**Carnegie Mellon**

# 15-826: Multimedia Databases and Data Mining

Lecture #28: Graph mining - patterns

*Christos Faloutsos*

---

**Carnegie Mellon**

# Must-read Material

- [Graph minining textbook] Deepayan Chakrabarti and Christos Faloutsos *Graph Mining: Laws, Tools and Case Studies*, Morgan Claypool, 2012
  - Part I (patterns)

15-826      (c) C. Faloutsos, 2017      2

---

**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

15-826      (c) C. Faloutsos, 2017      3

---

**Carnegie Mellon**

# Main outline

- Introduction
- Indexing
- Mining
  - Graphs – patterns
  - Graphs – generators and tools
  - Association rules
  - …

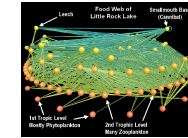15-826      (c) C. Faloutsos, 2017      4

Faloutsos

---

**Outline**

- Introduction – Motivation
- Problem#1: Patterns in graphs
- Problem#2: Scalability
- Conclusions

---

**Graphs - why should we care?**
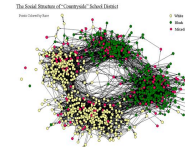
Food Web
[Martinez ' 91]

Friendship Network
[Moody ' 01]

Internet Map
[lumeta.com]

---

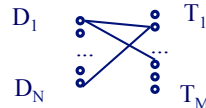**Graphs - why should we care?**

- IR: bi-partite graphs (doc-terms)

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

- web: hyper-text graph

- ... and more:

---

**Graphs - why should we care?**

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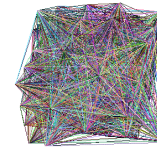- ....

Faloutsos

---

## Outline

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

15-826     (c) C. Faloutsos, 2017     9
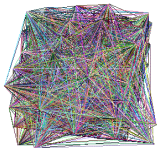
---

## 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?

15-826     (c) C. Faloutsos, 2017     10

---

## 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) C. Faloutsos, 2017     11

---

## 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?
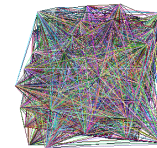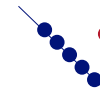  - To spot **anomalies** (rarities), we have to discover **patterns**
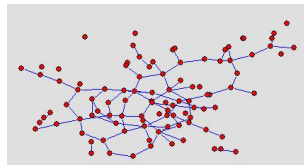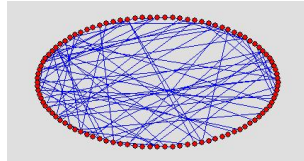  - **Large** datasets reveal patterns/anomalies that may be invisible otherwise…

15-826     (c) C. Faloutsos, 2017     12

3

## Slide 13

**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) C. Faloutsos, 2017      13

## Slide 14

**Carnegie Mellon**

# Graph mining

- Are real graphs random?

15-826      (c) C. Faloutsos, 2017      14

## Slide 15

**Carnegie Mellon**

# Laws and patterns

- Are real graphs random?
- A: NO!!
  - Diameter ('6 degrees', 'Kevin Bacon')
  - in- and out- degree distributions
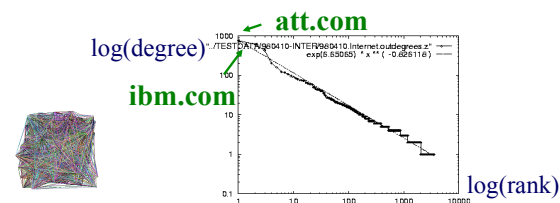  - other (surprising) patterns

- So, let's look at the data

15-826      (c) C. Faloutsos, 2017      15

## Slide 16

**Carnegie Mellon**

# Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

**internet domains**

att.com

log(degree)

ibm.com

log(rank)

15-826      (c) C. Faloutsos, 2017      16

Faloutsos

---

# 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) C. Faloutsos, 2017     17

---

# Solution# S.1

- Q: So what?

**internet domains**

att.com

log(degree)

ibm.com

-0.82

log(rank)

15-826     (c) C. Faloutsos, 2017     18

---

# Solution# S.1

- Q: So what?    = friends of friends (F.O.F.)
- A1: # of two-step-away pairs:

**internet domains**

att.com

log(degree)

ibm.com

-0.82

log(rank)

15-826     (c) C. Faloutsos, 2017     19

---

# Solution# S.1

- Q: So what?    = friends of friends (F.O.F.)
- A1: # of two-step-away pairs: 100^2 * N= 10 Trillion

**internet domains**

att.com

log(degree)

ibm.com

-0.82

log(rank)

15-826     (c) C. Faloutsos, 2017     20

---

5

Faloutsos

## Solution# S.1

- Q: So what?
- A1: # of two-step-away pairs: 100^    Trillion
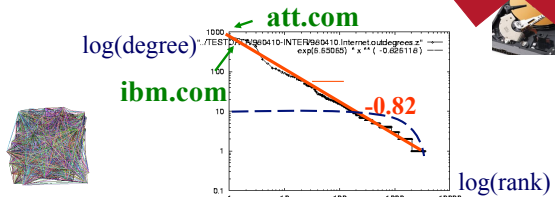
= friends of friends (F.O.F.)

**internet domains**

att.com

log(degree)

ibm.com

**-0.82**

log(rank)

15-826          (c) C. Faloutsos, 2017          21

---

**Gaussian trap**

## Solution# S.1

- Q: So what?
- A1: # of two-step-away pairs: O(d_max ^2) ~ 10M^2

= friends of friends (F.O.F.)

**internet domains**

⇩

~0.8PB ->
a data center(!)

att.com

log(degree)

ibm.com

**-0.82**

15-826          (c) C. Faloutsos, 2017     DCO @ CMU     22

---

**Gaussian trap**

## Solution# S.1

- Q: So what?
- A1: # of two-step-aw        ?) ~ 10M^2
  inte

**Such patterns ->
New algorithms**

~0.8PB ->
a data center(!)

**-0.82**

15-826          (c) C. Faloutsos, 2017          23

---

## Observation – big-data:

- $O(N^2)$ algorithms are ~intractable  - N=1B

- $N^2$ seconds = 31B years (>2x age of universe)

1B

1B

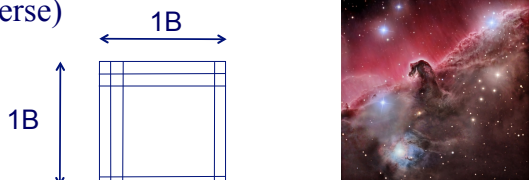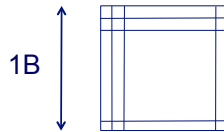15-826          (c) C. Faloutsos, 2017          24

6

**Slide 25**

Carnegie Mellon

# Observation – big-data:

- O($N^2$) algorithms are ~intractable  - N=1B

31M
- $N^2$ seconds = 31B years
- 1,000 machines

1B

**Slide 26**

Carnegie Mellon

# Observation – big-data:

- O($N^2$) algorithms are ~intractable  - N=1B

31K
- $N^2$ seconds = 31B years
- 1M machines

1B

Google Y!

**Slide 27**
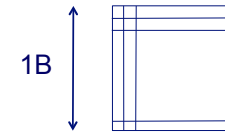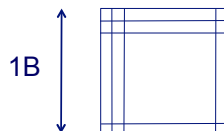
Carnegie Mellon

# Observation – big-data:

- O($N^2$) algorithms are ~intractable  - N=1B

3
- $N^2$ seconds = 31B years
- 10B machines ~ $10Trillion
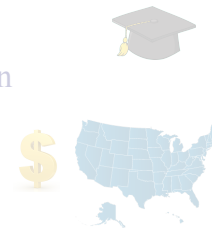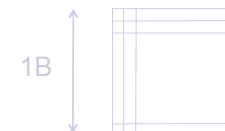
1B

**Slide 28**

Carnegie Mellon

# Observation – big-data:

- O($N^2$) algorithms are ~intractable  - N=1B

**And parallelism might not help**

3
- $N^2$ seconds = 31B years
- 10B machines ~ $10Trillion

1B

Faloutsos

---

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

Eigenvalue

Exponent = slope

*E = -0.48*

May 2001

**A x = λ x**

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

Rank of decreasing eigenvalue

15-826          (c) C. Faloutsos, 2017          29

---

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

Eigenvalue

Exponent = slope

*E = -0.48*

May 2001

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

Rank of decreasing eigenvalue

15-826          (c) C. Faloutsos, 2017          30

---

## But:

How about graphs from other domains?

15-826          (c) C. Faloutsos, 2017          31

---

## More power laws:

- web hit counts [w/ A. Montgomery]

Web Site Traffic

Count
(log scale)

Zipf

``ebay''

users

sites

in-degree (log scale)

15-826          (c) C. Faloutsos, 2017          32

8

## epinions.com

**Carnegie Mellon**

count



(out) degree

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

trusts-2000-people user

---

**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'

---

**Carnegie Mellon**

## Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
  - Static graphs
    - degree, diameter, eigen,
    - Triangles
  - Weighted graphs
  - Time evolving graphs

---

**Carnegie Mellon**

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

- Real social networks have a lot of triangles

---

**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) C. Faloutsos, 2017　　　　37

---

**Carnegie Mellon**

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

HEP-TH

ASN

Epinions

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

15-826　　　　　tsos, 2017　　　　38

---

**Carnegie Mellon**

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

HEP-TH

ASN

Epinions

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

15-826　　　　　tsos, 2017　　　　39

---

**Carnegie Mellon**

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

Reuters

DTPL
slope 1.68
slope 1.68

SN

DTPL
slope 1.74
slope 1.73

Epinions

DTPL
slope 1.61
slope 1.59

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

15-826　　　　(c) C. Faloutsos, 2017　　　　40

10

Faloutsos

details

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

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

15-826          (c) C. Faloutsos, 2017          41

---

Carnegie Mellon

details

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

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) C. Faloutsos, 2017          42

---

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) C. Faloutsos, 2017          43

---

Carnegie Mellon

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)

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

15-826          (c) C. Faloutsos, 2017          44

---

11

**Carnegie Mellon**

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)
[U Kang, Brendan Meeder, +, PAKDD'11]

15-826     (c) C. Faloutsos, 2017     45

**Carnegie Mellon**

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)
[U Kang, Brendan Meeder, +, PAKDD'11]

15-826     (c) C. Faloutsos, 2017     46

**Carnegie Mellon**

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)
[U Kang, Brendan Meeder, +, PAKDD'11]

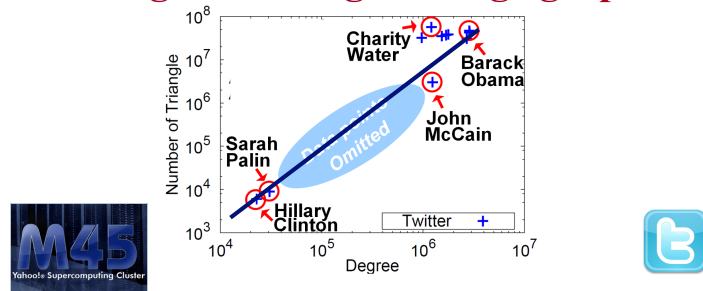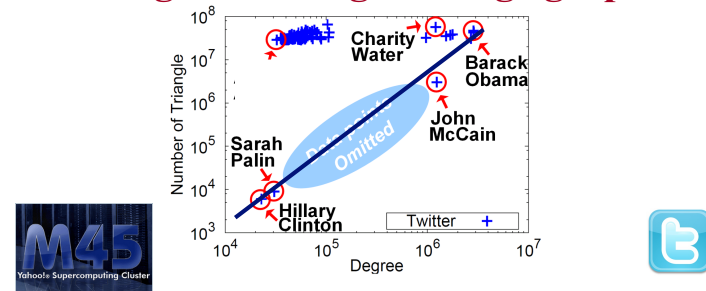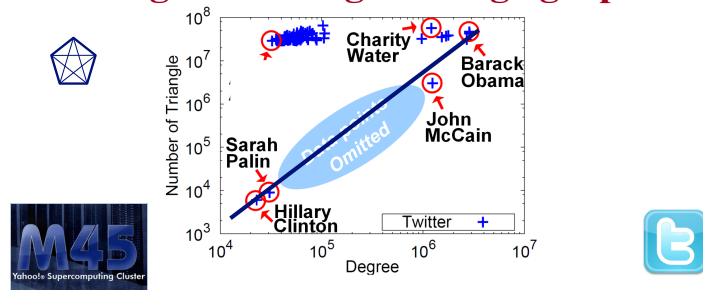15-826     (c) C. Faloutsos, 2017     47

**Carnegie Mellon**

# Triangle counting for large graphs?



Anomalous nodes in Twitter(~ 3 billion edges)
[U Kang, Brendan Meeder, +, PAKDD'11]

15-826     (c) C. Faloutsos, 2017     48

12

**CarnegieMellon**

# Any other 'laws' ?

Yes!

(c) C. Faloutsos, 2017

---

**CarnegieMellon**

# Any other 'laws' ?

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

(c) C. Faloutsos, 2017

---

**CarnegieMellon**

# Any other 'laws' ?

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

(c) C. Faloutsos, 2017

---

**CarnegieMellon**

# 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)

(c) C. Faloutsos, 2017

---

**Carnegie Mellon**

# Any other 'laws'?

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



15-826          (c) C. Faloutsos, 2017          53

---

**Carnegie Mellon**

# 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.

*Useful for fraud detection!*

15-826          (c) C. Faloutsos, 2017          54

---

**Carnegie Mellon**

# EigenSpokes

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

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

15-826          (c) C. Faloutsos, 2017          55

---

**Carnegie Mellon**

**details**

# EigenSpokes

- 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) C. Faloutsos, 2017          56

14

Faloutsos

**EigenSpokes**

details

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

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

$$N$$

$$N$$

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

---

Carnegie Mellon

**EigenSpokes**
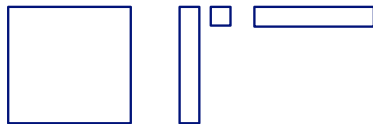
details

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

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

$$N$$

$$N$$

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

---

Carnegie Mellon

**EigenSpokes**

details

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

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

$$N$$

$$N$$

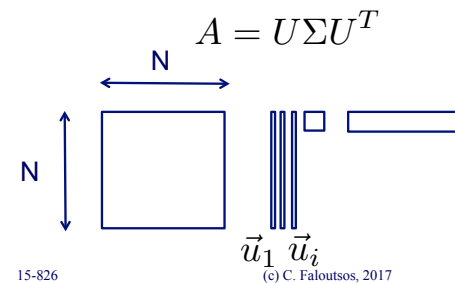$$\vec{u}_1 \; \vec{u}_i$$

---

Carnegie Mellon

**EigenSpokes**

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

2nd Principal component u2

u1

1st Principal component

## EigenSpokes

Carnegie Mellon

- 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) C. Faloutsos, 2017          61

---

Carnegie Mellon

## EigenSpokes - pervasiveness

- Present in mobile social graph
  - across time and space

- Patent citation graph

15-826          (c) C. Faloutsos, 2017          62

---

Carnegie Mellon

## EigenSpokes - explanation

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

15-826          (c) C. Faloutsos, 2017          63

---

Carnegie Mellon

## EigenSpokes - explanation

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

15-826          (c) C. Faloutsos, 2017          64

16

Faloutsos

---

**Carnegie Mellon**

## EigenSpokes - explanation

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

---

**Carnegie Mellon**

## EigenSpokes - explanation

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

spy plot of top 20 nodes

So what?
- Extract nodes with high *scores*
- high connectivity
- Good "communities"

$v_1$   $v_2$   $v_3$   $v_4$   $v_5$   $v_6$   $v_7$   $v_8$   $v_9$

---

**Carnegie Mellon**

## Bipartite Communities!

patents from same inventor(s)

`cut-and-paste' bibliography!

magnified bipartite community

*Useful for fraud detection!*

---

**Carnegie Mellon**

## Bipartite Communities!

IP – port scanners

victims

*Useful for fraud detection!*

## Slide 69

**Carnegie Mellon**

### Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
  - Static graphs
    - degree, diameter, eigen,
    - Triangles
  - Weighted graphs
  ➡ - Time evolving graphs
- Problem#2: Scalability
- Conclusions

15-826                     (c) C. Faloutsos, 2017                     69

## Slide 70

**Carnegie Mellon**

### 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) C. Faloutsos, 2017                     70

## Slide 71

**Carnegie Mellon**

### Observation W.1: Fortification

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

15-826                     (c) C. Faloutsos, 2017                     71

## Slide 72

**Carnegie Mellon**

### Observation W.1: Fortification

**More donors, more $ ?**

$10  'Reagan'
$5
'Clinton'
$7

15-826                     (c) C. Faloutsos, 2017                     72

Faloutsos

---

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

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

**More donors, even more $**

$10

$5

In-weights ($)

**Orgs-Candidates**

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

Edges (# donors)

(c) C. Faloutsos, 2017 73

---

## Outline

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

(c) C. Faloutsos, 2017 74

---

## Problem: Time evolution

- with Jure Leskovec (CMU -> Stanford)

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

(c) C. Faloutsos, 2017 75

---

## T.1 Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:
  - [diameter ~ O( $N^{1/3}$)]
  - diameter ~ O(log N)
  - diameter ~ O(log log N)
- What is happening in real data?

diameter

(c) C. Faloutsos, 2017 76

---

19

Faloutsos

---

## T.1 Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:
  - [diameter ~ $O(N^{1/3})$]
  - diameter ~ $O(\log N)$
  - diameter ~ $O(\log \log N)$
- What is happening in real data?
- Diameter **shrinks** over time

15-826          (c) C. Faloutsos, 2017          77

---

## T.1 Diameter – "Patents"

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



diameter

Effective diameter vs time [years]

Full graph / Post '85 subgraph / Post '85 subgraph, no past

15-826          (c) C. Faloutsos, 2017          78

---

## 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) C. Faloutsos, 2017          79

---

## 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) C. Faloutsos, 2017          80

---

20

Faloutsos



**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)

1.66

1999

1975

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

Number of edges

Number of nodes

N(t)

15-826          (c) C. Faloutsos, 2017          81



**Outline**

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

15-826          (c) C. Faloutsos, 2017          82

**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) C. Faloutsos, 2017          83

**[ 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

Time = 1914

t=1914

Diameter

Time

15-826          (c) C. Faloutsos, 2017          84

21

## Slide 85

# 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?

## Slide 86

# 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?

## Slide 87

# 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?

## Slide 88

# 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

Faloutsos

# 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) C. Faloutsos, 2017      89

---

# T.4 : popularity over time

# in links

1    2    3     lag: days after post

Post popularity drops-off – exponentially?

@t

@t + **lag**

15-826      (c) C. Faloutsos, 2017      90

---

# T.4 : popularity over time

# in links
(**log**)

Posts
= 541905.74 x   $R^2$=1.00

days after post
(**log**)

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

15-826      (c) C. Faloutsos, 2017      91

---

# T.4 : popularity over time

# in links
(**log**)

-1.6

Posts
= 541905.74 x$^{-1.60}$ $R^2$=1.00

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

Brown Noise

DFT of Brown Noise

15-826      (c) C. Faloutsos, 2017      92

23

## -1.5 slope

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



Figure 1 | The correspondence patterns of Darwin and Einstein.

93

## 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) C. Faloutsos, 2017       94

## Probably, power law (?)



??

15-826      (c) C. Faloutsos, 2017      95

## No Power Law!



15-826      (c) C. Faloutsos, 2017      96

## Slide 97

**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



15-826     (c) C. Faloutsos, 2017     97

## Slide 98

**Carnegie Mellon**

# Log-logistic distribution

- CDF(t)/(1- CDF(t)) == OR(t)
- For log-logistic: $\log[OR(t)] = \beta + \rho*\log(t)$

Odds ratio=

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

== power law



15-826     (c) C. Faloutsos, 2017     98

## Slide 99

**Carnegie Mellon**

# Log-logistic distribution

- CDF(t)/(1- CDF(t)) == OR(t)
- For log-logistic: $\log[OR(t)] = \beta + \rho*\log(t)$

OR(t)



- PDF looks like hyperbola;
- and, if clipped, like power-law

15-826     (c) C. Faloutsos, 2017     99

## Slide 100

**Carnegie Mellon**

# Log-logistic distribution

- CDF(t)/(1- CDF(t)) == OR(t)
- For log-logistic: $\log[OR(t)] = \beta + \rho*\log(t)$

OR(t)



Duration ( t )

15-826     100

25

Faloutsos

## Log-logistic distribution

- Logistic distribution: CDF -> sigmoid
- LOG-Logistic distribution:

$x \rightarrow ln(x)$



CDF(x) = 1/(1+exp(-x))     CDF(x) = 1/(1+1/x )

## Log-logistic distribution

- Logistic distribution: CDF -> sigmoid
- LOG-Logistic distribution:



CDF(x) = 1/(1+exp(-(x-m)/s))  CDF(x) = 1/(1+exp(-(ln(x)-m)/s))

## 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)

## Outliers:

26

## Slide 105

# Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
➡ - Problem#2: Scalability -PEGASUS
- Conclusions

15-826      (c) C. Faloutsos, 2017      105

## Slide 106

# 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) C. Faloutsos, 2017      106

## Slide 107

# 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** |  |

(arrow pointing to Diameter/ANF row)

15-826      (c) C. Faloutsos, 2017      107

## Slide 108

# 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) C. Faloutsos, 2017      108

## Slide 109

Carnegie Mellon

Count

Number of Nodes: $10^9$, $10^8$, $10^7$, $10^6$, $10^5$, $10^4$, $10^3$, $10^2$, $10^1$, $10^0$

Radius: 0, 5, 10, 15, 20, 25, 30

19+ [Barabasi+]

~1999, ~1M nodes

15-826    (c) C. Faloutsos, 2017    109

## Slide 110

Carnegie Mellon

Count

Number of Nodes: $10^9$, $10^8$, $10^7$, $10^6$, $10^5$, $10^4$, $10^3$, $10^2$, $10^1$, $10^0$

Radius: 0, 5, 10, 15, 20, 25, 30

??

19+ [Barabasi+]

~1999, ~1M nodes

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

15-826    (c) C. Faloutsos, 2017    110

## Slide 111

Carnegie Mellon

Count

Number of Nodes: $10^9$, $10^8$, $10^7$, $10^6$, $10^5$, $10^4$, $10^3$, $10^2$, $10^1$, $10^0$

Radius: 0, 5, 10, 15, 20, 25, 30

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) C. Faloutsos, 2017    111

## Slide 112

Carnegie Mellon

Count

Number of Nodes: $10^9$, $10^8$, $10^7$, $10^6$, $10^5$, $10^4$, $10^3$, $10^2$, $10^1$, $10^0$

Radius: 0, 5, 10, 15, 20, 25, 30

14 (dir.)

~7 (undir.)

19+? [Barabasi+]

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

15-826    (c) C. Faloutsos, 2017    112

### Slide 113

**Carnegie Mellon**



Count

~7 (undir.)

Radius

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

15-826      (c) C. Faloutsos, 2017      113

### Slide 114

**Carnegie Mellon**



S

Multi-Modal

Effective Diameter = 7.62

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

15-826      (c) C. Faloutsos, 2017      114

### Slide 115

**Carnegie Mellon**



GCC

C   google.com

Radius Plot of **GCC** of YahooWeb.

15-826      (c) C. Faloutsos, 2017      115

### Slide 116

**Carnegie Mellon**



S

Multi-Modal

Effective Diameter = 7.62

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

15-826      (c) C. Faloutsos, 2017      116

29

Faloutsos

Conjecture:

EN

DE

BR

~7

Effective Diameter = 7.62

Multi-Modal

YahooWeb

Number of Nodes

Radius

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

Carnegie Mellon

Conjecture:

~7

Effective Diameter = 7.62

Multi-Modal

YahooWeb

Number of Nodes

Radius

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

Carnegie Mellon

details

7.6x faster

5.1x

3.8x

3.2x

Run time in hours

HADI-plain
HADI-BSE
HADI-BL
HADI-OPT

KR-2B    KR-1.1B    ER-2B    ER-1.1B
Data

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** |  |

**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) C. Faloutsos, 2017      121

---

**Carnegie Mellon**

*details*

# Generalized Iterated Matrix Vector Multiplication (GIMV)

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

Matrix – vector
Multiplication
(iterated)

15-826      (c) C. Faloutsos, 2017      122

---

**Carnegie Mellon**

# Example: GIM-V At Work

• Connected Components – 4 observations:



Count

YahooWeb

Giant
Connected
Component

Size

15-826      (c) C. Faloutsos, 2017      123

---

**Carnegie Mellon**

# Example: GIM-V At Work

• Connected Components



Count

YahooWeb

Giant
Connected
Component

1) 10K x
larger
than next

Size

15-826      (c) C. Faloutsos, 2017      124

Faloutsos

# Example: GIM-V At Work

- Connected Components

Count

YahooWeb

2) ~0.7B singleton nodes

**Giant Connected Component**

Size

(c) C. Faloutsos, 2017

---

# Example: GIM-V At Work

- Connected Components

Count

YahooWeb

3) SLOPE!

**Giant Connected Component**

Size

(c) C. Faloutsos, 2017

---

# Example: GIM-V At Work

- Connected Components

Count

YahooWeb

300-size cmpt X 500. Why?

1100-size cmpt X 65. Why?

**Giant Connected Component**

4) Spikes!

Size

(c) C. Faloutsos, 2017

---

# Example: GIM-V At Work

- Connected Components

Count

YahooWeb

suspicious financial-advice sites (not existing now)

**Giant Connected Component**

Size

(c) C. Faloutsos, 2017

Faloutsos

# GIM-V At Work

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



Stable tail slope
after the gelling point

15-826                    (c) C. Faloutsos, 2017                    129

---

# Outline

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

15-826                    (c) C. Faloutsos, 2017                    130

---

# OVERALL CONCLUSIONS – low level:

- Several new **patterns** (fortification, shrinking diameter, triangle-laws, conn. components, etc)
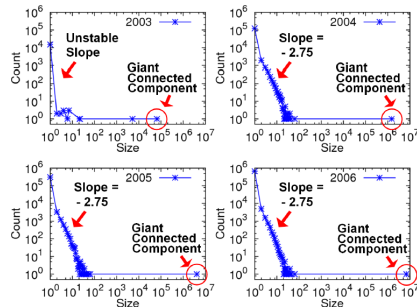
- Log-logistic distribution: ubiquitus

- New **tools**:

  – anomaly detection (OddBall), belief propagation, immunization

- **Scalability**: PEGASUS / hadoop

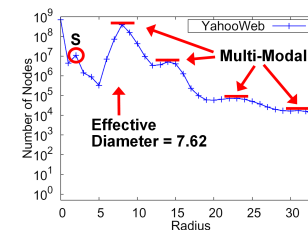15-826         (c) C. Faloutsos, 2017         131

---

# OVERALL CONCLUSIONS – high level

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



15-826                    (c) C. Faloutsos, 2017                    132

Faloutsos

**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) C. Faloutsos, 2017          133

**References**

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

15-826          (c) C. Faloutsos, 2017          134

**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) C. Faloutsos, 2017          135

**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) C. Faloutsos, 2017          136

34

**Carnegie Mellon**

# References

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

15-826        (c) C. Faloutsos, 2017      137

**Carnegie Mellon**

# 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) C. Faloutsos, 2017      138

**Carnegie Mellon**

# 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) C. Faloutsos, 2017      139

**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

15-826        (c) C. Faloutsos, 2017      140