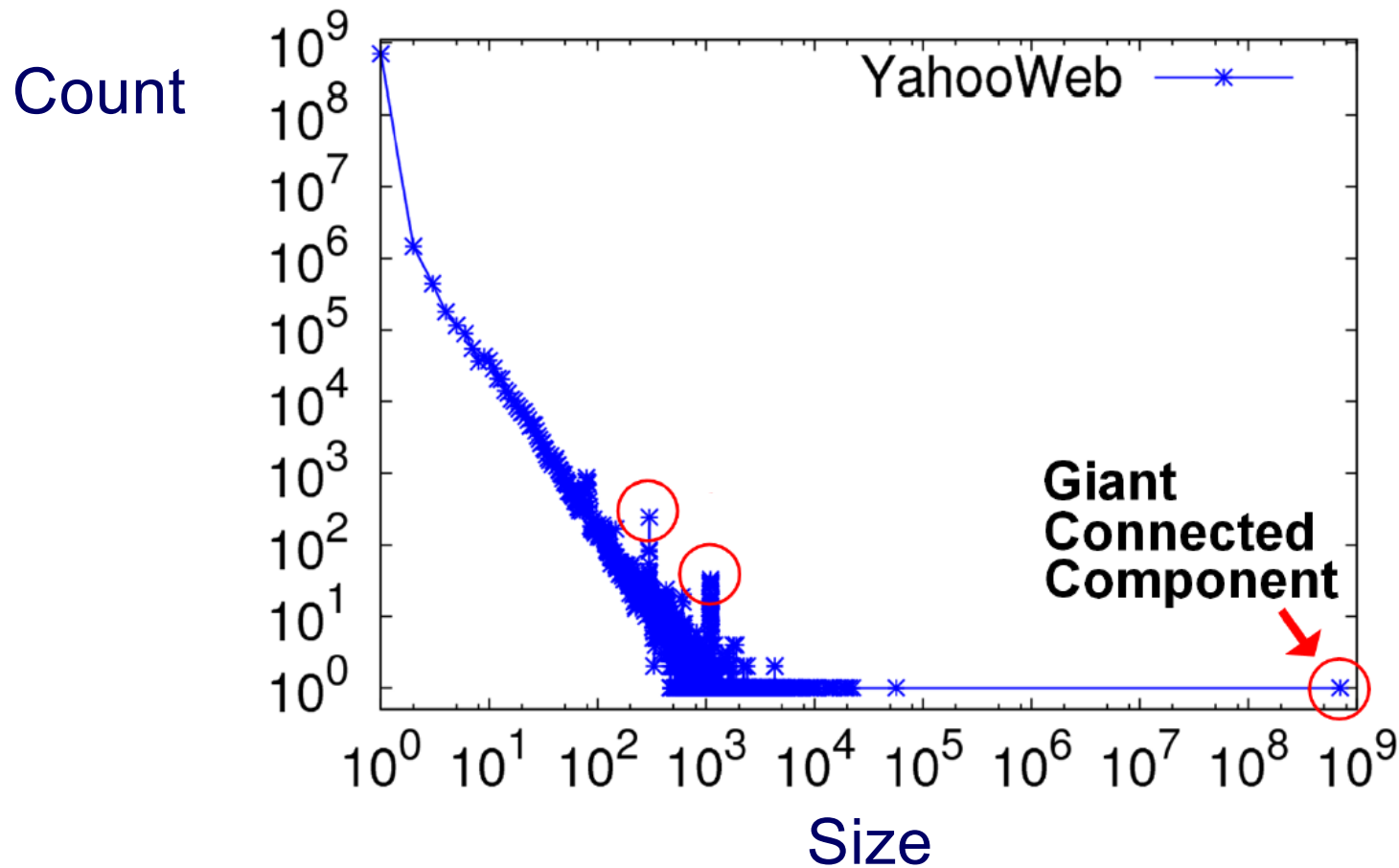
ibm.com



log(rank)

# Solution# S.1

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

**internet domains**

# S2: connected component sizes
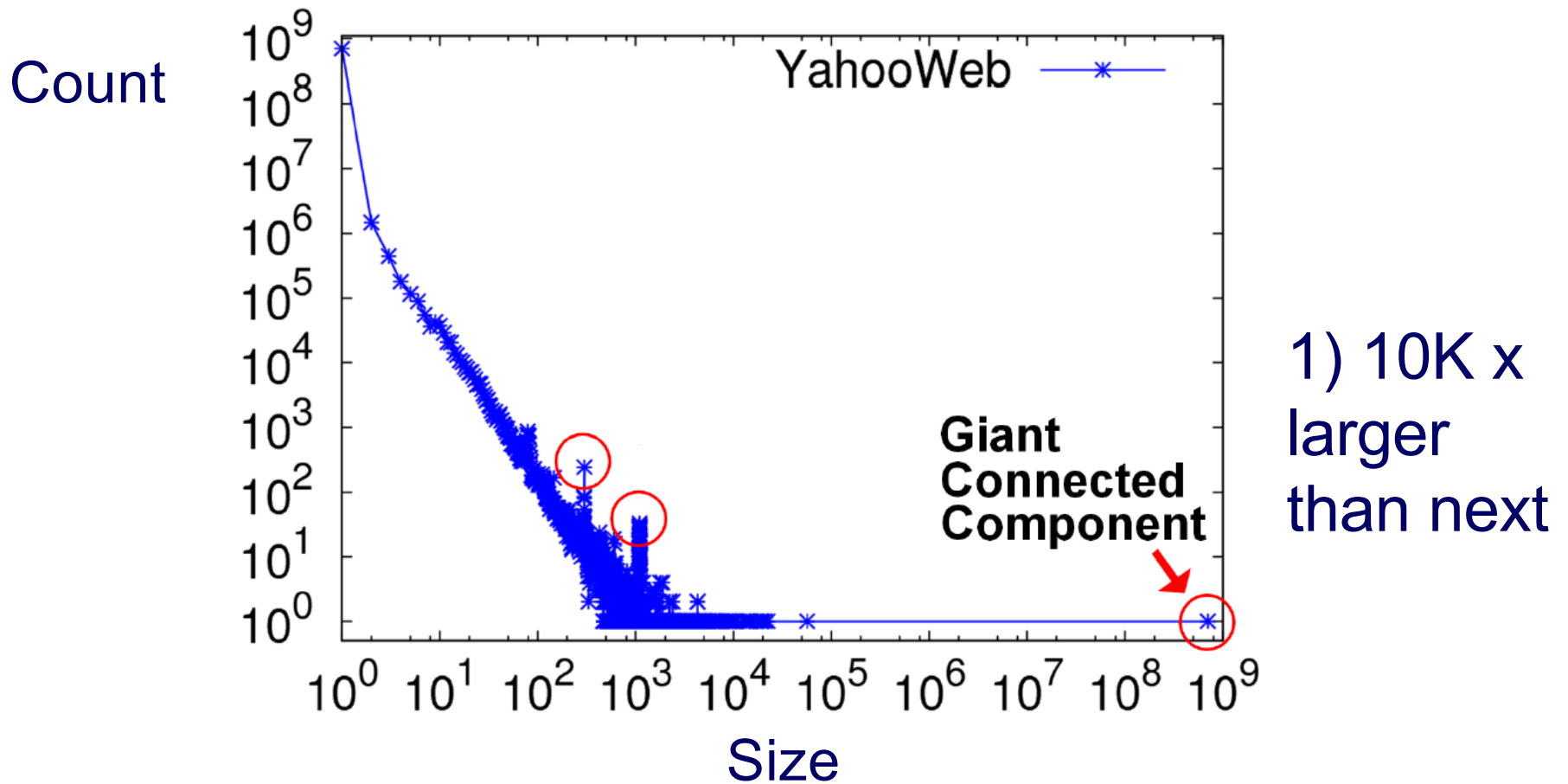
- Connected Components – 4 observations:



Count

YahooWeb

1.4B nodes
6B edges

Giant Connected Component

Size

# S2: connected component sizes

- Connected Components



1) 10K x larger than next

Christos Faloutsos

# S2: connected component sizes

- Connected Components

Count

2) ~0.7B singleton nodes

# S2: connected component sizes

- Connected Components

**Count**

**3) SLOPE!**

YahooWeb

**Giant Connected Component**

Size

# S2: connected component sizes

- ## Connected Components



Count

YahooWeb

300-size cmpt X 500. Why?

1100-size cmpt X 65. Why?

Giant Connected Component

4) Spikes!

Size
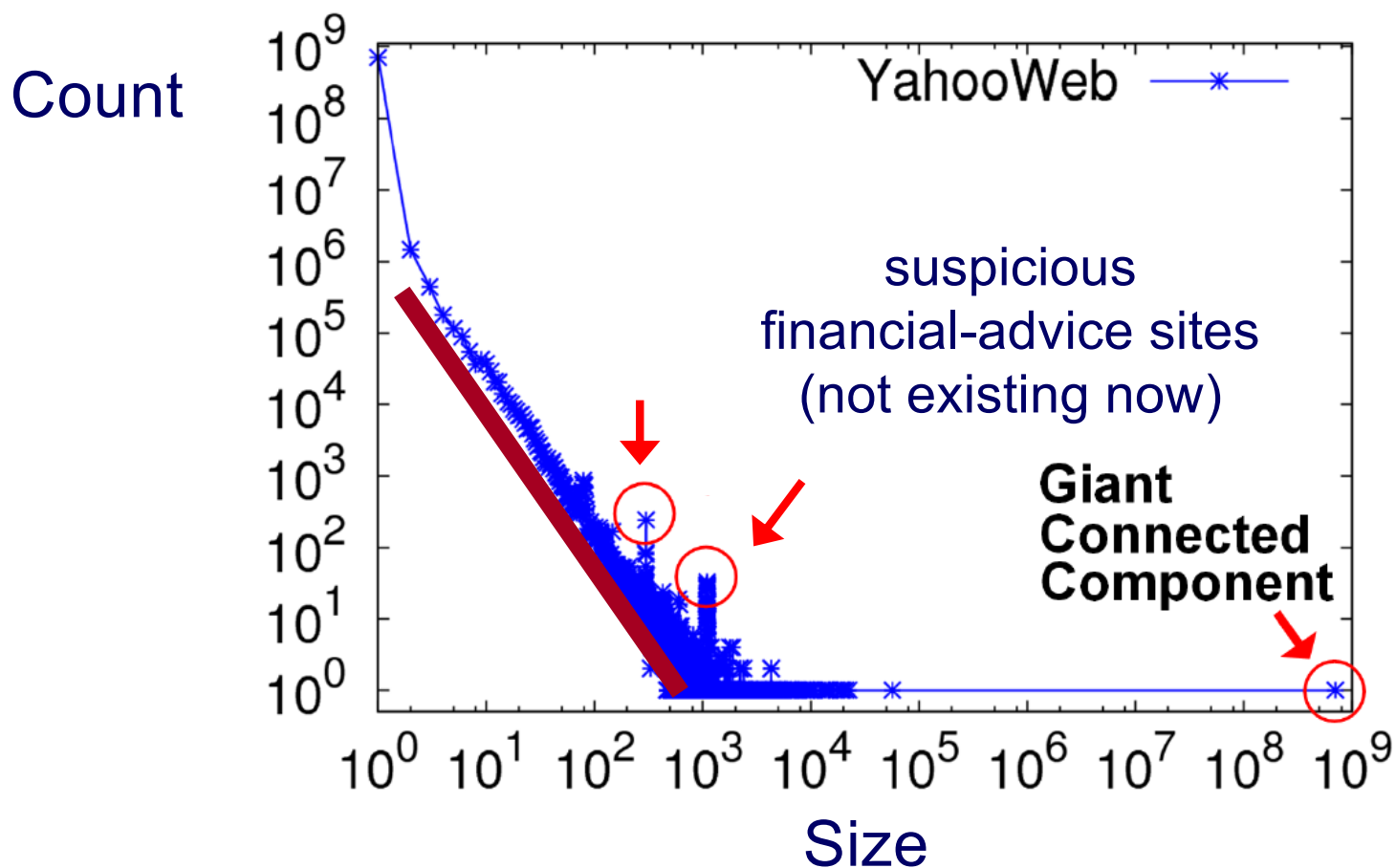
# S2: connected component sizes
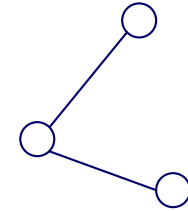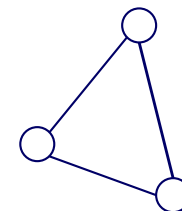
- Connected Components

# Roadmap



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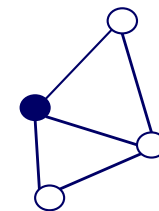- Conclusions

# Solution# S.3: Triangle 'Laws'

- Real social networks have a lot of triangles
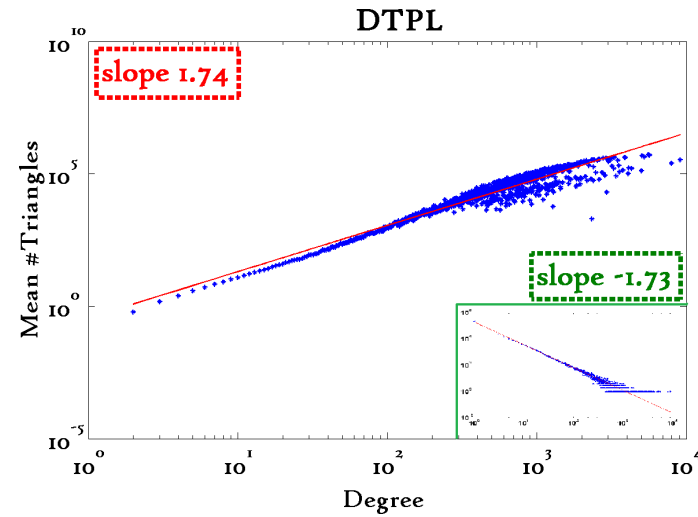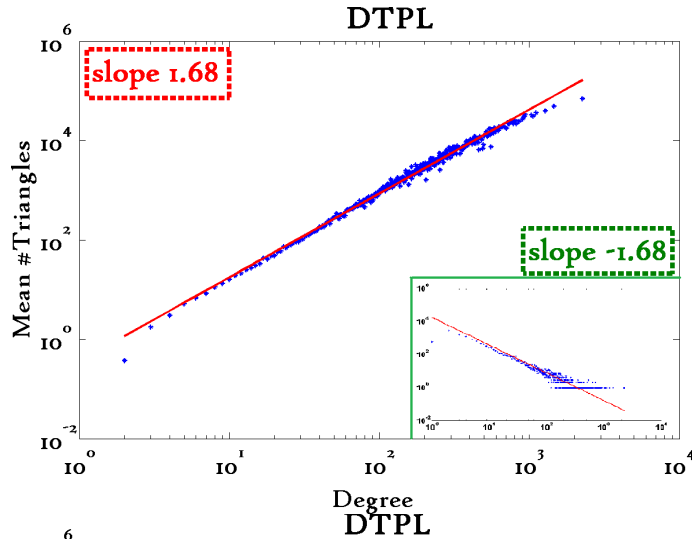
# Solution# S.3: Triangle 'Laws'

- Real social networks have a lot of triangles
  - Friends of friends are friends

- Any patterns?
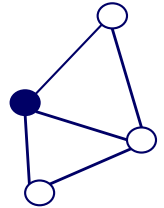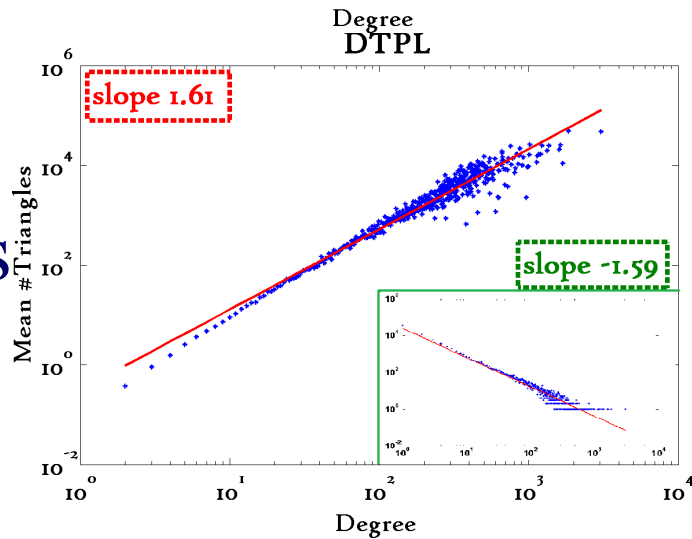  - 2x the friends, 2x the triangles ?

# Triangle Law: #S.3
## [Tsourakakis ICDM 2008]



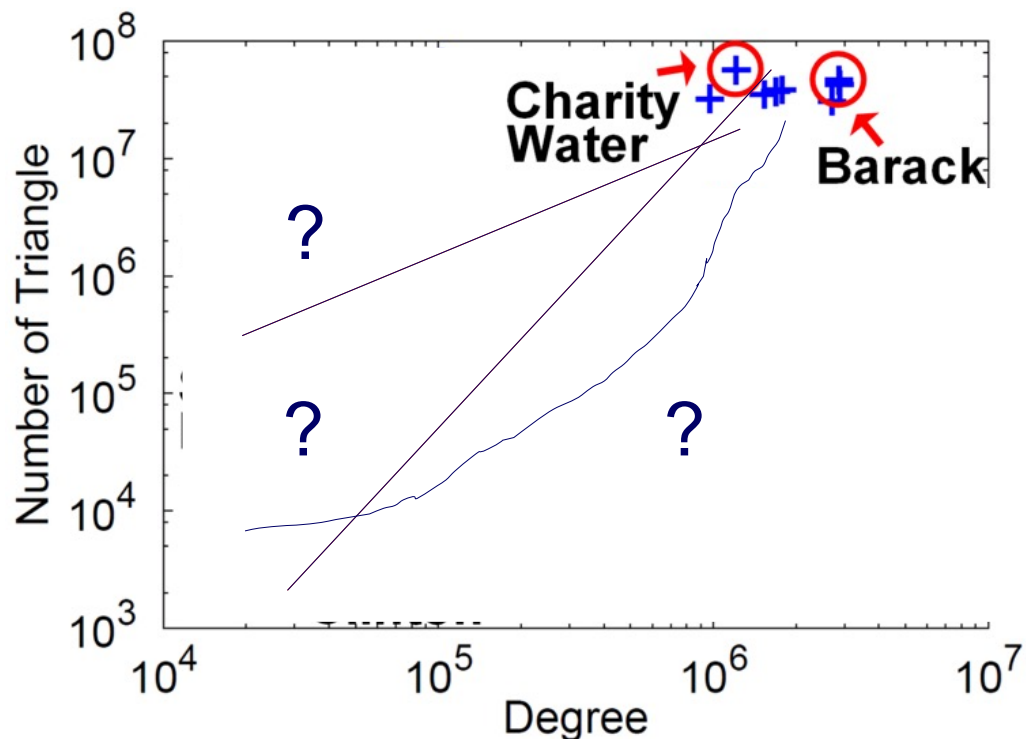Reuters

SN

Epinions

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

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)

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

Christos 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]

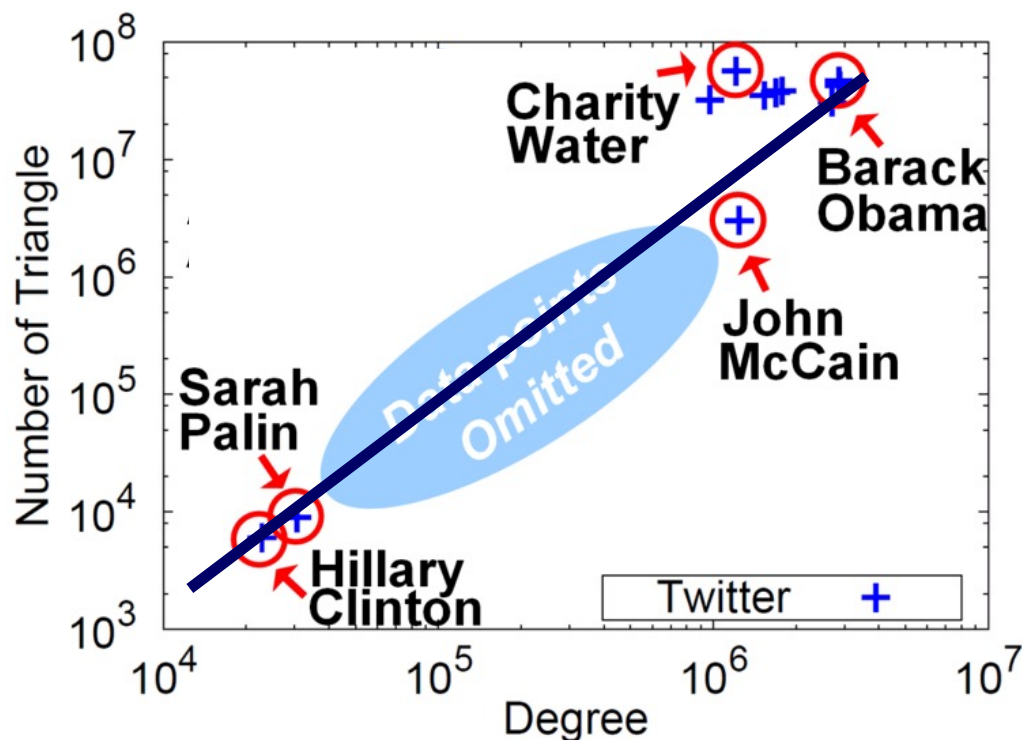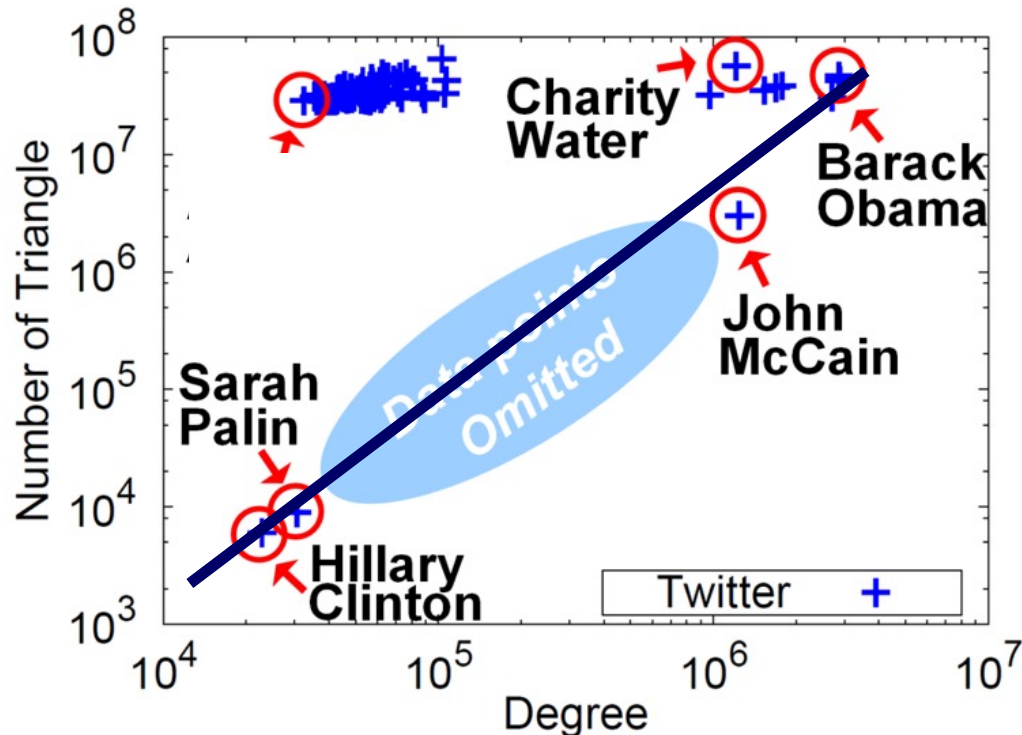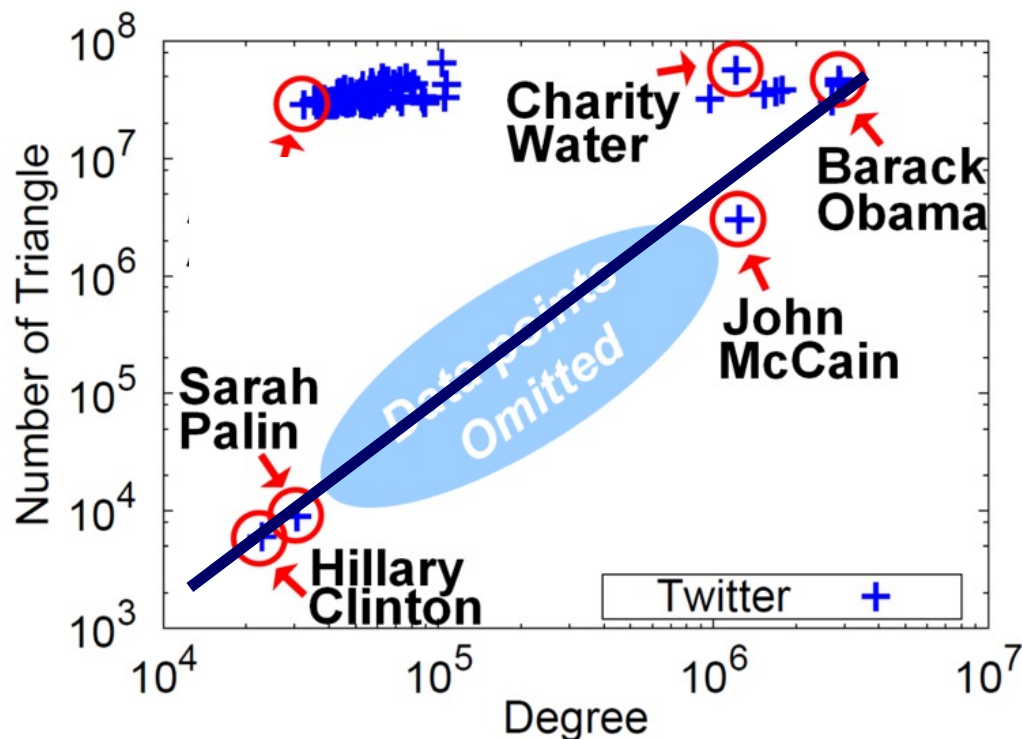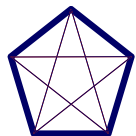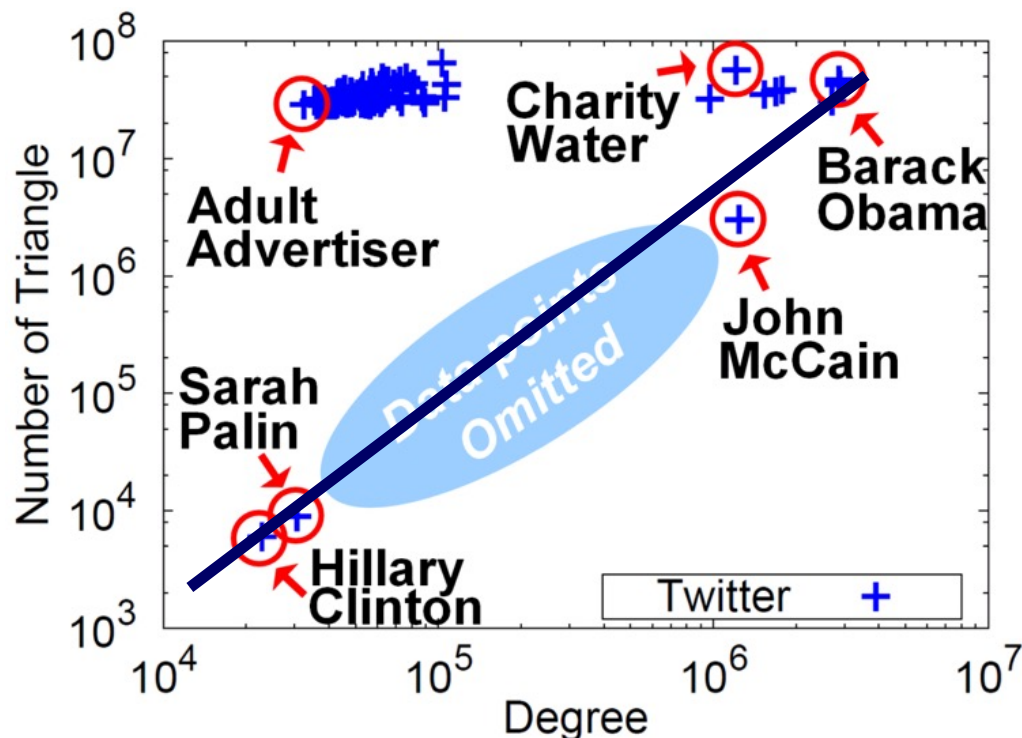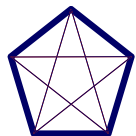# Triangle counting for large graphs?



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

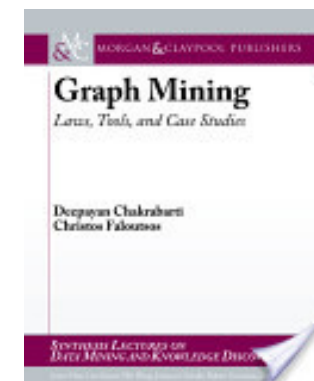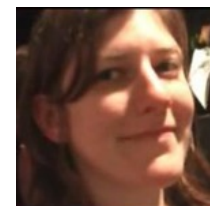| | Unweighted | Weighted |
|---|---|---|
| **Static** | ✓ **L01.** Power-law degree distribution [Faloutsos et al. `99, Kleinberg et al. `99, Chakrabarti et al. `04, Newman `04]<br>✓ **L02.** Triangle Power Law (TPL) [Tsourakakis `08]<br>✓ **L03.** Eigenvalue Power Law (EPL) [Siganos et al. `03]<br>**L04.** Community structure [Flake et al. `02, Girvan and Newman `02] | **L10.** Snapshot Power Law (SPL) [McGlohon et al. `08] |
| **Dynamic** | **L05.** Densification Power Law (DPL) [Leskovec et al. `05]<br>**L06.** Small and shrinking diameter [Albert and Barabási `99, Leskovec et al. `05]<br>**L07.** Constant size 2nd and 3rd connected components [McGlohon et al. `08]<br>**L08.** Principal Eigenvalue Power Law ($\lambda_1$PL) [Akoglu et al. `08]<br>**L09.** Bursty/self-similar edge/weight additions [Gomez and Santonja `98, Gribble et al. `98, Crovella and | **L11.** Weight Power Law (WPL) [McGlohon et al. `08] |

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

# MORE Graph Patterns

|  | Unweighted | Weighted |
|---|---|---|
| Static | L01. Power-law degree distribution [Faloutsos et al. `99, Kleinberg et al. `99, Chakrabarti et al. `04, Newman `04] <br> L02. Triangle Power Law (TPL) [Tsourakakis `08] <br> L03. Eigenvalue Power Law (EPL) [Siganos et al. `03] <br> L04. Community structure [Flake et al. `02, Girvan and Newman `02] | L10. Snapshot Power Law (SPL) [McGlohon et al. `08] |
| Dynamic | L05. Densification Power Law (DPL) [Leskovec et al. `05] <br> L06. Small and shrinking diameter [Albert and Barabási `99, Leskovec et al. `05] <br> L07. Constant size 2$^{nd}$ and 3$^{rd}$ connected components [McGlohon et al. `08] <br> L08. Principal Eigenvalue Power Law ($\lambda_1$PL) [Akoglu et al. `08] <br> L09. Bursty/self-similar edge/weight additions [Gomez and Santonja `98, Gribble et al. `98, Crovella and Bestavros `99, McGlohon et al. `08] | L11. Weight Power Law (WPL) [McGlohon et al. `08] |

• 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

  Patterns  anomalies

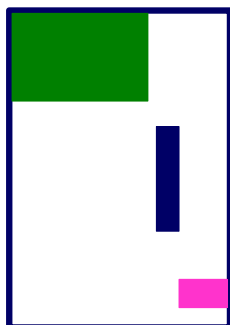- Part#2: time-evolving graphs; tensors
- Conclusions

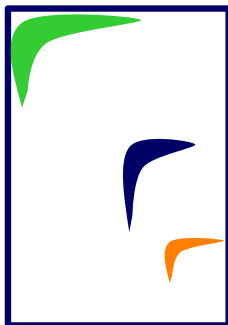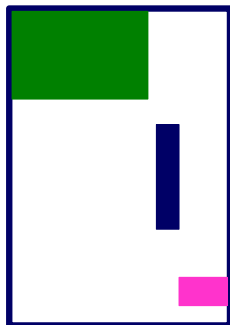# How to find 'suspicious' groups?

- 'blocks' are normal, right?

idols

fans

# Except that:

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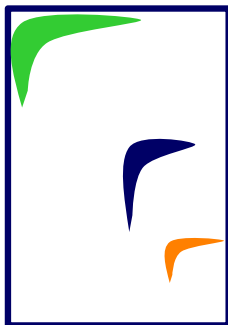- 'hyperbolic' communities are more realistic [Araujo+, PKDD'14]

# **Except that:**

- 'blocks' are usually suspicious
- 'hyperbolic' communities are more realistic [Araujo+, PKDD'14]
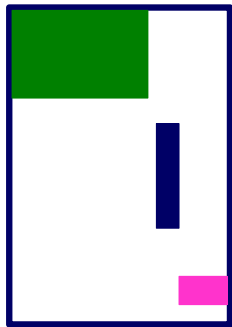
Q: Can we spot blocks, easily?

# 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

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



M idols

N fans

$\vec{v_1}$

$\sim$    $\vec{u_1}$    +    +    $\vec{u_i}$

# Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks



'meat-eaters' 'steaks'   'vegetarians' 'plants'   'kids' 'cookies'

M products

N users
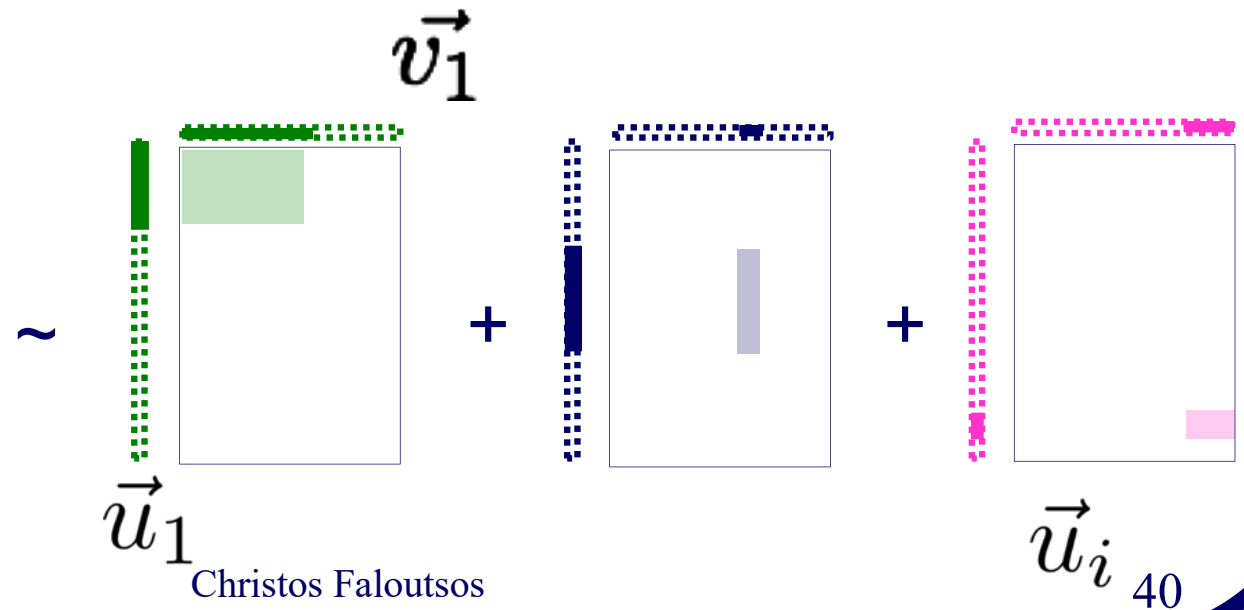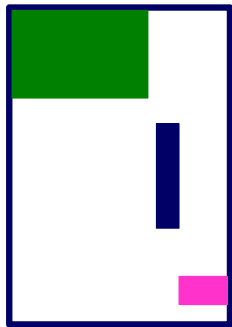
$\vec{v_1}$

$\vec{u_1}$   $\vec{u_i}$

~   +   +

# Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks



'cancer'  'alzheimer'  'Parkinson'

M timestamps

N genes

$\vec{v_1}$

$\sim$  $\vec{u_1}$  $+$  $+$  $\vec{u_i}$
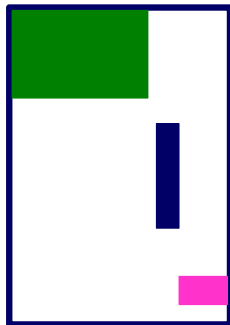
PNC 2023

Christos Faloutsos

40

# Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks

M idols

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

$\vec{v_1}$

N fans

$\sim$ $\vec{u_1}$ $+$ $+$ $\vec{u_i}$

Christos Faloutsos

41

# Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks

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

M idols

N fans

$\vec{v_1}$

$\sim$ $\vec{u_1}$ $+$ $+$ $\vec{u_i}$

Christos Faloutsos

42

# Crush intro to SVD

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

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

$\vec{v_1}$

M idols

N fans

$\sim$ $\vec{u_1}$ + + $\vec{u_i}$

Christos Faloutsos

43

# Dataset

- Tencent Weibo 

- 117 million nodes (with profile and UGC data)

- 3.33 billion directed edges

# Real Data

## "Rays"



## "Block"



Christos Faloutsos

# 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
    - No labels – accounting application
    - With labels: Belief Propagation
- Part#2: time-evolving graphs; tensors
- Conclusions

# AutoAudit: Mining Accounting and Time-Evolving Graphs

*IEEE Big Data, 2020*

Meng-Chieh Lee[1], Yue Zhao[2], Aluna Wang[2], Pierre Jinghong Liang[2], Leman Akoglu[2], Vincent S. Tseng[1], Christos Faloutsos[2]

# 'Smurfing'

'Alan'

'Bob'

How to spot it?

# 'Smurfing'



Receiver

Sender

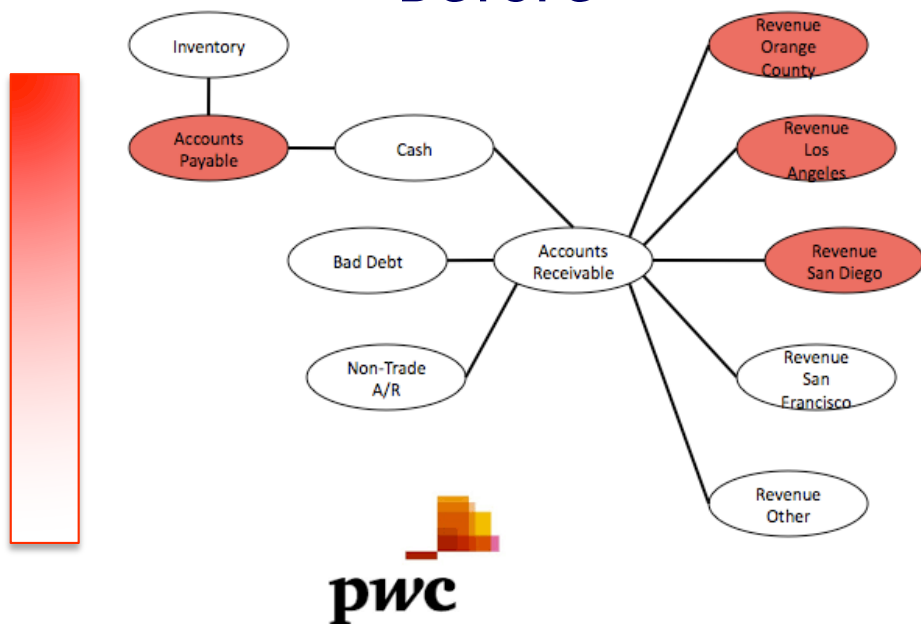'Alan'    'Bob'

'Alan'

'Bob'

'Alan'

'Bob'

'smurfs'

# AutoAudit: Experiments



"Smurfing"

Reordering



Effectiveness in Accounting Dataset

Effectiveness in Czech Dataset

← Ideal: 100%

accuracy

- AA_Smurf
- AA_Smurf w/o Purity
- AA_Smurf w/o MDL
- Number of Intermediaries
- FlowScope

# Roadmap



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

# Network Effect Tools: SNARE

- Some accounts are sort-of-suspicious – how to combine weak signals?

**Before**

# Network Effect Tools: SNARE

- A: Belief Propagation.

Before

# Network Effect Tools: SNARE

- A: Belief Propagation.

Mary McGlohon, Stephen Bay, Markus G. Anderle, David M. Steier, Christos Faloutsos: *SNARE: a link analytic system for graph labeling and risk detection*. KDD 2009: 1265-1274

# Network Effect Tools: SNARE

- Produces improvement over simply using flags
  - Up to 6.5 lift
  - Improvement especially for low false positive rate

**Results for accounts data (ROC Curve)**

Ideal →

True positive rate

SNARE

Baseline (flags only)

False positive rate

# 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

Christos Faloutsos

# Part 2: Time evolving graphs; tensors

Christos Faloutsos

# Graphs over time -> tensors!

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

johnson

smith

# 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

Tue

Mon

# Graphs over time -> tensors!

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



time

caller

callee

Christos Faloutsos

# Answer : tensor factorization

- Recall: (SVD) matrix factorization: finds blocks



'meat-eaters' 'steaks'    'vegetarians' 'plants'    'kids' 'cookies'

M products

N users

$\vec{v_1}$

$\vec{u_1}$

$\vec{u_i}$

# Crush intro to SVD

- Recall: (SVD) matrix factorization: finds blocks



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

M idols

N fans

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

Christos Faloutsos

# Answer: tensor factorization

- PARAFAC decomposition



politicians     artists     athletes

verb

subject

object

# 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 -> tensors

- PARAFAC finds patterns

- Surprising temporal patterns

Christos Faloutsos

# Roadmap



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

# *TgraphSpot*: Fast and Effective Anomaly Detection for Time-Evolving Graphs
## *IEEE BigData, 2022*

Mirela Cazzolato[1,2], Saranya Vijayakumar[1], Xinyi Zheng[1], Namyong Park[1], Meng-Chieh Lee[1], Pedro Fidalgo[3,4], Bruno Lages[3], Agma J. M. Traina[2], Christos Faloutsos[1]

Open source:
https://github.com/mtcazzolato/tgraph-spot

Video: https://youtu.be/jI1adN-BQuo?t=1537

# Authors

Mirela Cazzolato

Saranya Vijayakumar

Xinyi Zheng

Namyong Park

Meng-Chieh Jeremy Lee

Pedro Fidalgo

Bruno Lages

Agma Traina

Christos Faloutsos

75

Christos Faloutsos

# Problem definition



(source, destination, timestamp, duration)

# Problem definition

(source, destination, timestamp, $amount)

$5

$5

$9

$2

$5

$5

Christos Faloutsos

# System Overview - current



Feature extraction

Select nodes for further investigation

Feature visualization

Deep Dive: EgoNet

Video: https://youtu.be/jI1adN-BQuo?t=1537

# Discovery #1

in-degree



Weighted in-degree (= in-seconds)

# Discovery #1



100 in-calls
100 seconds

# Discovery #1

# Q: Why?

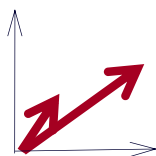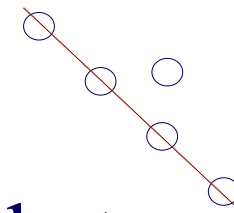- Q: Why would people call hotel-like numbers, for 1second?

# Roadmap



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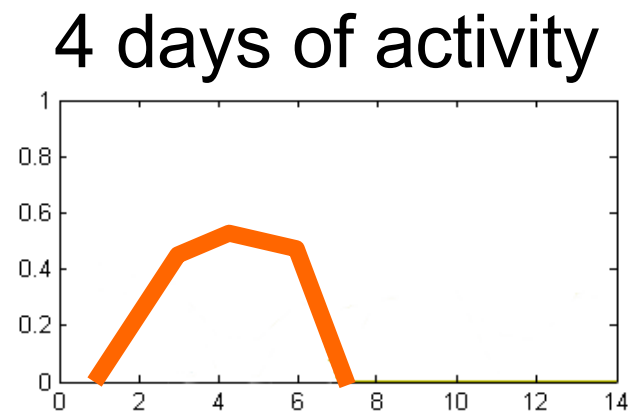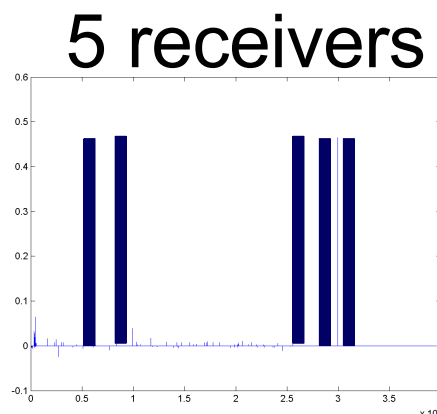- Part#2: time-evolving graphs; tensors
- Visualization
- Conclusions

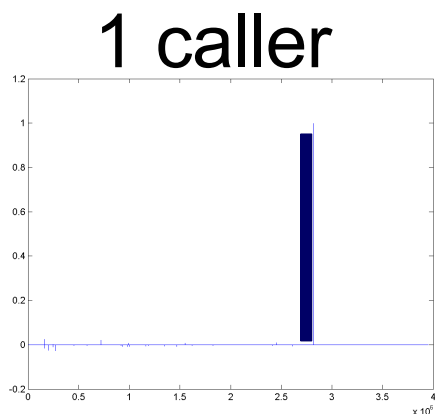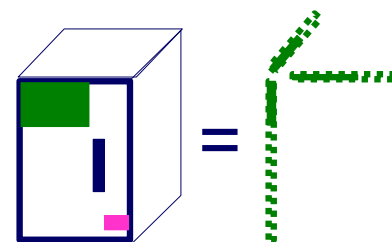# CONCLUSION#1 – Big data

- **Patterns** ⤫ **Anomalies**

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

# CONCLUSION#2 – tensors
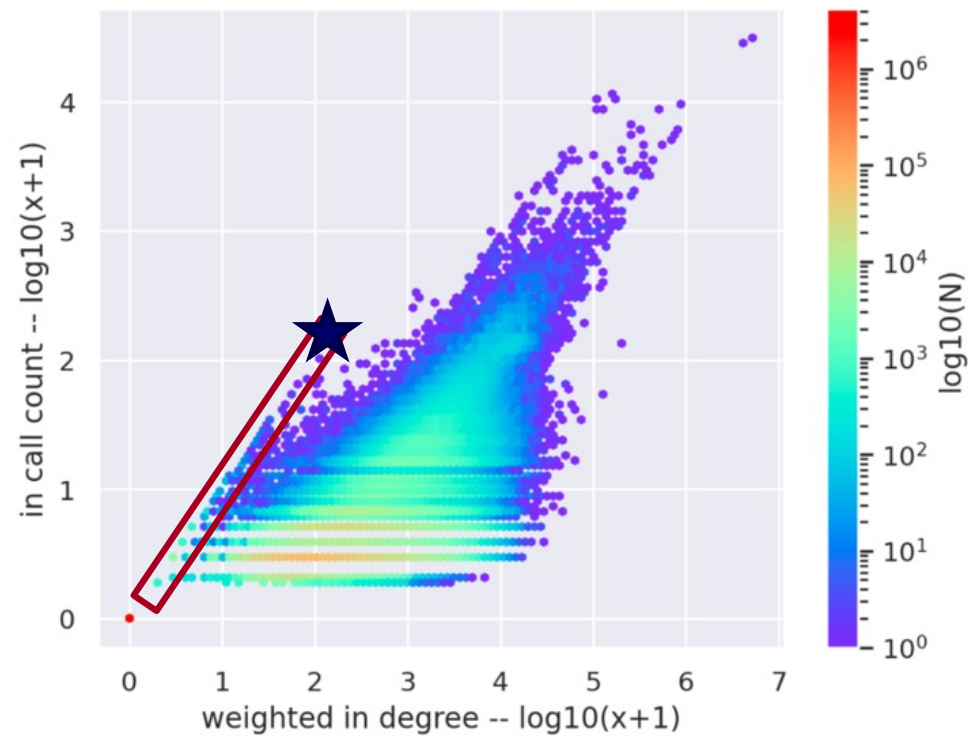
- powerful tool



## 1 caller          5 receivers       4 days of activity

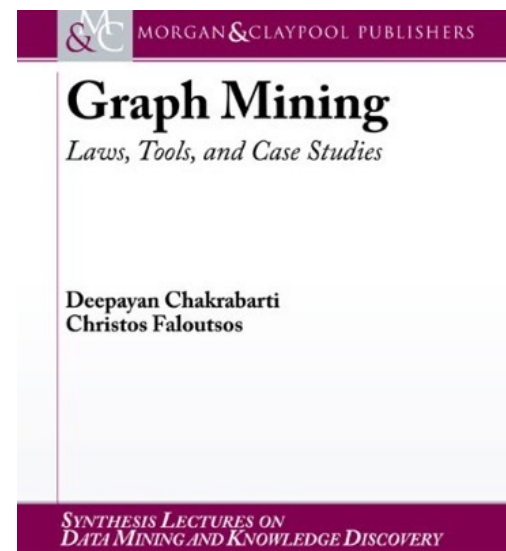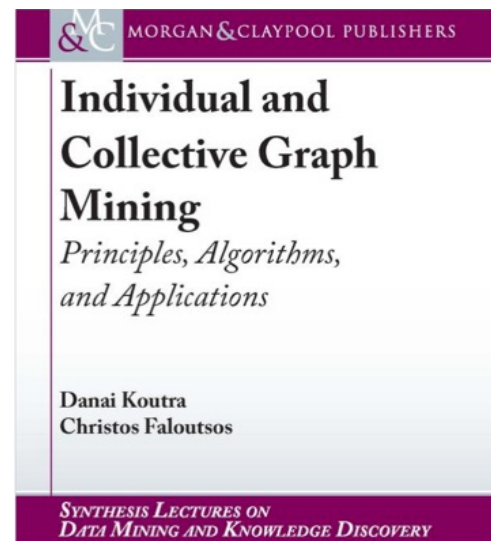Christos Faloutsos

# CONCLUSION#3 - visualization

in-degree



Weighted in-degree (= in-seconds)

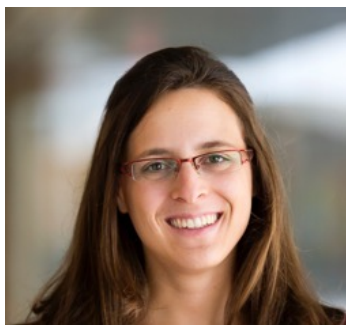# References

- D. Chakrabarti, C. Faloutsos: *Graph Mining – Laws, Tools and Case Studies*, Morgan Claypool 2012
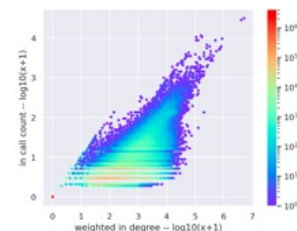- http://www.morganclaypool.com/doi/abs/10.2200/S004 49ED1V01Y201209DMK006

# References

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

# *TgraphSpot*: Fast and Effective Anomaly Detection for Time-Evolving Graphs
## *IEEE BigData, 2022*

Mirela Cazzolato[1,2], Saranya Vijayakumar[1], Xinyi Zheng[1],
Namyong Park[1], Meng-Chieh Lee[1], Pedro Fidalgo[3,4],
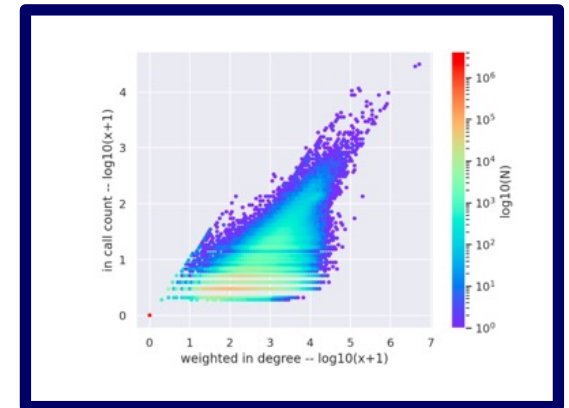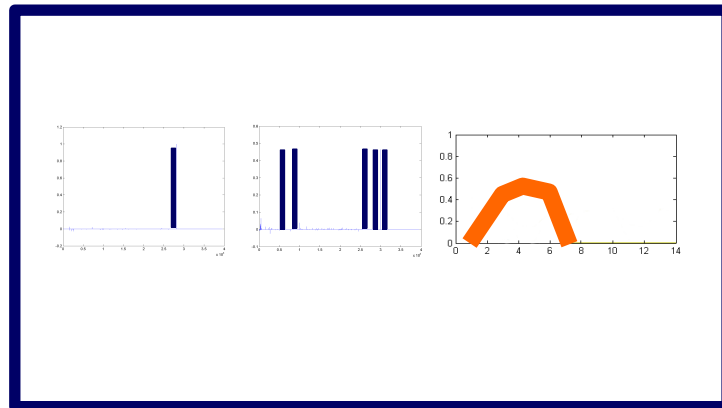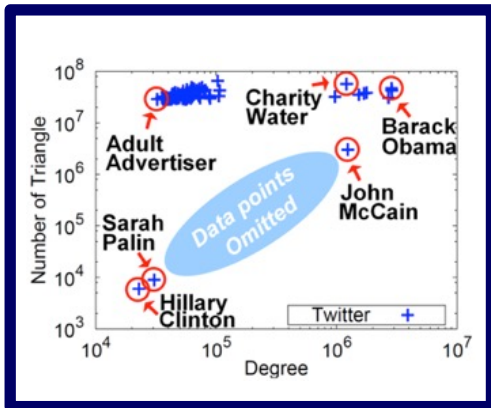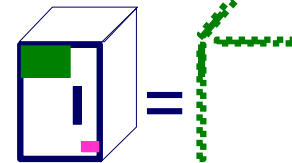Bruno Lages[3], Agma J. M. Traina[2], Christos Faloutsos[1]

Open source:
`https://github.com/mtcazzolato/tgraph-spot`

Video: https://youtu.be/jI1adN-BQuo?t=1537

# TAKE HOME MESSAGE:

# Cross-disciplinarity

# Thank you!

## Cross-disciplinarity

Christos Faloutsos