

## 15-826: Multimedia (Databases) and Data Mining

Anomaly detection in large graphs

Christos Faloutsos



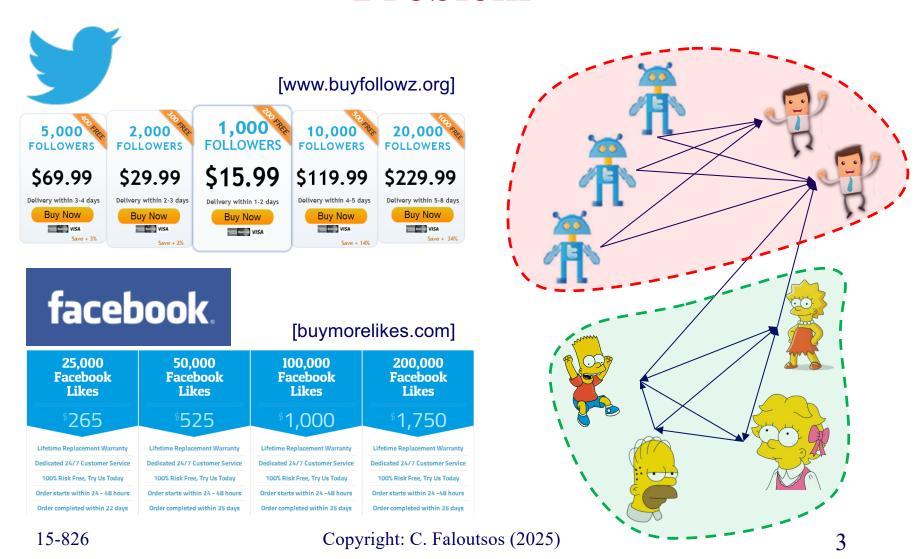
#### Roadmap



- Introduction Motivation
- Unsupervised
  - static graphs
  - Time-evolving graphs
- Supervised: BP
- Conclusions

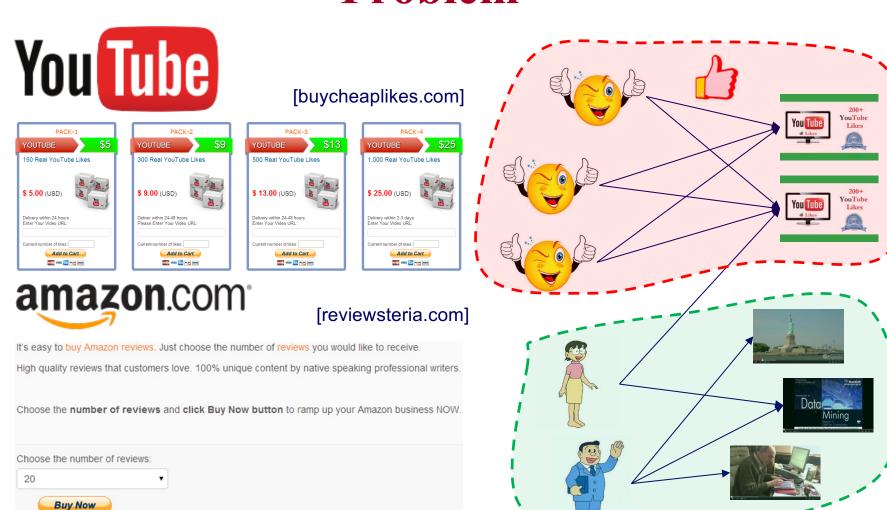


## Fraud Detection: Graph Analysis Problem





## Fraud Detection: Graph Analysis Problem



15-826

UISA E PRO COSSO

Copyright: C. Faloutsos (2025)

4



#### Other applications:

- Cyber-security (DDoS attacks distributed denial of service)
- Abnormal individuals / social groups
  - dysfunctional departments in a company
  - Onset of depression/hijacking wrt a FB user

**—** ...



#### Roadmap

- Introduction Motivation
- Unsupervised
  - static graphs

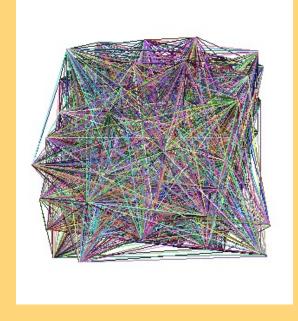


- Eigenspokes, and extensions
- Spot strange, individual nodes: oddBall
- Time-evolving graphs
- Supervised: BP
- Conclusions



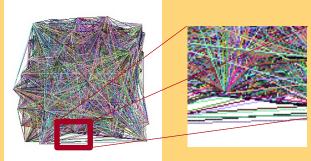
#### **Problem**

#### Given:



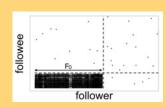
#### Find:

- 1) Outliers
- 2) Lock-step
- 3) B.P.





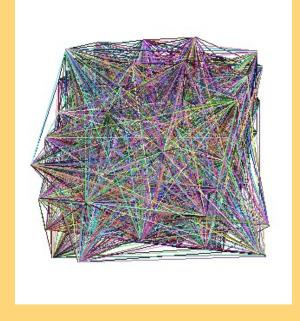


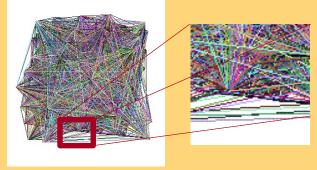


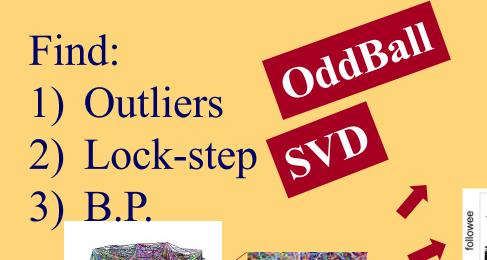


#### **Solutions**

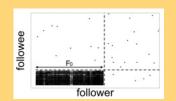
#### Given:









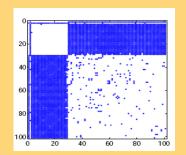


semi-supervised

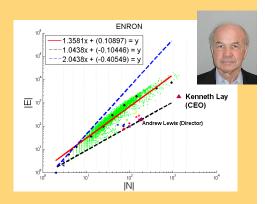


#### Solutions, cont'd

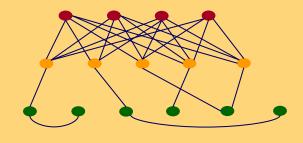
• Spectral/tensor methods: spot lockstep behavior



OddBall: strange, single nodes



• BP: good for the supervised setting (ie., some labels are given)





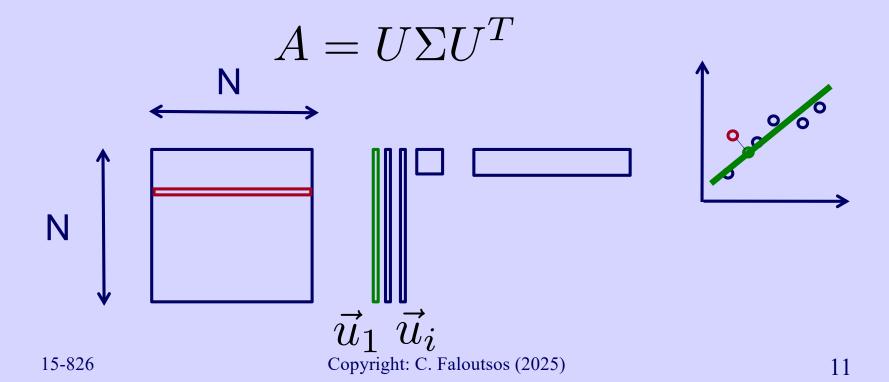


B. Aditya Prakash, Mukund Seshadri, Ashwin Sridharan, Sridhar Machiraju and Christos Faloutsos: *EigenSpokes: Surprising Patterns and Scalable Community Chipping in Large Graphs*, PAKDD 2010, Hyderabad, India, 21-24 June 2010.



#### **EigenSpokes**

- Eigenvectors of adjacency matrix
  - equivalent to singular vectors (symmetric, undirected graph)

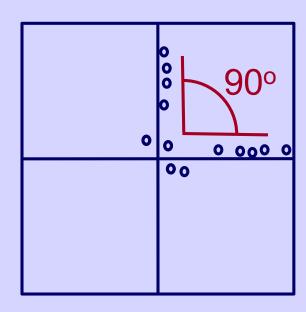


#### **EigenSpokes**

**u**2

- EE plot:
- Scatter plot of scores of u1 vs u2
- One would expect
  - Many points @origin

- A few ttered ~r tu mly



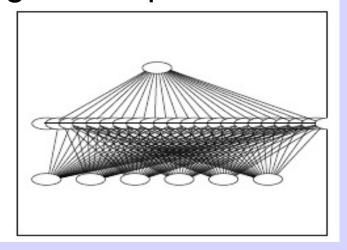
**u**1

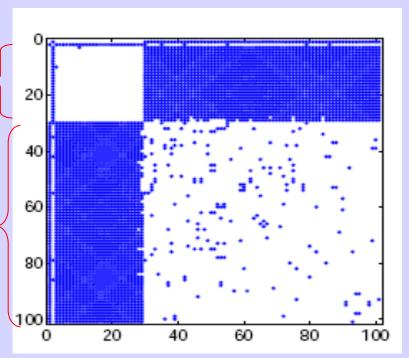
#### **Bipartite Communities!**

patents from same inventor(s)

`cut-and-paste' bibliography!

magnified bipartite community





15-826



#### Roadmap

- Introduction Motivation
- Unsupervised
  - static graphs



- Eigenspokes, and extensions
- Spot strange, individual nodes: oddBall
- Time-evolving graphs
- Supervised: BP
- Conclusions



### Inferring Strange Behavior from Connectivity Pattern in Social Networks



Meng Jiang, Peng Cui, Shiqiang Yang (Tsinghua, Beijing)

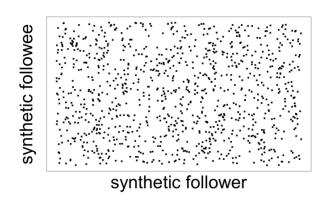
Alex Beutel, Christos Faloutsos (CMU)



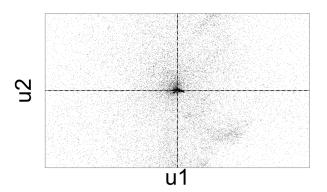


- Case #0: No lockstep behavior in random power law graph of 1M nodes, 3M edges
- Random ← "Scatter"

#### Adjacency Matrix



#### **Spectral Subspace Plot**

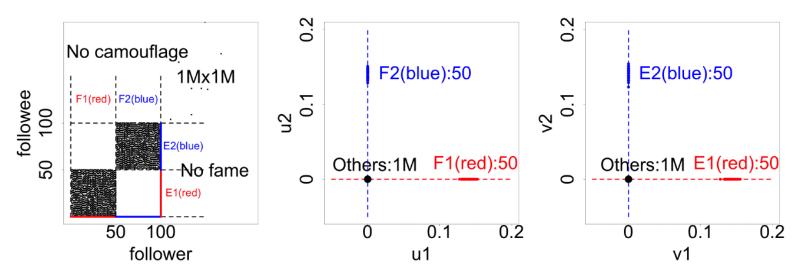




- Case #1: non-overlapping lockstep
- "Blocks" ← "Rays"

#### **Adjacency Matrix**

#### **Spectral Subspace Plot**



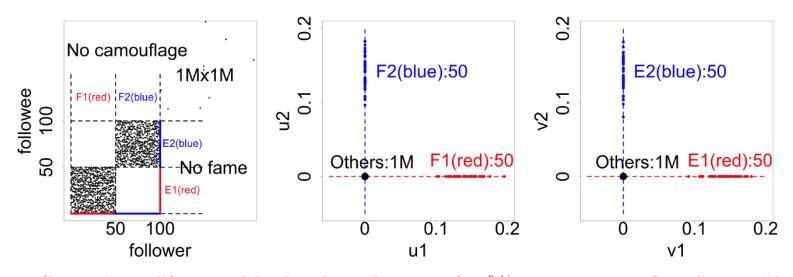
Rule 1 (short "rays"): two blocks, high density (90%), no "camouflage", no "fame" 15-826 Copyright: C. Faloutsos (2025)



- Case #2: non-overlapping lockstep
- "Blocks; low density" ← Elongation

#### **Adjacency Matrix**

#### **Spectral Subspace Plot**



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

15-826 Copyright: C. Faloutsos (2025)

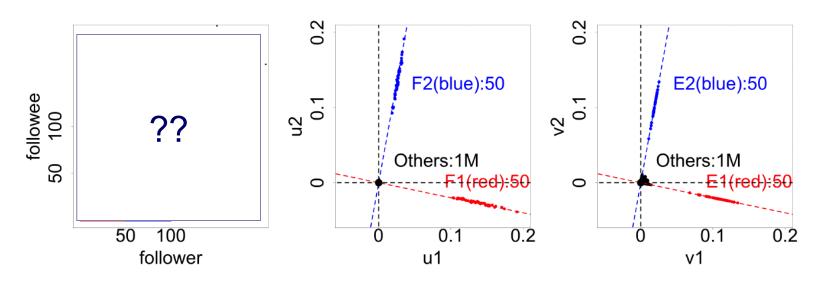
18



• Case #3: non-overlapping lockstep

#### **Adjacency Matrix**

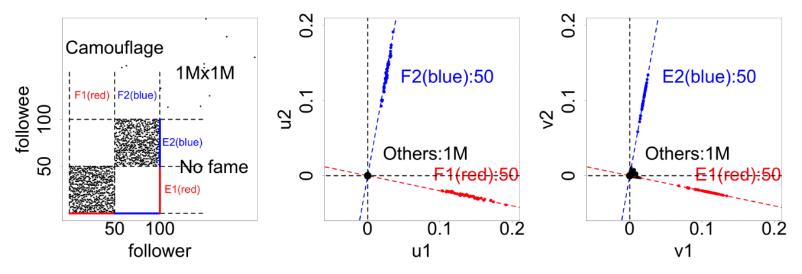
#### Spectral Subspace Plot



15-826 Copyright: C. Faloutsos (2025)



- Case #3: non-overlapping lockstep
- "Camouflage" (or "Fame") ← Tilting
   "Rays"
   Adjacency Matrix
   Spectral Subspace Plot



Rule 3 (tilting "rays"): two blocks, with "camouflage", no "fame"
15-826 Copyright: C. Faloutsos (2025)

20



• Case #4:

? lockstep

• "?"

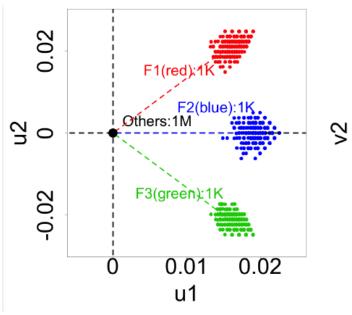


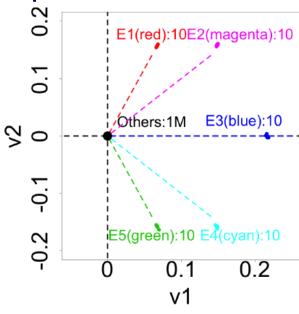
"Pearls"

#### **Adjacency Matrix**

#### **Spectral Subspace Plot**





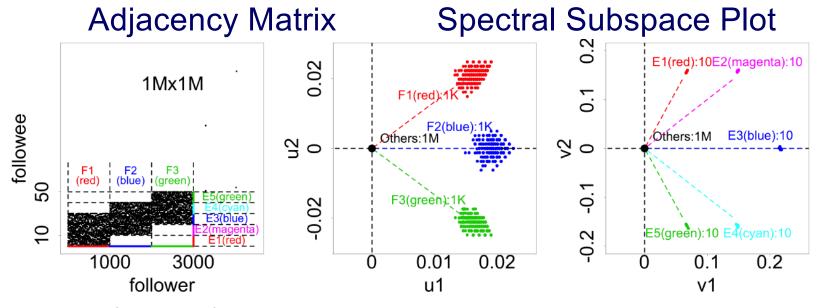


Copyright: C. Faloutsos (2025)



Case #4: overlapping lockstep

• "Staircase" "Pearls"



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

15-826



#### Roadmap

- Introduction Motivation
- Unsupervised
  - static graphs
    - Eigenspokes, and extensions
    - Spot strange, individual nodes: oddBall
  - Time-evolving graphs
- Supervised: BP
- Conclusions





# OddBall: Spotting Anomalies in Weighted Graphs



Carnegie Mellon University School of Computer Science

PAKDD 2010, Hyderabad, India

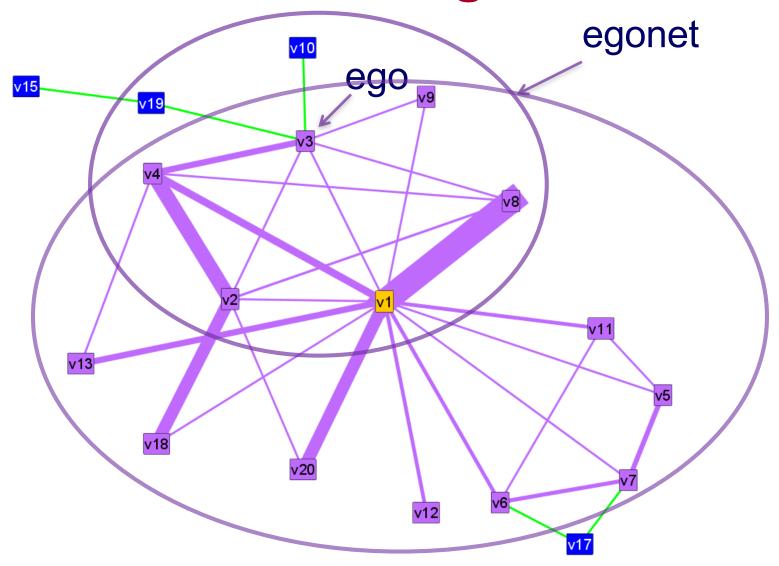
#### Main idea

For each node,

- extract 'ego-net' (=1-step-away neighbors)
- Extract features (#edges, total weight, etc etc)
- Compare with the rest of the population

#### Carnegie Mellon

### What is an egonet?

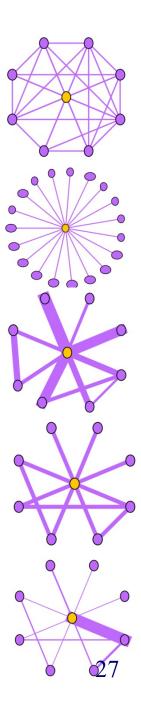


15-826

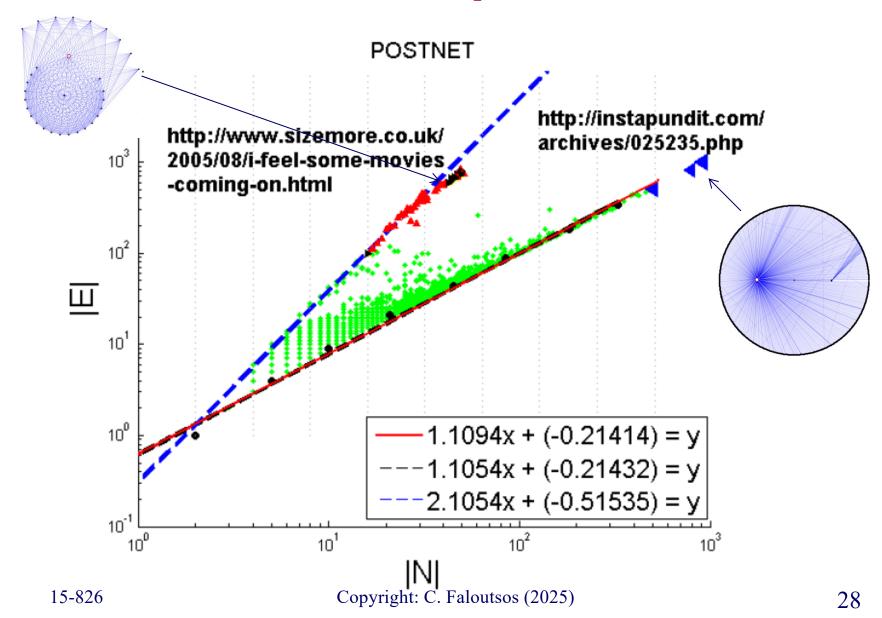


#### **Selected Features**

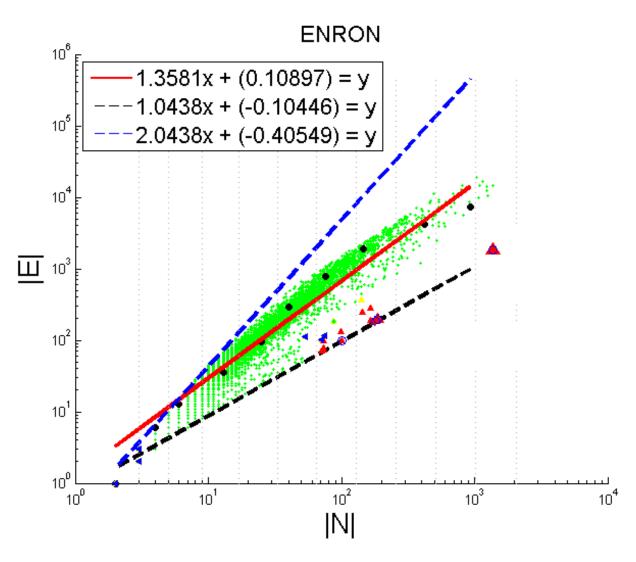
- $N_i$ : number of neighbors (degree) of ego i
- $E_i$ : number of edges in egonet i
- W<sub>i</sub>: total weight of egonet i
- $\lambda_{w,i}$ : principal eigenvalue of the weighted adjacency matrix of egonet I





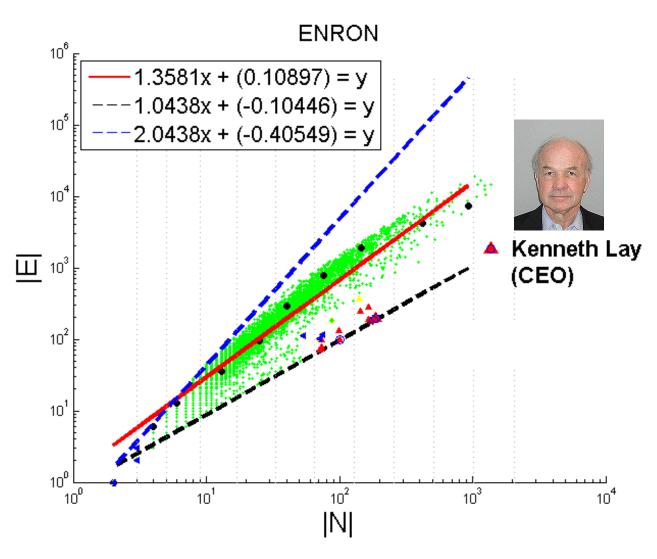




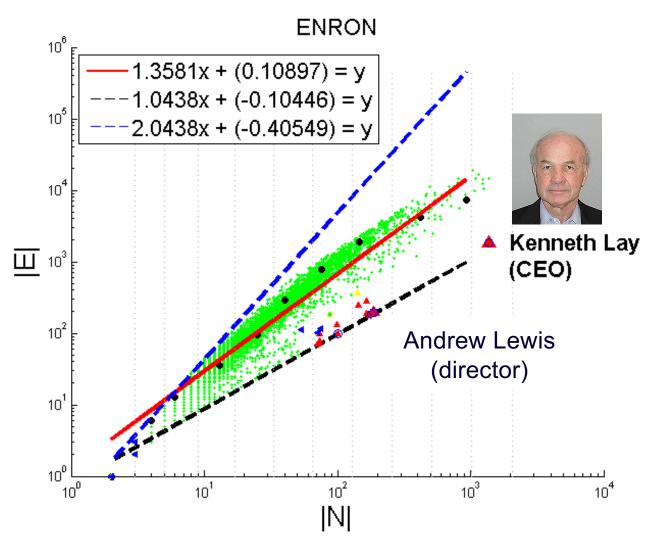


15-826

Copyright: C. Faloutsos (2025)



15-826 Copyright: C. Faloutsos (2025)



15-826

Copyright: C. Faloutsos (2025)



#### Roadmap

- Introduction Motivation
- Unsupervised
  - static graphs
  - Time-evolving graphs
    - 'copyCatch' in FaceBook
    - Patterns of IAT
    - [tensors]
- Supervised: BP
- Conclusions



#### Fraud

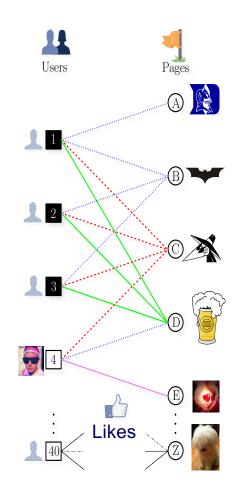
- Given
  - Who 'likes' what page, and when
- Find
  - Suspicious users and suspicious products



CopyCatch: Stopping Group Attacks by Spotting Lockstep Behavior in Social Networks, Alex Beutel, Wanhong Xu, Venkatesan Guruswami, Christopher Palow, Christos Faloutsos *WWW*, 2013.

#### Fraud

- Given
  - Who 'likes' what page, and when
- Find
  - Suspicious users and suspicious products



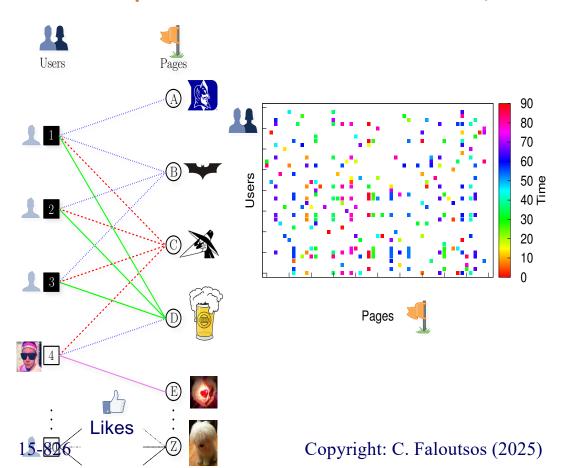
CopyCatch: Stopping Group Attacks by Spotting Lockstep Behavior in Social Networks, Alex Beutel, Wanhong Xu, Venkatesan Guruswami, Christopher Palow, Copyright: C. Faloutsos (2025)

Christos Faloutsos WWW, 2013.



## **Graph Patterns and Lockstep**Our intuition **Behavior**

Lockstep behavior: Same Likes, same time

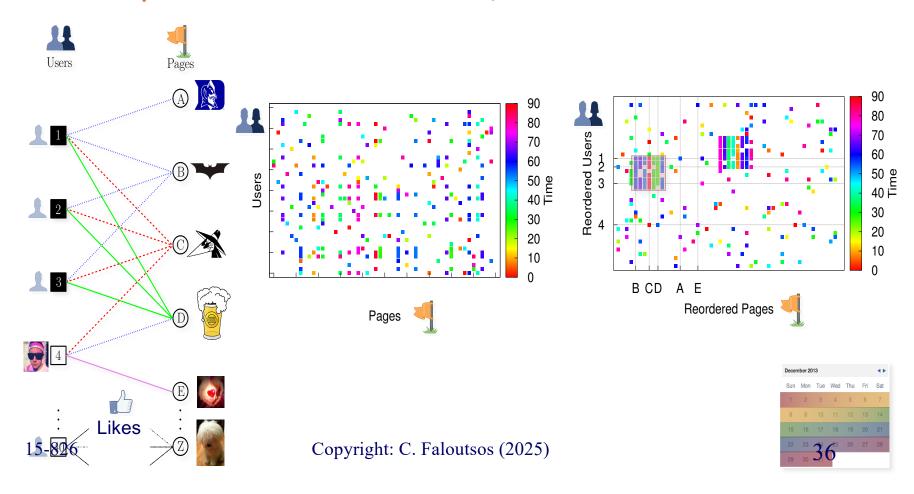






## Graph Patterns and Lockstep Our intuition Behavior

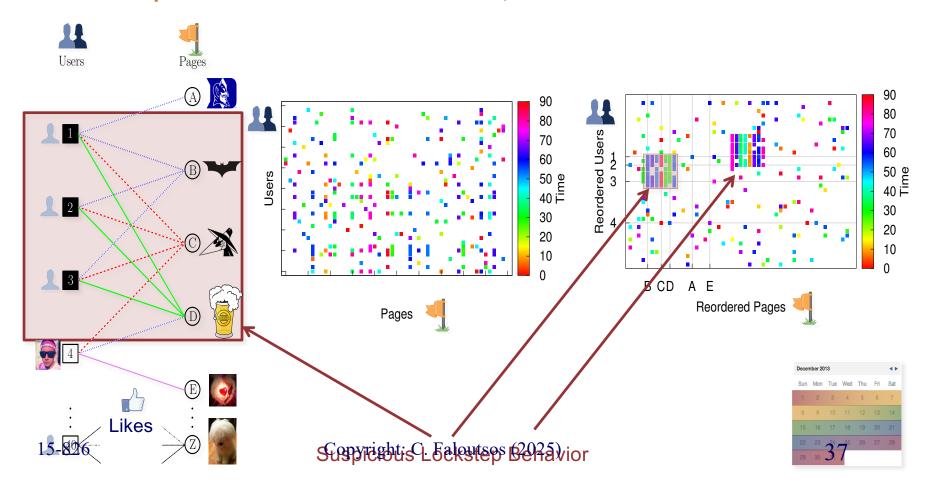
Lockstep behavior: Same Likes, same time





## Graph Patterns and Lockstep Our intuition Behavior

Lockstep behavior: Same Likes, same time

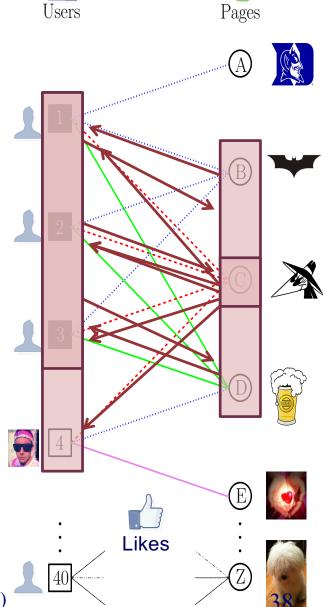




#### **MapReduce Overview**

Use Hadoop to search for many clusters in parallel:

- Start with randomly seed
- 2. Update set of Pages and center Like times for each cluster
- 3. Repeat until convergence





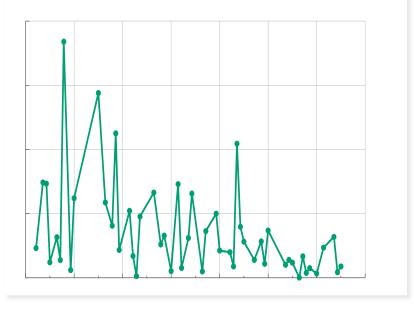


#### **Deployment at Facebook**

CopyCatch runs regularly (along with many other security mechanisms, and a large Site Integrity team)

3 months of CopyCatch @ Facebook

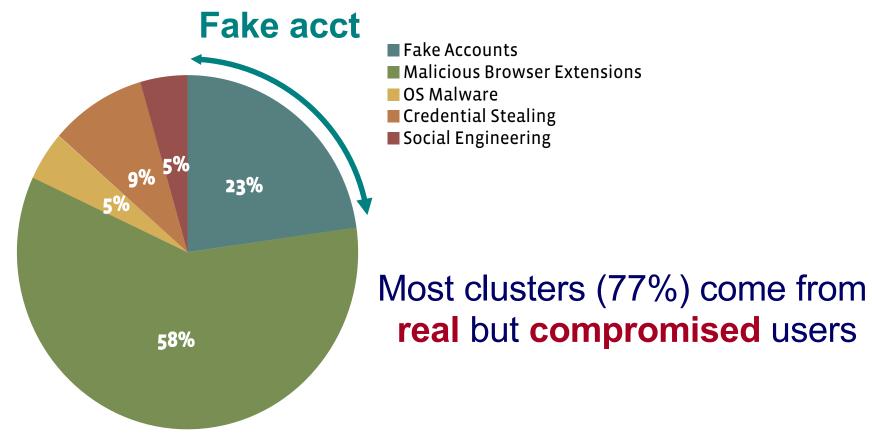
#users caught







#### **Deployment at Facebook**



Manually labeled 22 randomly selected clusters from February 2013



#### Roadmap

- Introduction Motivation
- Unsupervised
  - static graphs
  - Time-evolving graphs
    - 'copyCatch' in FaceBook
    - Patterns of IAT
    - [tensors]
- Supervised: BP
- 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

<sup>\*</sup>alceufc@icmc.usp.br

#### **Pattern Mining: Datasets**

#### Reddit Dataset

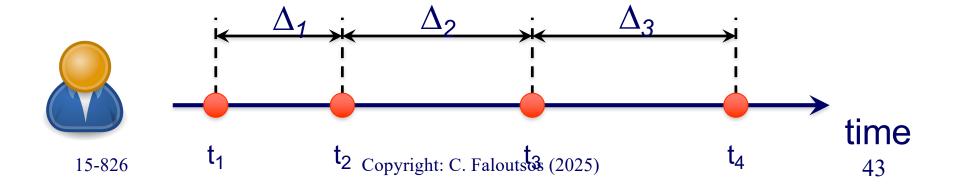
Time-stamp from comments 21,198 users 20 Million time-stamps

#### Twitter Dataset

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

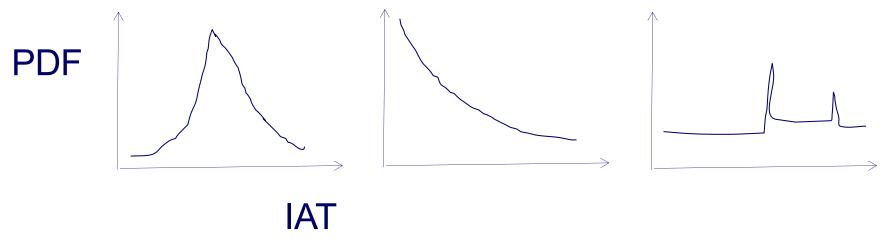
#### For each user we have:

Sequence of postings time-stamps:  $T = (t_1, t_2, t_3, ...)$ Inter-arrival times (IAT) of postings:  $(\Delta_1, \Delta_2, \Delta_3, ...)$ 



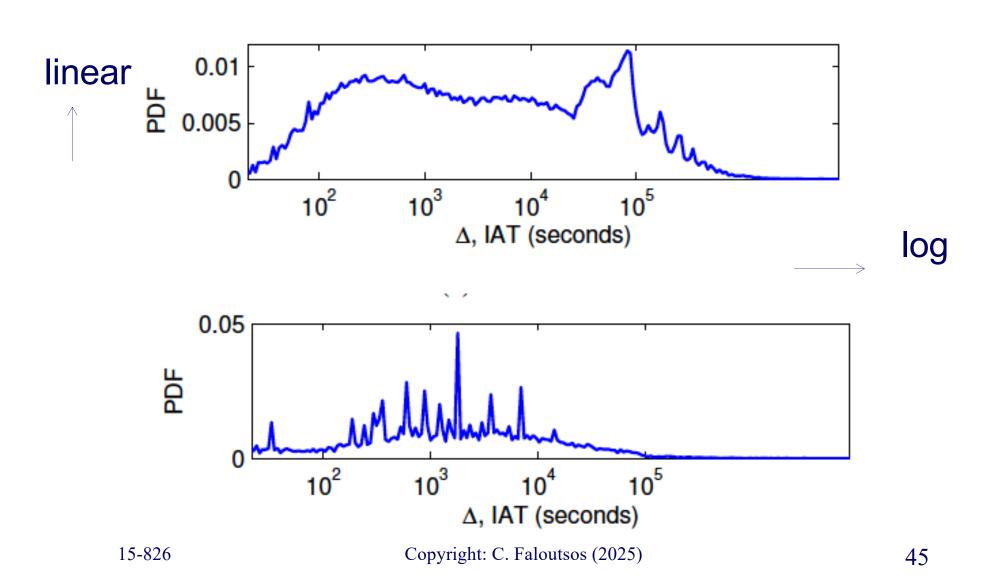
#### Guess the IAT pdf

- Gaussian (10' +/- a bit)?
- Power law? Slope?
- Log-logistic?
- Periodic (spike for 1day, 1 week)?

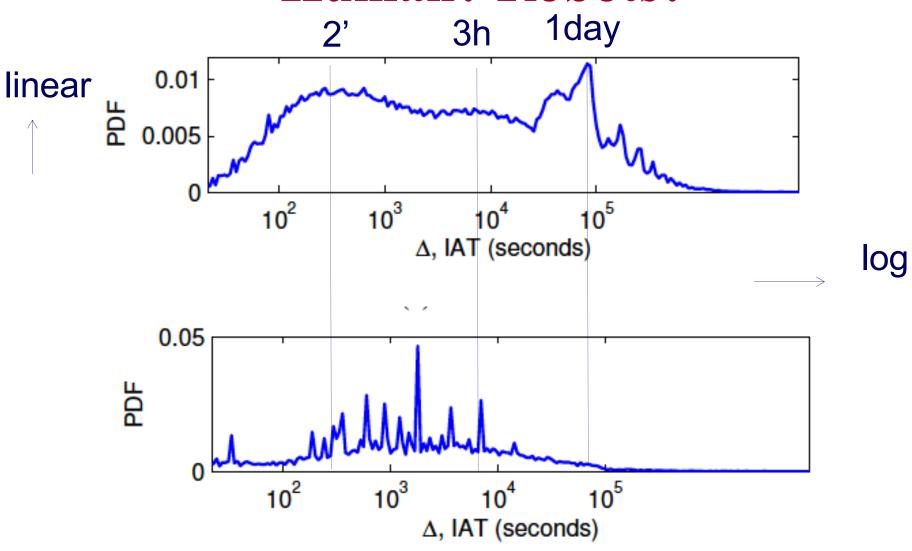


15-826

#### **Human? Robots?**



#### **Human? Robots?**



15-826

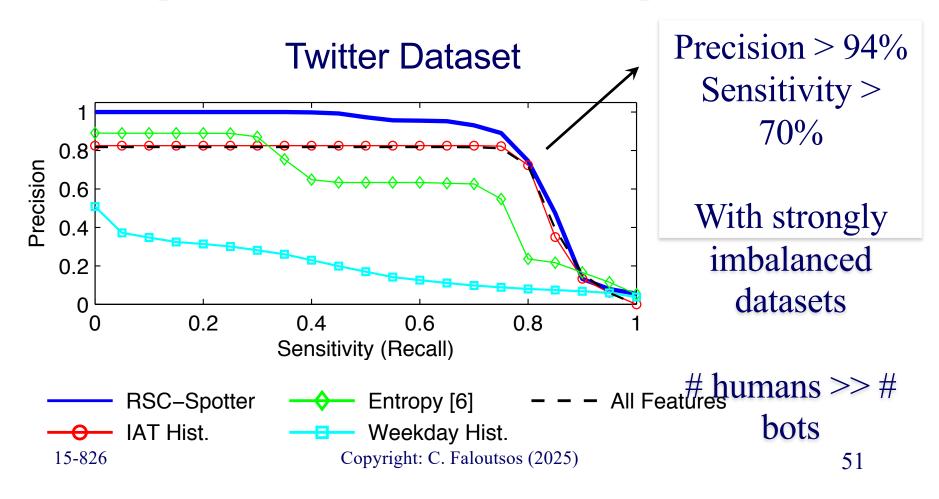
Copyright: C. Faloutsos (2025)

46

## **Experiments: Can RSC-Spotter Detect Bots?**

Precision vs. Sensitivity Curves

Good performance: curve close to the top

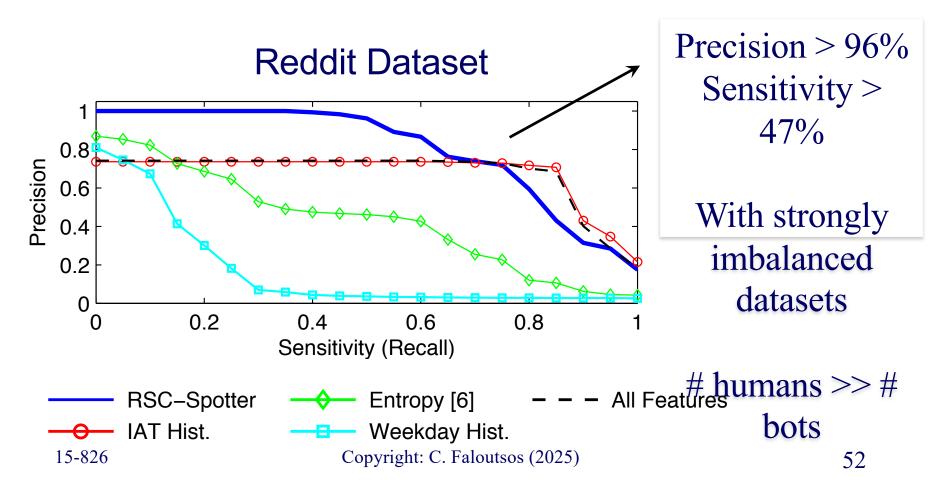




## **Experiments: Can RSC-Spotter Detect Bots?**

Precision vs. Sensitivity Curves

Good performance: curve close to the top





#### Roadmap

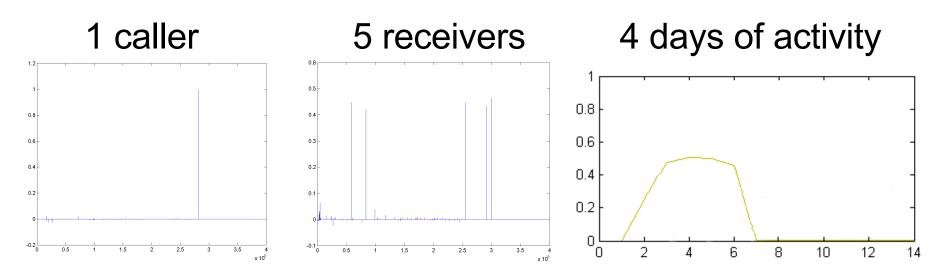
- Introduction Motivation
- Unsupervised
  - static graphs
  - Time-evolving graphs
    - 'copyCatch' in FaceBook
    - Patterns of IAT
    - [tensors]
- Supervised: BP
- Conclusions





## Anomaly detection in timeevolving graphs

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



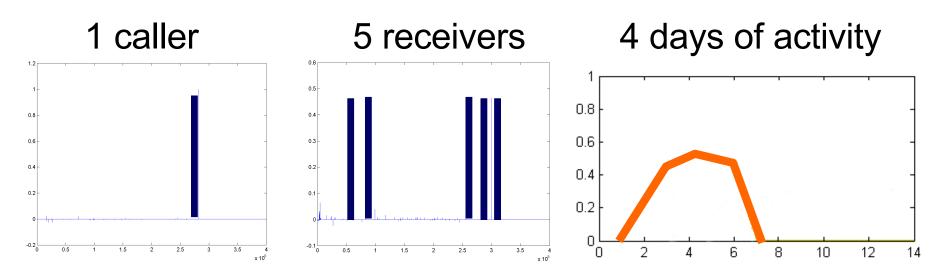
~200 calls to EACH receiver on EACH day!

15-826



## Anomaly detection in timeevolving graphs

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



~200 calls to EACH receiver on EACH day!

15-826



## Anomaly detection in timeevolving 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
- Unsupervised
  - static graphs
  - Time-evolving graphs
- Supervised: BP



- Ebay fraud ('NetProbe')
- Malware (Polonium)
- Conclusions

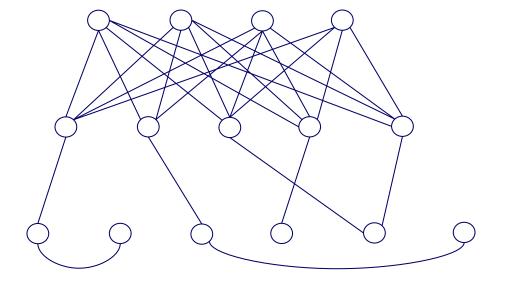


#### E-bay Fraud detection

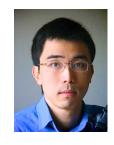


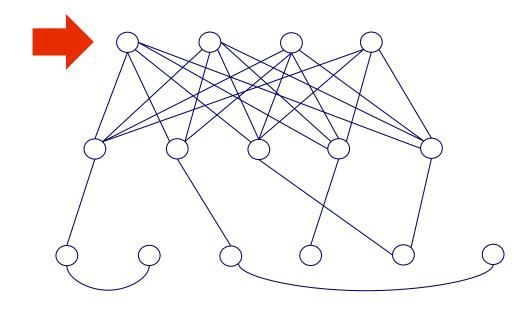


w/ Polo Chau & Shashank Pandit, CMU [www'07]

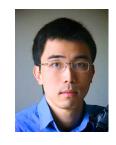


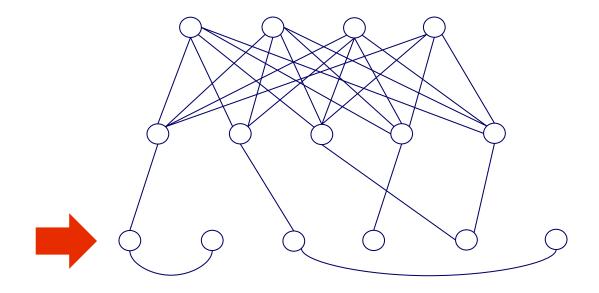








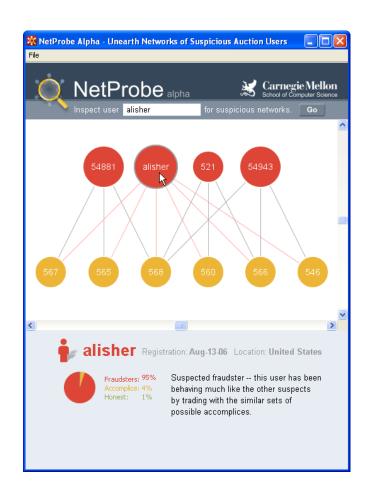


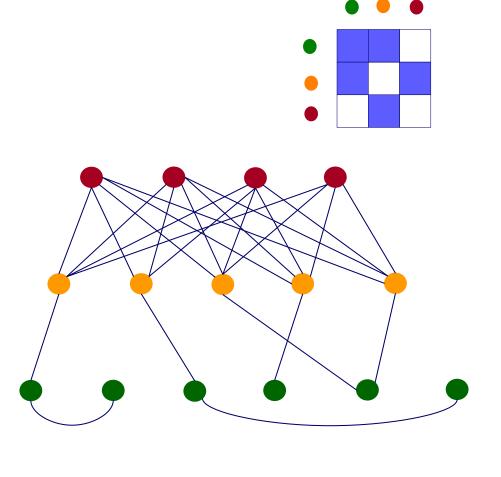


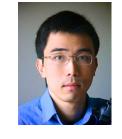












#### Popular press



And less desirable attention:

• E-mail from 'Belgium police' ('copy of your code?')



#### Roadmap

- Introduction Motivation
- Unsupervised
- Supervised: BP (Belief Propagation)
  - Ebay fraud ('NetProbe')
  - Malware (Polonium)



- Intuition behind BP
- Conclusions



## Unifying Guilt-by-Association Approaches: Theorems and Fast Algorithms



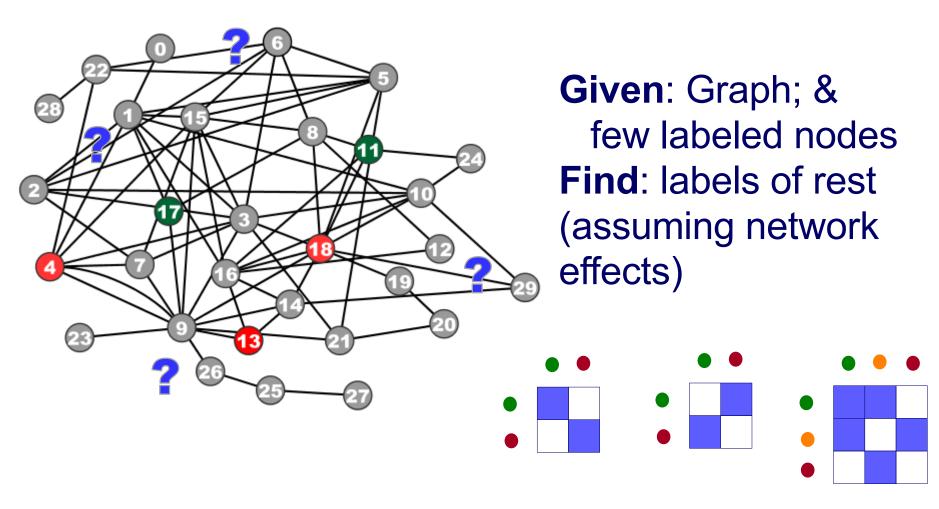
Danai Koutra
U Kang
Hsing-Kuo Kenneth Pao

Tai-You Ke Duen Horng (Polo) Chau Christos Faloutsos

ECML PKDD, 5-9 September 2011, Athens, Greece



## Problem Definition: GBA techniques



15-826

#### Are they related?

- RWR (Random Walk with Restarts)
  - google's pageRank ('if my friends are important, I'm important, too')
- SSL (Semi-supervised learning)
  - minimize the differences among neighbors
- BP (Belief propagation)
  - send messages to neighbors, on what you believe about them



### Are they related? YES!

- RWR (Random Walk with Restarts)
  - google's pageRank ('if my friends are important, I'm important, too')
- SSL (Semi-supervised learning)
  - minimize the differences among neighbors
- BP (Belief propagation)
  - send messages to neighbors, on what you believe about them



#### **Correspondence of Methods**

Method	Matrix	Unknow n		known
RWR	$[\mathbf{I} - \mathbf{c}  \mathbf{A}\mathbf{D}^{-1}]  \times$	X	=	(1-c) <b>y</b>
SSL	$[I + a(D - \overline{A})] \times$	X	=	$\mathbf{y}$
FABP	$[I + a D - c'A] \times$	$\mathbf{b_h}$	=	$\phi_h$
	dl d2 d3 d3 d1 0 1 0 adjacency matrix	? final labels/ beliefs		prior labels/beliefs

#### Cast



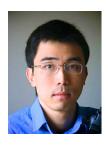
Akoglu, Leman



Araujo, Miguel



Beutel, Alex



Chau, Polo



Kang, U



Koutra, Danai



Lee, Jay Yoon



Prakash, Aditya



Papalexakis, Vagelis



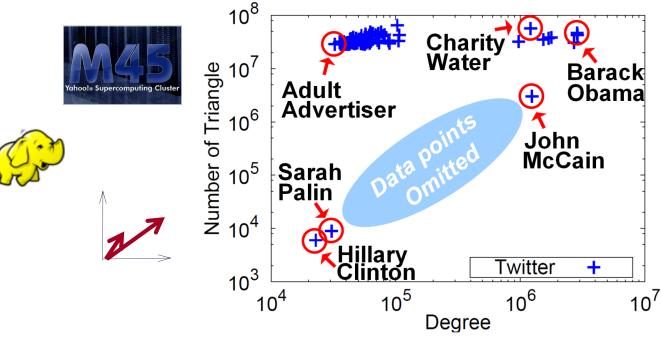
Shah, Neil



#### **Insights:**



• Large datasets reveal patterns/outliers that are invisible otherwise

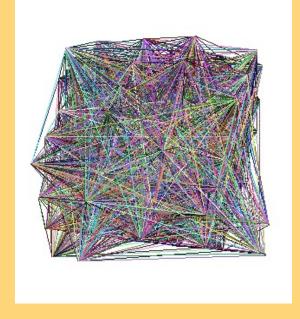


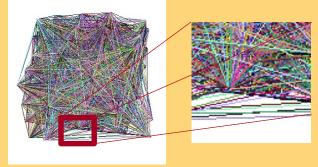
15-826

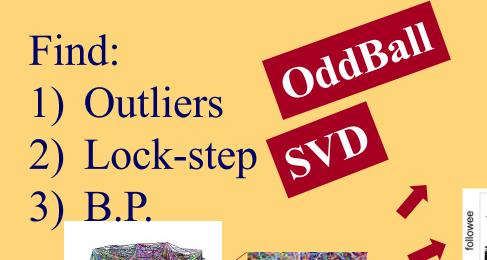


#### **Solutions**

#### Given:

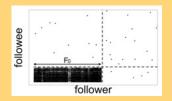








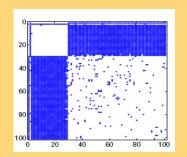




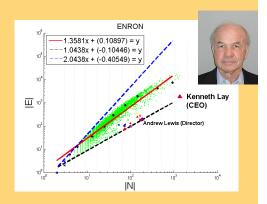
# un-supervised

#### Solutions, cont'd

Spectral/tensor methods: spot lockstep behavior



OddBall: strange, single nodes



semi-supervised

BP: good for the supervised setting (ie., some labels are given)

