**Carnegie Mellon**

# 15-826: Multimedia Databases and Data Mining

Lecture #26: Graph mining - patterns

*Christos Faloutsos*

---

**Carnegie Mellon**

## Must-read Material

- Michalis Faloutsos, Petros Faloutsos and Christos Faloutsos, On Power-Law Relationships of the Internet Topology, SIGCOMM 1999.
- R. Albert, H. Jeong, and A.-L. Barabasi, Diameter of the World Wide Web Nature, 401, 130-131 (1999).
- Reka Albert and Albert-Laszlo Barabasi Statistical mechanics of complex networks, Reviews of Modern Physics, 74, 47 (2002).
- Jure Leskovec, Jon Kleinberg, Christos Faloutsos Graphs over Time: Densification Laws, Shrinking Diameters and Possible Explanations, KDD 2005, Chicago, IL, USA

---

**Carnegie Mellon**

## Must-read Material (cont'd)

- D. Chakrabarti and C. Faloutsos, Graph Mining: Laws, Generators and Algorithms, in ACM Computing Surveys, 38 (1), 2006
- J. Leskovec, D. Chakrabarti, J. Kleinberg, and C. Faloutsos, Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication, in PKDD 2005, Porto, Portugal

**Carnegie Mellon**

## Outline

➡ • Introduction – Motivation
• Problem#1: Patterns in graphs
• Problem#2: Tools
• Problem#3: Scalability
• Conclusions

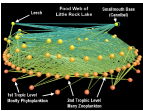15-826                    (c) 2011  C. Faloutsos                    4

---

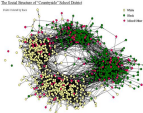**Carnegie Mellon**

## Graphs - why should we care?

**Linked in**
**f** **t**

Food Web
[Martinez '91]

Friendship Network
[Moody '01]

Internet Map
[lumeta.com]

15-826                    (c) 2011  C. Faloutsos                    5

---

**Carnegie Mellon**

## Graphs - why should we care?

• IR: bi-partite graphs (doc-terms)

$D_1$ ... $T_1$
$D_N$ ... $T_M$

• web: hyper-text graph

• ... and more:

15-826                    (c) 2011  C. Faloutsos                    6

Faloutsos

## Graphs - why should we care?

- 'viral' marketing
- web-log ('blog') news propagation
- computer network security: email/IP traffic and anomaly detection
- ....

## Outline

- Introduction – Motivation
- ➡ Problem#1: Patterns in graphs
  - Static graphs
  - Weighted graphs
  - Time evolving graphs
- Problem#2: Tools
- Problem#3: Scalability
- Conclusions

## Problem #1 - network and graph mining

- What does the Internet look like?
- What does FaceBook look like?

- What is 'normal'/'abnormal'?
- which patterns/laws hold?

**Carnegie Mellon**

# Problem #1 - network and graph mining



- What does the Internet look like?
- What does FaceBook look like?

- What is 'normal'/'abnormal'?
- which patterns/laws hold?
  – To spot **anomalies** (rarities), we have to discover **patterns**
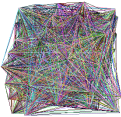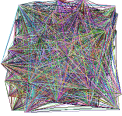
15-826       (c) 2011  C. Faloutsos      10

---

**Carnegie Mellon**

# Problem #1 - network and graph mining



- What does the Internet look like?
- What does FaceBook look like?

- What is 'normal'/'abnormal'?
- which patterns/laws hold?
  – To spot **anomalies** (rarities), we have to discover **patterns**
  – **Large** datasets reveal patterns/anomalies that may be invisible otherwise…

15-826       (c) 2011  C. Faloutsos      11

---

**Carnegie Mellon**

# Are real graphs random?

- random (Erdos-Renyi) graph – 100 nodes, avg degree = 2
- before layout
- after layout
- No obvious patterns

(generated with: pajek
http://vlado.fmf.uni-lj.si/pub/networks/pajek/ )



15-826       (c) 2011  C. Faloutsos      12

**Carnegie Mellon**

# Graph mining

- Are real graphs random?

15-826      (c) 2011 C. Faloutsos      13

---

**Carnegie Mellon**

# Laws and patterns

- Are real graphs random?
- A: NO!!
  - Diameter
  - in- and out- degree distributions
  - other (surprising) patterns

- So, let's look at the data

15-826      (c) 2011 C. Faloutsos      14

---

**Carnegie Mellon**

# Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

**internet domains**



15-826      (c) 2011 C. Faloutsos      15

# Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

**internet domains**

att.com

log(degree)

ibm.com

**-0.82**

log(rank)

15-826     (c) 2011 C. Faloutsos     16

---

# Solution# S.2: Eigen Exponent $E$

Eigenvalue

Exponent = slope

$E = -0.48$

May 2001

Rank of decreasing eigenvalue

- A2: power law in the eigenvalues of the adjacency matrix

15-826     (c) 2011 C. Faloutsos     17

---

# Solution# S.2: Eigen Exponent $E$

Eigenvalue

Exponent = slope

$E = -0.48$

May 2001

Rank of decreasing eigenvalue

- [Mihail, Papadimitriou '02]: slope is ½ of rank exponent

15-826     (c) 2011 C. Faloutsos     18

### But:

How about graphs from other domains?

### More power laws:

- web hit counts [w/ A. Montgomery]



Web Site Traffic

Count (log scale)

Zipf

``ebay''

in-degree (log scale)

users

sites

### epinions.com

count

Original graph
R-MAT graph

- who-trusts-whom [Richardson + Domingos, KDD 2001]

trusts-2000-people user

(out) degree

---

**Carnegie Mellon**

## And numerous more

- # of sexual contacts
- Income [Pareto] –'80-20 distribution'
- Duration of downloads [Bestavros+]
- Duration of UNIX jobs ('mice and elephants')
- Size of files of a user
- …
- 'Black swans'

15-826          (c) 2011  C. Faloutsos          22

---

**Carnegie Mellon**

## Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
  - Static graphs
    - degree, diameter, eigen,
  ➡ - triangles
    - cliques
  - Weighted graphs
  - Time evolving graphs
- Problem#2: Tools

15-826          (c) 2011  C. Faloutsos          23

---

**Carnegie Mellon**

## Solution# S.3: Triangle 'Laws'

- Real social networks have a lot of triangles

15-826          (c) 2011  C. Faloutsos          24

---

**Carnegie Mellon**

# Solution# S.3: Triangle 'Laws'

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

15-826                    (c) 2011 C. Faloutsos                    25

---

**Carnegie Mellon**

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

HEP-TH

ASN

Epinions

X-axis: # of participating triangles
Y: count (~ pdf)

15-826                    Faloutsos                    26

---

**Carnegie Mellon**

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

HEP-TH

ASN

Epinions

X-axis: # of participating triangles
Y: count (~ pdf)

15-826                    Faloutsos                    27

## Triangle Law: #S.4
### [Tsourakakis ICDM 2008]

Reuters

SN

Epinions

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

15-826                    C. Faloutsos                    28

---

## Triangle Law: Computations
### [Tsourakakis ICDM 2008]

details

But: triangles are expensive to compute
(3-way join; several approx. algos)
Q: Can we do that quickly?

15-826                (c) 2011  C. Faloutsos                29

---

## Triangle Law: Computations
### [Tsourakakis ICDM 2008]

details

But: triangles are expensive to compute
(3-way join; several approx. algos)
Q: Can we do that quickly?
A: Yes!
   **#triangles = 1/6 Sum ( $\lambda_i^3$ )**
   (and, because of skewness (S2) ,
   we only need the top few eigenvalues!

15-826                (c) 2011  C. Faloutsos                30

**Carnegie Mellon**

details

# Triangle Law: Computations
## [Tsourakakis ICDM 2008]
Wikipedia graph 2006-Nov-04
≈ 3,1M nodes ≈ 37M edges



(1021x, 97.4%)

(1277x, 94.7%)

(1329x, 92.8%)

**1000x+ speed-up, >90% accuracy**

15-826                    (c) 2011 C. Faloutsos                    31

---

**Carnegie Mellon**

# Triangle counting for large graphs?

Anomalous nodes in Twitter(~ 3 billion edges)

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

15-826                    (c) 2011 C. Faloutsos                    32

---

**Carnegie Mellon**

# Triangle counting for large graphs?



Charity
Water

Adult
Advertiser

Barack
Obama

John
McCain

Sarah
Palin

Data points
Omitted

Hillary
Clinton

Twitter    +

Anomalous nodes in Twitter(~ 3 billion edges)

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

15-826                    (c) 2011 C. Faloutsos                    33

**Carnegie Mellon**

# Any other 'laws'?

Yes!

15-826                    (c) 2011 C. Faloutsos                    34

---

**Carnegie Mellon**

# Any other 'laws'?

Yes!

- Small diameter (~ constant!) –
  - six degrees of separation / 'Kevin Bacon'
  - small worlds [Watts and Strogatz]

15-826                    (c) 2011 C. Faloutsos                    35

---

**Carnegie Mellon**

# Any other 'laws'?

- Bow-tie, for the web [Kumar+ '99]
- IN, SCC, OUT, 'tendrils'
- disconnected components

15-826                    (c) 2011 C. Faloutsos                    36

## Any other 'laws'?

- power-laws in communities (bi-partite cores) [Kumar+, '99]

Log(count)

n:1

n:3     n:2

Log(m)

2:3 core
(m:n core)

## Any other 'laws'?

- "Jellyfish" for Internet [Tauro+ '01]
- core: ~clique
- ~5 concentric layers
- many 1-degree nodes

## EigenSpokes

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)

$$A = U\Sigma U^T$$

15-826       (c) 2011 C. Faloutsos      40

---

**EigenSpokes**

details

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

$$A = U\Sigma U^T$$

N

N

$\vec{u}_1 \ \vec{u}_i$

15-826       (c) 2011 C. Faloutsos      41
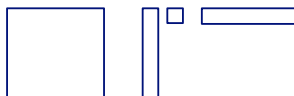
---

**EigenSpokes**

details

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

$$A = U\Sigma U^T$$

N

N

$\vec{u}_1 \ \vec{u}_i$

15-826       (c) 2011 C. Faloutsos      42

14

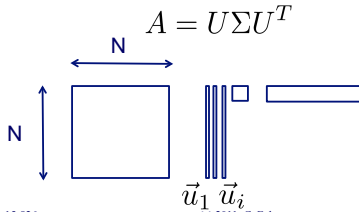### EigenSpokes

details

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

$$A = U\Sigma U^T$$

N

N

$\vec{u}_1 \ \vec{u}_i$

15-826

(c) 2011 C. Faloutsos

43
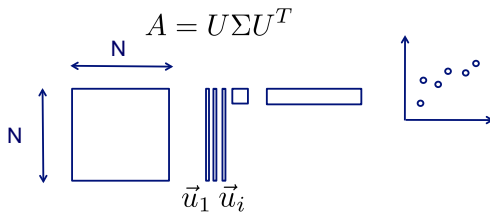
---

### EigenSpokes

details

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

$$A = U\Sigma U^T$$

N

N

$\vec{u}_1 \ \vec{u}_i$

15-826

(c) 2011 C. Faloutsos

44

---

### EigenSpokes

- EE plot:
- Scatter plot of scores of u1 vs u2
- One would expect
  - Many points @ origin
  - A few scattered ~randomly

2<sup>nd</sup> Principal component
u2

u1

1<sup>st</sup> Principal component

15-826

(c) 2011 C. Faloutsos

45

Faloutsos

---

**EigenSpokes**

- EE plot:
- Scatter plot of scores of u1 vs u2
- One would expect
  - Many points @ origin
  - A few scattered ~randomly

u2

90°

u1

15-826 (c) 2011 C. Faloutsos 46

---

**EigenSpokes - pervasiveness**
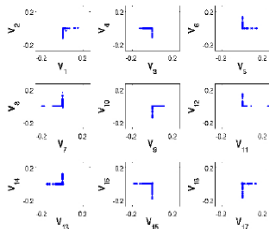
- Present in mobile social graph
  - across time and space
- Patent citation graph

15-826 (c) 2011 C. Faloutsos 47

---

**EigenSpokes - explanation**

Near-cliques, or near-bipartite-cores, loosely connected

15-826 (c) 2011 C. Faloutsos 48

**Carnegie Mellon**

# EigenSpokes - explanation

Near-cliques, or near-
bipartite-cores, loosely
connected

15-826       (c) 2011 C. Faloutsos       49

---

**Carnegie Mellon**

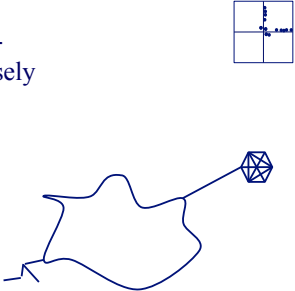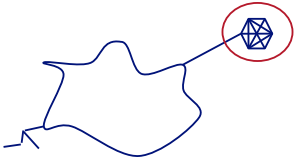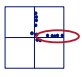# EigenSpokes - explanation

Near-cliques, or near-
bipartite-cores, loosely
connected

15-826       (c) 2011 C. Faloutsos       50

---

**Carnegie Mellon**

# EigenSpokes - explanation

Near-cliques, or near-
bipartite-cores, loosely
connected

spy plot of top 20 nodes

So what?

- Extract nodes with high s*cores*
- high connectivity
- Good "communities"

15-826       (c) 2011 C. Faloutsos       51

## Bipartite Communities!



patents from same inventor(s)

`cut-and-paste' bibliography!

magnified bipartite community

15-826      (c) 2011 C. Faloutsos      52

---

## Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
  - Static graphs
    - degree, diameter, eigen,
    - triangles
    - cliques
  - Weighted graphs
  - Time evolving graphs
- Problem#2: Tools

15-826      (c) 2011 C. Faloutsos      53

---

## Observations on weighted graphs?

- A: yes - even more 'laws'!



M. McGlohon, L. Akoglu, and C. Faloutsos
*Weighted Graphs and Disconnected Components: Patterns and a Generator.*
*SIG-KDD* 2008

15-826      (c) 2011 C. Faloutsos      54

**Carnegie Mellon**

# Observation W.1: Fortification

*Q: How do the weights
of nodes relate to degree?*

15-826       (c) 2011 C. Faloutsos       55

---

**Carnegie Mellon**

# Observation W.1: Fortification

**More donors,
more $ ?**

$10    'Reagan'

$5

$7    'Clinton'

15-826       (c) 2011 C. Faloutsos       56

---

**Carnegie Mellon**

# Observation W.1: fortification:
## Snapshot Power Law

- Weight: super-linear on in-degree
- exponent 'iw': 1.01 < iw < 1.26

**More donors,
<u>even</u> more $**

**Orgs-Candidates**

$10    In-weights ($)

e.g. John Kerry,
$10M received,
from 1K donors

$5

Edges (# donors)

15-826       (c) 2011 C. Faloutsos       57

19

## Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
  – Static graphs
  – Weighted graphs
  ➡ – Time evolving graphs
- Problem#2: Tools
- …

15-826      (c) 2011 C. Faloutsos      58

---

## Problem: Time evolution

- with Jure Leskovec (CMU -> Stanford)

- and Jon Kleinberg (Cornell – sabb. @ CMU)

15-826      (c) 2011 C. Faloutsos      59

---

## T.1 Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:
  – diameter ~ $O(\log N)$
  – diameter ~ $O(\log \log N)$
- What is happening in real data?

15-826      (c) 2011 C. Faloutsos      60

**Carnegie Mellon**

## T.1 Evolution of the Diameter

- Prior work on Power Law graphs hints
  at **slowly growing diameter**:
  - diameter ~ O(log N)
  - diameter ~ O(log log N)
- What is happening in real data?
- Diameter **shrinks** over time

15-826                    (c) 2011  C. Faloutsos                    61

---

**Carnegie Mellon**

## T.1 Diameter – "Patents"

- Patent citation
  network
- 25 years of data
- @1999
  - 2.9 M nodes
  - 16.5 M edges

diameter

time [years]

15-826                    (c) 2011  C. Faloutsos                    62

---

**Carnegie Mellon**

## T.2 Temporal Evolution of the Graphs

- N(t) … nodes at time t
- E(t) … edges at time t
- Suppose that
    N(t+1) = 2 * N(t)
- Q: what is your guess for
    E(t+1) =? 2 * E(t)

15-826                    (c) 2011  C. Faloutsos                    63

**Carnegie Mellon**

## T.2 Temporal Evolution of the Graphs

- N(t) … nodes at time t
- E(t) … edges at time t
- Suppose that
    N(t+1) = 2 * N(t)
- Q: what is your guess for
    E(t+1) = 2 * E(t)
- A: over-doubled!
    – But obeying the ``Densification Power Law''

15-826                    (c) 2011  C. Faloutsos                    64

---

**Carnegie Mellon**

## T.2 Densification – Patent Citations

- Citations among patents granted
- @1999
    – 2.9 M nodes
    – 16.5 M edges
- Each year is a datapoint

E(t)

$10^8$

1999

$10^7$

Number of edges

1.66

$10^6$

1975

$10^5$

Edges
= 0.0002 x$^{1.66}$ R$^2$=0.99

$10^5$        $10^6$        $10^7$
Number of nodes

N(t)

15-826                    (c) 2011  C. Faloutsos                    65

---

**Carnegie Mellon**

## Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
    – Static graphs
    – Weighted graphs
    – Time evolving graphs
- Problem#2: Tools
- …

15-826                    (c) 2011  C. Faloutsos                    66

**CarnegieMellon**

# More on Time-evolving graphs

M. McGlohon, L. Akoglu, and C. Faloutsos
*Weighted Graphs and Disconnected*
*Components: Patterns and a Generator.*
*SIG-KDD* 2008

15-826        (c) 2011 C. Faloutsos      67

---

**CarnegieMellon**

# [ Gelling Point ]

- Most real graphs display a gelling point
- After gelling point, they exhibit typical behavior. This is marked by a spike in diameter.



**IMDB**

t=1914

Diameter

Time

15-826        (c) 2011 C. Faloutsos      68

---

**CarnegieMellon**

# Observation T.3: NLCC behavior

*Q: How do NLCC's emerge and join with the GCC?*

(``NLCC'' = non-largest conn. components)
– Do they continue to grow in size?
– or do they shrink?
– or stabilize?



15-826        (c) 2011 C. Faloutsos      69

**Carnegie Mellon**

## Observation T.3: NLCC behavior

*Q: How do NLCC's emerge and join with the GCC?*

(``NLCC'' = non-largest conn. components)
– Do they continue to grow in size?
– or do they <u>shrink</u>?
– or stabilize?

15-826                    (c) 2011  C. Faloutsos                    70

---

**Carnegie Mellon**

## Observation T.3: NLCC behavior

*Q: How do NLCC's emerge and join with the GCC?*

(``NLCC'' = non-largest conn. components)
YES – Do they continue to grow in size?
YES – or do they shrink?
YES – or stabilize?

15-826                    (c) 2011  C. Faloutsos                    71

---

**Carnegie Mellon**

## Observation T.3: NLCC behavior

• After the gelling point, the GCC takes off, but NLCC's remain ~constant (actually, **oscillate**).

**IMDB**

CC size

Time-stamp

15-826                    (c) 2011  C. Faloutsos                    72

Faloutsos

25

**Timing for Blogs**

- with Mary McGlohon (CMU->Google)
- Jure Leskovec (CMU->Stanford)
- Natalie Glance (now at Google)
- Mat Hurst (now at MSR)

[SDM'07]

15-826 (c) 2011 C. Faloutsos 73

---

**T.4 : popularity over time**

# in links

lag: days after post

1  2  3

Post popularity drops-off – exponentially?

@t

@t + **lag**

15-826 (c) 2011 C. Faloutsos 74

---

**T.4 : popularity over time**

# in links
(**log**)

days after post
(**log**)

Post popularity drops-off – exponentially?
POWER LAW!
Exponent?

15-826 (c) 2011 C. Faloutsos 75

---

**Carnegie Mellon**

## T.4 : popularity over time

# in links
(**log**)



-1.6

days after post
(**log**)

Post popularity drops-off – exponentially?
POWER LAW!
Exponent? -1.6
• close to -1.5: Barabasi's stack model
• and like the zero-crossings of a random walk

15-826          (c) 2011  C. Faloutsos          76

---

**Carnegie Mellon**

## -1.5 slope

J. G. Oliveira & A.-L. Barabási Human Dynamics: The
   Correspondence Patterns of Darwin and Einstein.
   *Nature* **437,** 1251 (2005) . [PDF]



[1] Figure 1 | The correspondence patterns of Darwin and Einstein.     77

---

**Carnegie Mellon**

## T.5: duration of phonecalls

*Surprising Patterns for the Call
   Duration Distribution of Mobile
   Phone Users*

Pedro O. S. Vaz de Melo, Leman
   Akoglu, Christos Faloutsos, Antonio
   A. F. Loureiro

PKDD 2010

15-826          (c) 2011  C. Faloutsos          78

### Slide 79

**Carnegie Mellon**

# Probably, power law (?)



Plot: count vs Duration (s), with y-axis from $10^{-2}$ to $10^2$ and x-axis from $10^0$ to $10^3$. Marked "??"

15-826                    (c) 2011  C. Faloutsos                    79

### Slide 80

**Carnegie Mellon**

# No Power Law!



Plot: count vs Duration (s), legend: data, TLAC, log-normal, exponential

15-826                    (c) 2011  C. Faloutsos                    80

### Slide 81

**Carnegie Mellon**

# 'TLaC: Lazy Contractor'

- The longer a task (phonecall) has taken,
- The even longer it will take

Odds ratio=

*Casualties(<x): Survivors(>=x)*

== power law



Plot: count vs Duration (s), legend: data, TLAC, log-normal, exponential



Plot vs duration (s), y-axis $10^{-4}$ to $10^4$, legend: data, TLAC, log-normal, exponential

15-826                    (c) 2011  C. Faloutsos                    81

**Carnegie Mellon**

# Data Description

- Data from a private mobile operator of a large city
  - 4 months of data
  - 3.1 million users
  - more than 1 billion phone records
- Over 96% of 'talkative' users obeyed a TLAC distribution ('talkative': >30 calls)

15-826        (c) 2011 C. Faloutsos       82

---

**Carnegie Mellon**

# Outliers:



15-826        (c) 2011 C. Faloutsos       83

---

**Carnegie Mellon**

# Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
- Problem#2: Tools
  - OddBall (anomaly detection)
  - Belief Propagation
  - Immunization
- Problem#3: Scalability
- Conclusions

15-826        (c) 2011 C. Faloutsos       84

**CarnegieMellon**

## OddBall: Spotting Anomalies in Weighted Graphs

Leman Akoglu, Mary McGlohon, Christos Faloutsos

*Carnegie Mellon University*

*School of Computer Science*

PAKDD 2010, Hyderabad, India

---

**CarnegieMellon**

## 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

15-826        (c) 2011 C. Faloutsos       86

---

**CarnegieMellon**

## What is an egonet?



15-826        (c) 2011 C. Faloutsos       87

**Selected Features**

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

15-826 (c) 2011 C. Faloutsos



**Near-Clique/Star**

POSTNET

http://www.sizemore.co.uk/
2005/08/i-feel-some-movies
-coming-on.html

http://instapundit.com/
archives/025235.php

$1.1094x + (-0.21414) = y$
$1.1054x + (-0.21432) = y$
$2.1054x + (-0.51535) = y$

15-826 (c) 2011 C. Faloutsos 89



**Near-Clique/Star**

ENRON

$1.3581x + (0.10897) = y$
$1.0438x + (-0.10446) = y$
$2.0438x + (-0.40549) = y$

15-826 (c) 2011 C. Faloutsos 90

CarnegieMellon

## Near-Clique/Star

ENRON

▲ Kenneth Lay (CEO)

15-826        (c) 2011 C. Faloutsos        91

CarnegieMellon

## Near-Clique/Star

ENRON

▲ Kenneth Lay (CEO)

Andrew Lewis (director)

15-826        (c) 2011 C. Faloutsos        92

CarnegieMellon

## Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
- Problem#2: Tools
  - OddBall (anomaly detection)
  - Belief Propagation
  - Immunization
- Problem#3: Scalability
- Conclusions

15-826        (c) 2011 C. Faloutsos        93

CarnegieMellon

# E-bay Fraud detection

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

15-826     (c) 2011 C. Faloutsos     94

CarnegieMellon

# E-bay Fraud detection

15-826     (c) 2011 C. Faloutsos     95

CarnegieMellon

# E-bay Fraud detection

15-826     (c) 2011 C. Faloutsos     96

### Slide 97

**CarnegieMellon**

## E-bay Fraud detection - NetProbe



15-826                    (c) 2011  C. Faloutsos                    97

### Slide 98

**CarnegieMellon**

## Popular press

USA TODAY

The Washington Post

Los Angeles Times

And less desirable attention:
• E-mail from 'Belgium police' ('copy of your code?')

15-826                    (c) 2011  C. Faloutsos                    98

### Slide 99

**CarnegieMellon**

## Outline

• Introduction – Motivation
• Problem#1: Patterns in graphs
• Problem#2: Tools
  – OddBall (anomaly detection)
  ➡ – Belief Propagation – antivirus app
  – Immunization
• Problem#3: Scalability
• Conclusions

15-826                    (c) 2011  C. Faloutsos                    99

**Carnegie Mellon**

## Polonium: Tera-Scale Graph Mining and Inference for Malware Detection

*SDM 2011, Mesa, Arizona*

**Polo Chau**
Machine Learning Dept

**Carey Nachenberg**
Vice President & Fellow

**Jeffrey Wilhelm**
Principal Software Engineer

**Adam Wright**
Software Engineer

**Prof. Christos Faloutsos**
Computer Science Dept

---

**Carnegie Mellon**

## Polonium: The Data

60+ terabytes of data *anonymously* contributed by participants of worldwide *Norton Community Watch* program

50+ million machines

900+ million executable files

Constructed a machine-file bipartite graph (0.2 TB+)

1 billion nodes (machines and files)

37 billion edges

15-826        (c) 2011  C. Faloutsos        101

---

**Carnegie Mellon**

## Polonium: Key Ideas

- Use Belief Propagation to propagate domain knowledge in machine-file graph to detect malware
- Use "guilt-by-association" (i.e., homophily)
  - E.g., files that appear on machines with many bad files are more likely to be bad
- Scalability: handles 37 billion-edge graph

15-826        (c) 2011  C. Faloutsos        102

## Polonium: One-Interaction Results

**Ideal**

**84.9%** True Positive Rate
**1%** False Positive Rate

**True Positive Rate**
% of malware
correctly identified

1.0

.8

.6

.4

.2

0

0    .2    .4    .6    .8    1.0

**False Positive Rate**
% of non-malware wrongly labeled as malware

15-826

---

## Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
- Problem#2: Tools
  - OddBall (anomaly detection)
  - Belief propagation
  - ➡ Immunization
- Problem#3: Scalability -PEGASUS
- Conclusions

15-826                    (c) 2011  C. Faloutsos                    104

---

## Immunization and epidemic thresholds

- Q1: which nodes to immunize?
- Q2: will a virus vanish, or will it create an epidemic?

15-826                    (c) 2011  C. Faloutsos                    105

**Carnegie Mellon**

# Q1: Immunization:

- Given
  - a network,
  - k vaccines, and
  - the virus details
- Which nodes to immunize?

15-826          (c) 2011  C. Faloutsos          106

**Carnegie Mellon**

# Q1: Immunization:

- Given
  - a network,
  - k vaccines, and
  - the virus details
- Which nodes to immunize?

15-826          (c) 2011  C. Faloutsos          107

**Carnegie Mellon**

# Q1: Immunization:

- Given
  - a network,
  - k vaccines, and
  - the virus details
- Which nodes to immunize?

15-826          (c) 2011  C. Faloutsos          108

**Q1: Immunization:**

•Given
  •a network,
  •k vaccines, and
  •the virus details
•Which nodes to immunize?

A: immunize the ones that
maximally raise
the `epidemic threshold'
[Tong+, ICDM'10]

15-826        (c) 2011  C. Faloutsos        109



details

**Q2: will a virus take over?**

• Flu-like virus (no immunity, 'SIS')
• Mumps (life-time immunity, 'SIR')
• Pertussis (finite-length immunity, 'SIRS')

$\beta$: attack prob
$\delta$: heal prob

15-826        (c) 2011  C. Faloutsos        110



details

**Q2: will a virus take over?**

• Flu-like virus (no immunity, 'SIS')
• Mumps (life-time immunity, 'SIR')
• Pertussis (finite-length immunity, 'SIRS')

$\beta$: attack prob
$\delta$: heal prob

A: depends on connectivity
   (avg degree? Max degree?
   variance?  Something else?

15-826        (c) 2011  C. Faloutsos        111

**Carnegie Mellon**

details

# Epidemic threshold $\tau$

What should $\tau$ depend on?
- avg. degree? and/or highest degree?
- and/or variance of degree?
- and/or third moment of degree?
- and/or diameter?

15-826         (c) 2011 C. Faloutsos         112

---

**Carnegie Mellon**

details

# Epidemic threshold

- [Theorem] We have no epidemic, if

$$\beta/\delta < \tau = 1/\lambda_{1,A}$$

15-826         (c) 2011 C. Faloutsos         113

---

**Carnegie Mellon**

details

# Epidemic threshold

- [Theorem] We have no epidemic, if

recovery prob.     epidemic threshold

$$\beta/\delta < \tau = 1/\lambda_{1,A}$$

attack prob.       largest eigenvalue of adj. matrix $A$

Proof: [Wang+03] (for SIS=flu only)

15-826         (c) 2011 C. Faloutsos         114

**Carnegie Mellon**

details

## A2: will a virus take over?

- For **all** typical virus propagation models (flu, mumps, pertussis, HIV, etc)
- The **only** connectivity measure that matters, is

$1/\lambda_1$

the first eigenvalue of the
adj. matrix
[Prakash+, '10, arxiv]

15-826      (c) 2011 C. Faloutsos      115

---

**Carnegie Mellon**

## Thresholds for some models

- $s$ = effective strength
- $s < 1$ : below threshold

| Models | Effective Strength (s) | Threshold (tipping point) |
|---|---|---|
| SIS, SIR, SIRS, SEIR | $s = \lambda \cdot \left(\frac{\beta}{\delta}\right)$ | |
| SIV, SEIV | $s = \lambda \cdot \left(\frac{\beta\gamma}{\delta(\gamma+\theta)}\right)$ | $s = 1$ |
| $SI_1I_2V_1V_2$ (**H.I.V.**) | $s = \lambda \cdot \left(\frac{\beta_1 v_2 + \beta_2 \varepsilon}{v_2(\varepsilon + v_1)}\right)$ | |

---

**Carnegie Mellon**

## A2: will a virus take over?

Fraction of infected

SIRS Infected (log-log)

under1
under2
over1
over2

Above: take-over

Below: exp. extinction

Graph:
Portland, OR
31M links
1.5M nodes

Time ticks

15-826      (c) 2011 C. Faloutsos      117

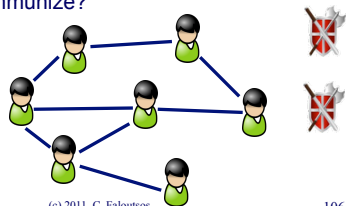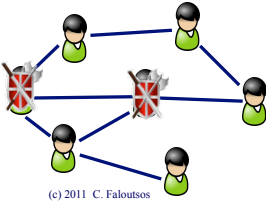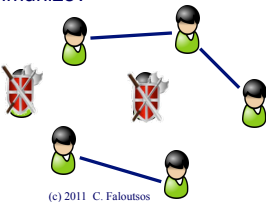**Carnegie Mellon**

## Q1: Immunization:

• Given
   • a network,
   • k vaccines, and
   • the virus details
• Which nodes to immunize?

A: immunize the ones that
   maximally raise
   the `epidemic threshold'
   [Tong+, ICDM'10]

15-826        (c) 2011 C. Faloutsos        118

---

**Carnegie Mellon**

## Q1: Immunization:

• Given
   • a network,
   • k vaccines, and
   • the virus details
• Which nodes to immunize?

A: immunize the ones that

**Max eigen-drop $\Delta\lambda$ for any virus!**

15-826        (c) 2011 C. Faloutsos        119

---

**Carnegie Mellon**

## Outline

• Introduction – Motivation
• Problem#1: Patterns in graphs
• Problem#2: Tools
   – OddBall (anomaly detection)
   – Belief propagation
   – Immunization
➡ • Problem#3: Scalability -PEGASUS
• Conclusions

15-826        (c) 2011 C. Faloutsos        120

# Scalability

- Google: > 450,000 processors in clusters of ~2000 processors each [Barroso, Dean, Hölzle, *"Web Search for a Planet: The Google Cluster Architecture"* IEEE Micro 2003]
- Yahoo: 5Pb of data [Fayyad, KDD'07]
- Problem: machine failures, on a daily basis
- How to parallelize data mining tasks, then?
- A: map/reduce – hadoop (open-source clone) http://hadoop.apache.org/

15-826                        (c) 2011  C. Faloutsos                        121

# Outline – Algorithms & results

|  | Centralized | Hadoop/ PEGASUS |
|---|---|---|
| Degree Distr. | old | old |
| Pagerank | old | old |
| Diameter/ANF | old | **HERE** |
| Conn. Comp | old | **HERE** |
| Triangles | **done** | **HERE** |
| Visualization | **started** | |

15-826                        (c) 2011  C. Faloutsos                        122

# HADI for diameter estimation

- *Radius Plots for Mining Tera-byte Scale Graphs* **U Kang**, Charalampos Tsourakakis, Ana Paula Appel, Christos Faloutsos, Jure Leskovec, SDM'10
- Naively: diameter needs **O(N**2)** space and up to O(N**3) time – **prohibitive** (N~1B)
- Our HADI: linear on E (~10B)
  - Near-linear scalability wrt # machines
  - Several optimizations -> 5x faster

15-826                        (c) 2011  C. Faloutsos                        123

## Slide 124

**Carnegie Mellon**

Count

Number of Nodes

$10^9$
$10^8$
$10^7$
$10^6$
$10^5$
$10^4$
$10^3$
$10^2$
$10^1$
$10^0$

0  5  10  15  20  25  30
Radius

Radius

19+ [Barabasi+]

~1999, ~1M nodes

15-826          (c) 2011 C. Faloutsos          124

## Slide 125

**Carnegie Mellon**

Count

Number of Nodes

$10^9$
$10^8$
$10^7$
$10^6$
$10^5$
$10^4$
$10^3$
$10^2$
$10^1$
$10^0$

0  5  10  15  20  25  30
Radius

Radius

??

19+ [Barabasi+]

~1999, ~1M nodes

YahooWeb graph  (120Gb, 1.4B nodes, 6.6 B edges)
• Largest publicly available graph ever studied.

15-826          (c) 2011 C. Faloutsos          125

## Slide 126

**Carnegie Mellon**

Count

Number of Nodes

$10^9$
$10^8$
$10^7$
$10^6$
$10^5$
$10^4$
$10^3$
$10^2$
$10^1$
$10^0$

0  5  10  15  20  25  30
Radius

Radius

**14 (dir.)**

**~7 (undir.)**

19+? [Barabasi+]

YahooWeb graph  (120Gb, 1.4B nodes, 6.6 B edges)
• Largest publicly available graph ever studied.

15-826          (c) 2011 C. Faloutsos          126

Faloutsos

**Carnegie Mellon**

Count $10^9$
$10^8$
$10^7$
$10^6$
$10^5$
$10^4$
$10^3$
$10^2$
$10^1$
$10^0$

Number of Nodes

**14 (dir.)**

**~7 (undir.)**

**19+? [Barabasi+]**

0    5    10   15   20   25   30
Radius

**Radius**

YahooWeb graph  (120Gb, 1.4B nodes, 6.6 B edges)
•7 degrees of separation (!)
•Diameter: shrunk

15-826                    (c) 2011  C. Faloutsos                    127

---

**Carnegie Mellon**

Count $10^9$
$10^8$
$10^7$
$10^6$
$10^5$
$10^4$
$10^3$
$10^2$
$10^1$
$10^0$

Number of Nodes

**~7 (undir.)**

0    5    10   15   20   25   30
Radius

**Radius**

YahooWeb graph  (120Gb, 1.4B nodes, 6.6 B edges)
Q: Shape?

15-826                    (c) 2011  C. Faloutsos                    128

---

**Carnegie Mellon**

$10^9$
$10^8$
$10^7$
$10^6$
$10^5$
$10^4$
$10^3$
$10^2$
$10^1$
$10^0$

Number of Nodes

YahooWeb

**S**

**Multi-Modal**

**Effective
Diameter = 7.62**

0    5    10   15   20   25   30
Radius

YahooWeb graph  (120Gb, 1.4B nodes, 6.6 B edges)
• effective diameter: surprisingly small.
• Multi-modality (?!)

15-826                    (c) 2011  C. Faloutsos                    129

Radius Plot of **GCC** of YahooWeb.

15-826        (c) 2011 C. Faloutsos       130



YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)
• effective diameter: surprisingly small.
• Multi-modality: probably mixture of cores .

15-826        (c) 2011 C. Faloutsos       131



YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)
• effective diameter: surprisingly small.
• Multi-modality: probably mixture of cores .

15-826        (c) 2011 C. Faloutsos       132

44

Faloutsos

**Carnegie Mellon**

Conjecture:

~7

YahooWeb graph (120Gb, 1.4B nodes, 6.6 B edges)
• effective diameter: surprisingly small.
• Multi-modality: probably mixture of cores .

15-826      (c) 2011 C. Faloutsos      133

---

**Carnegie Mellon**

details

Running time - Kronecker and Erdos-Renyi
Graphs with billions edges.

---

**Carnegie Mellon**

## Outline – Algorithms & results

|  | **Centralized** | **Hadoop/ PEGASUS** |
|---|---|---|
| Degree Distr. | old | old |
| Pagerank | old | old |
| Diameter/ANF | old | **HERE** |
| Conn. Comp | old | **HERE** |
| Triangles |  | **HERE** |
| Visualization | **started** |  |

15-826      (c) 2011 C. Faloutsos      135

**Carnegie Mellon**

### Generalized Iterated Matrix Vector Multiplication (GIMV)

*PEGASUS: A Peta-Scale Graph Mining System - Implementation and Observations*.
U Kang, Charalampos E. Tsourakakis,
and Christos Faloutsos.
(ICDM) 2009, Miami, Florida, USA.
Best Application Paper (runner-up)**.**

15-826                 (c) 2011  C. Faloutsos                 136

---

**Carnegie Mellon**

### Generalized Iterated Matrix Vector Multiplication (GIMV)

details

• PageRank
• proximity (RWR)
• Diameter
• Connected components
• (eigenvectors,
•  Belief Prop.
• … )

Matrix – vector
Multiplication
(iterated)

15-826                 (c) 2011  C. Faloutsos                 137

---

**Carnegie Mellon**

### Example: GIM-V At Work

• Connected Components – 4 observations:



15-826                 (c) 2011  C. Faloutsos                 138

**Carnegie Mellon**

# Example: GIM-V At Work

- Connected Components

Count

YahooWeb

1) 10K x
larger
than next

**Giant Connected Component**

Size

15-826     (c) 2011 C. Faloutsos     139

---

**Carnegie Mellon**

# Example: GIM-V At Work

- Connected Components

Count

YahooWeb

2) ~0.7B
singleton
nodes

**Giant Connected Component**

Size

15-826     (c) 2011 C. Faloutsos     140

---

**Carnegie Mellon**

# Example: GIM-V At Work

- Connected Components

Count

YahooWeb

3) SLOPE!

**Giant Connected Component**

Size

15-826     (c) 2011 C. Faloutsos     141

Faloutsos

---

**Carnegie Mellon**

## Example: GIM-V At Work

- Connected Components

Count

10^9
10^8
10^7
10^6
10^5
10^4
10^3
10^2
10^1
10^0

YahooWeb

300-size cmpt X 500. Why?

1100-size cmpt X 65. Why?

**Giant Connected Component**

4) Spikes!

10^0 10^1 10^2 10^3 10^4 10^5 10^6 10^7 10^8 10^9

Size

15-826        (c) 2011  C. Faloutsos        142

---

**Carnegie Mellon**

## Example: GIM-V At Work

- Connected Components

Count

10^9
10^8
10^7
10^6
10^5
10^4
10^3
10^2
10^1
10^0

YahooWeb

suspicious financial-advice sites (not existing now)

**Giant Connected Component**

10^0 10^1 10^2 10^3 10^4 10^5 10^6 10^7 10^8 10^9

Size

15-826        (c) 2011  C. Faloutsos        143

---

**Carnegie Mellon**

## GIM-V At Work

- Connected Components over Time
- **LinkedIn: 7.5M nodes and 58M edges**

2003 — Unstable Slope — **Giant Connected Component**

2004 — Slope = - 2.75 — **Giant Connected Component**

2005 — Slope = - 2.75 — **Giant Connected Component**

2006 — Slope = - 2.75 — **Giant Connected Component**

Stable tail slope after the gelling point

15-826        (c) 2011  C. Faloutsos        144

**CarnegieMellon**

### Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
- Problem#2: Tools
- Problem#3: Scalability
➡ • Conclusions

15-826                    (c) 2011  C. Faloutsos                    145

---

**CarnegieMellon**

### OVERALL CONCLUSIONS – low level:

- Several new **patterns** (fortification, shrinking diameter, triangle-laws, conn. components, etc)

- New **tools**:
    – anomaly detection (OddBall), belief propagation, immunization

- **Scalability**: PEGASUS / hadoop

15-826                    (c) 2011  C. Faloutsos                    146
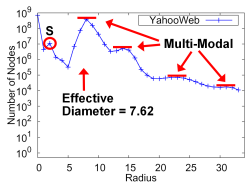
---

**CarnegieMellon**

### OVERALL CONCLUSIONS – high level

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



15-826                    (c) 2011  C. Faloutsos                    147

**Carnegie Mellon**

# References

- Leman Akoglu, Christos Faloutsos: *RTG: A Recursive Realistic Graph Generator Using Random Typing*. ECML/PKDD (1) 2009: 13-28

- Deepayan Chakrabarti, Christos Faloutsos: *Graph mining: Laws, generators, and algorithms*. ACM Comput. Surv. 38(1): (2006)

15-826      (c) 2011 C. Faloutsos      148

**Carnegie Mellon**

# References

- Deepayan Chakrabarti, Yang Wang, Chenxi Wang, Jure Leskovec, Christos Faloutsos: *Epidemic thresholds in real networks*. ACM Trans. Inf. Syst. Secur. 10(4): (2008)

- Deepayan Chakrabarti, Jure Leskovec, Christos Faloutsos, Samuel Madden, Carlos Guestrin, Michalis Faloutsos: *Information Survival Threshold in Sensor and P2P Networks*. INFOCOM 2007: 1316-1324

15-826      (c) 2011 C. Faloutsos      149

**Carnegie Mellon**

# References

- Christos Faloutsos, Tamara G. Kolda, Jimeng Sun: *Mining large graphs and streams using matrix and tensor tools*. Tutorial, SIGMOD Conference 2007: 1174

15-826      (c) 2011 C. Faloutsos      150

**CarnegieMellon**

# References

- T. G. Kolda and J. Sun. *Scalable Tensor Decompositions for Multi-aspect Data Mining*. In: ICDM 2008, pp. 363-372, December 2008.

15-826        (c) 2011 C. Faloutsos      151

**CarnegieMellon**

# References

- Jure Leskovec, Jon Kleinberg and Christos Faloutsos *Graphs over Time: Densification Laws, Shrinking Diameters and Possible Explanations*, KDD 2005 (Best Research paper award).
- Jure Leskovec, Deepayan Chakrabarti, Jon M. Kleinberg, Christos Faloutsos: *Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication*. PKDD 2005: 133-145

15-826        (c) 2011 C. Faloutsos      152

**CarnegieMellon**

# References

- Jimeng Sun, Yinglian Xie, Hui Zhang, Christos Faloutsos. *Less is More: Compact Matrix Decomposition for Large Sparse Graphs*, SDM, Minneapolis, Minnesota, Apr 2007.
- Jimeng Sun, Spiros Papadimitriou, Philip S. Yu, and Christos Faloutsos, *GraphScope: Parameter-free Mining of Large Time-evolving Graphs* ACM SIGKDD Conference, San Jose, CA, August 2007

15-826        (c) 2011 C. Faloutsos      153

**CarnegieMellon**

# References

- Jimeng Sun, Dacheng Tao, Christos Faloutsos: *Beyond streams and graphs: dynamic tensor analysis*. KDD 2006: 374-383

15-826       (c) 2011 C. Faloutsos       154

**CarnegieMellon**

# References

- Hanghang Tong, Christos Faloutsos, and Jia-Yu Pan, *Fast Random Walk with Restart and Its Applications*, ICDM 2006, Hong Kong.
- Hanghang Tong, Christos Faloutsos, *Center-Piece Subgraphs: Problem Definition and Fast Solutions*, KDD 2006, Philadelphia, PA

15-826       (c) 2011 C. Faloutsos       155

**CarnegieMellon**

# References

- Hanghang Tong, Christos Faloutsos, Brian Gallagher, Tina Eliassi-Rad: Fast best-effort pattern matching in large attributed graphs. KDD 2007: 737-746
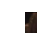
15-826       (c) 2011 C. Faloutsos       156

**Carnegie Mellon**

## (Project info)

`www.cs.cmu.edu/~pegasus`

Project Pegasus

Chau,
Polo

Koutra,
Danae

Prakash,
Aditya

Akoglu,
Leman

Kang, U

McGlohon,
Mary

Tong,
Hanghang