

15-826: Multimedia Databases and Data Mining

Lecture #28: Graph mining - patterns
Christos Faloutsos

Must-read Material

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

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

Must-read Material (cont' d)

- D. Chakrabarti and C. Faloutsos, Graph Mining: Laws, Generators and Algorithms, in ACM Computing Surveys, 38(1), 2006

Carnegie Mellon

Main outline




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

15-826 (c) C. Faloutsos, 2016 5

Carnegie Mellon

Outline




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

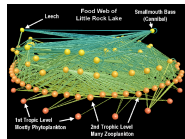
15-826 (c) C. Faloutsos, 2016 6

Carnegie Mellon


Graphs - why should we care?



Friendship Network
[Moody '01]



Food Web
[Martinez '91]




Internet Map
[lumeta.com]

15-826 (c) C. Faloutsos, 2016 7

Carnegie Mellon

Graphs - why should we care?

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



- web: hyper-text graph
- ... and more:

15-826 (c) C. Faloutsos, 2016 8

Graphs - why should we care?

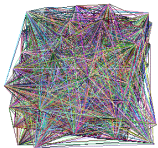
- ‘viral’ marketing
- web-log (‘blog’) news propagation
- computer network security: email/IP traffic and anomaly detection
-

Outline



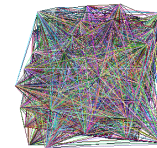
- Introduction – Motivation
- ➔ • Problem#1: Patterns in graphs
 - Static graphs
 - Weighted graphs
 - Time evolving graphs
- Problem#2: Scalability
- Conclusions

Problem #1 - network and graph mining



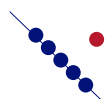
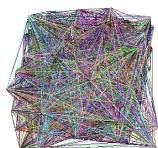
- What does the Internet look like?
- What does FaceBook look like?
- What is ‘normal’ / ‘abnormal’ ?
- which patterns/laws hold?

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**

Problem #1 - network and graph mining



15-826

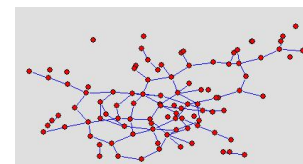
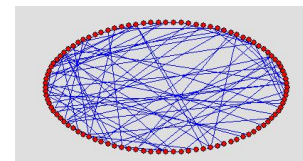
- What does the Internet look like?
- What does FaceBook look like?
- What is 'normal' / 'abnormal'?
- which patterns/laws hold?
 - To spot **anomalies** (rarities), we have to discover **patterns**
 - **Large** datasets reveal patterns/anomalies that may be invisible otherwise...

(c) C. Faloutsos, 2016

13

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, 2016

14

Graph mining

- Are real graphs random?

15-826

(c) C. Faloutsos, 2016

15

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

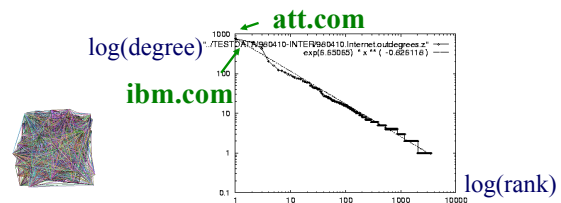
(c) C. Faloutsos, 2016

16

Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

internet domains



15-826

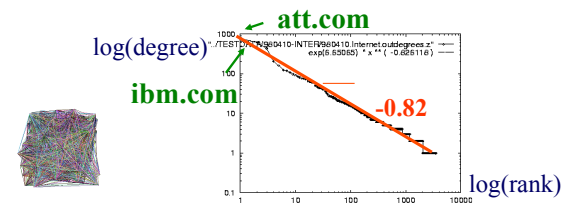
(c) C. Faloutsos, 2016

17

Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

internet domains



15-826

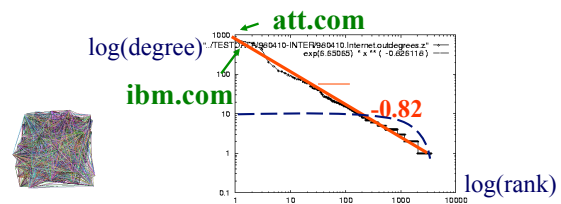
(c) C. Faloutsos, 2016

18

Solution# S.1

- Q: So what?

internet domains



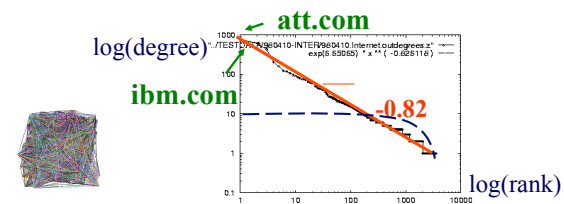
15-826

(c) C. Faloutsos, 2016

19

Solution# S.1

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



15-826

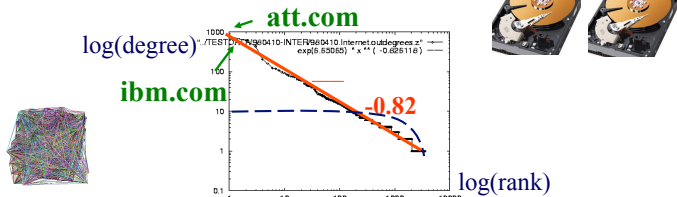
(c) C. Faloutsos, 2016

20

Carnegie Mellon

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

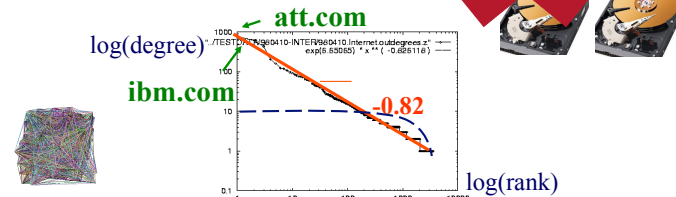


15-826 (c) C. Faloutsos, 2016 21

Carnegie Mellon

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



15-826 (c) C. Faloutsos, 2016 22

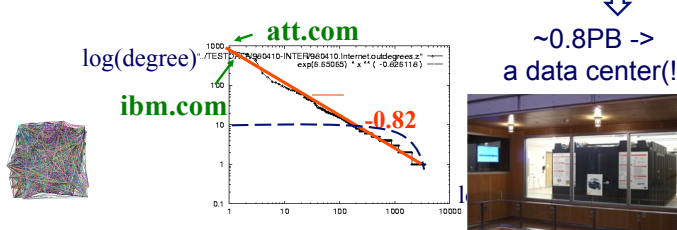
Carnegie Mellon

Solution# S.1

Gaussian trap

- Q: So what? = friends of friends (F.O.F.)
- A1: # of two-step-away pairs: $O(d_{\max}^2) \sim 10M^2$ internet domains

↓
~0.8PB -> a data center(!)



15-826 (c) C. Faloutsos, 2016 DCO @ CMU 23

Carnegie Mellon

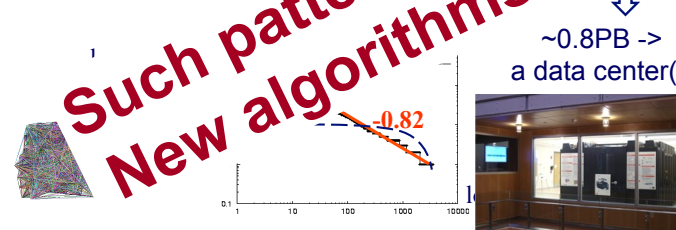
Solution# S.1

Gaussian trap

- Q: So what?
- A1: # of two-step-away pairs: $O(d_{\max}^2) \sim 10M^2$ internet domains

↓
~0.8PB -> a data center(!)

Such patterns -> New algorithms

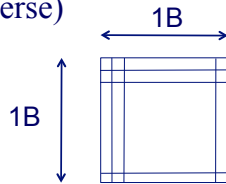


15-826 (c) C. Faloutsos, 2016 24

Carnegie Mellon

Observation – big-data:

- $O(N^2)$ algorithms are ~intractable - $N=1B$
- N^2 seconds = 31B years (>2x age of universe)



15-826

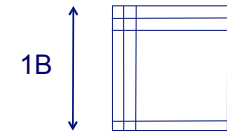
(c) C. Faloutsos, 2016

25

Carnegie Mellon

Observation – big-data:

- $O(N^2)$ algorithms are ~intractable - $N=1B$
- N^2 seconds = ~~31B~~^{31M} years
- 1,000 machines



15-826

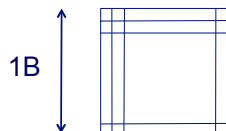
(c) C. Faloutsos, 2016

26

Carnegie Mellon

Observation – big-data:

- $O(N^2)$ algorithms are ~intractable - $N=1B$
- N^2 seconds = ~~31B~~^{31K} years
- 1M machines



Google Y!

15-826

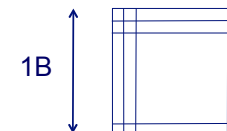
(c) C. Faloutsos, 2016

27

Carnegie Mellon

Observation – big-data:

- $O(N^2)$ algorithms are ~intractable - $N=1B$
- N^2 seconds = ~~31B~~³ years
- 10B machines ~ \$10Trillion



15-826

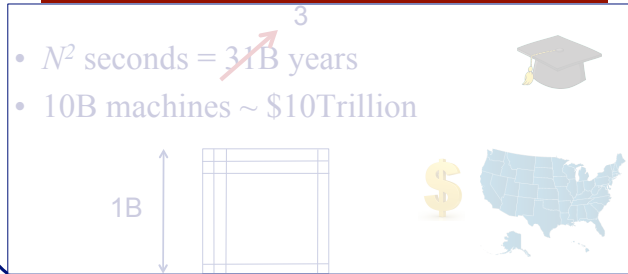
(c) C. Faloutsos, 2016

28

Observation – big-data:

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

And parallelism might not help

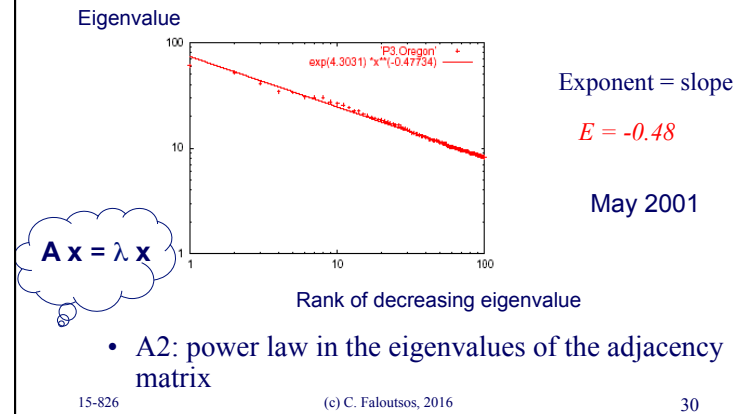


15-826

(c) C. Faloutsos, 2016

29

Solution# S.2: Eigen Exponent E

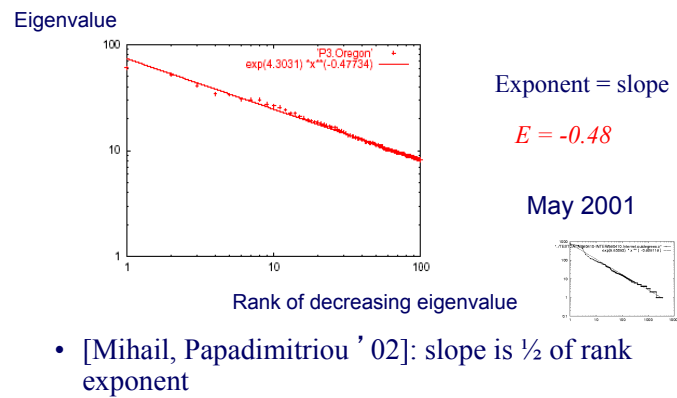


15-826

(c) C. Faloutsos, 2016

30

Solution# S.2: Eigen Exponent E



15-826

(c) C. Faloutsos, 2016

31

But:

How about graphs from other domains?

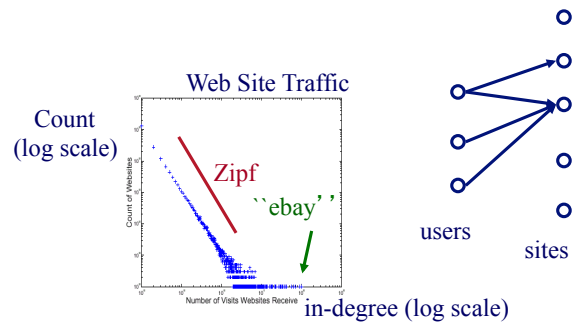
15-826

(c) C. Faloutsos, 2016

32

More power laws:

- web hit counts [w/ A. Montgomery]



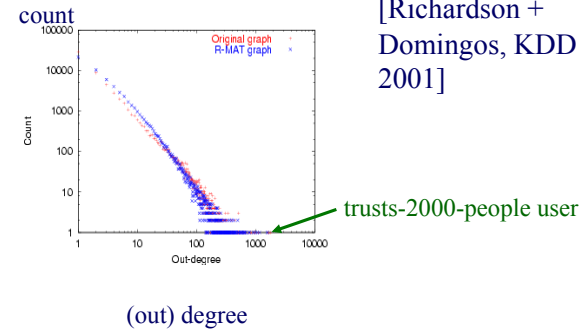
15-826

(c) C. Faloutsos, 2016

33

epinions.com

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



15-826

(c) C. Faloutsos, 2016

34

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) C. Faloutsos, 2016

35

Outline



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



15-826

(c) C. Faloutsos, 2016

36

Solution# S.3: Triangle 'Laws'



- Real social networks have a lot of triangles

15-826

(c) C. Faloutsos, 2016

37

Solution# S.3: Triangle 'Laws'



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

15-826

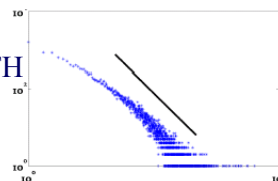
(c) C. Faloutsos, 2016

38

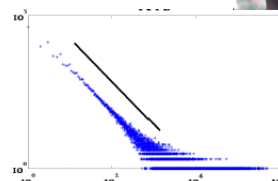
Triangle Law: #S.3 [Tsourakakis ICDM 2008]



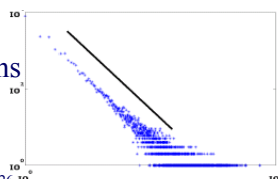
HEP-TH



ASN



Epinions



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



15-826

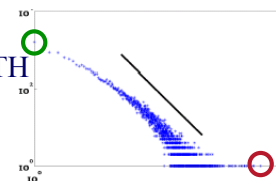
tsos, 2016

39

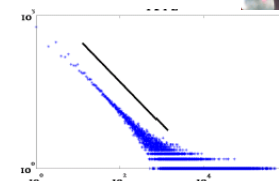
Triangle Law: #S.3 [Tsourakakis ICDM 2008]



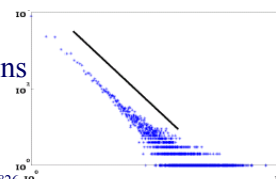
HEP-TH



ASN



Epinions



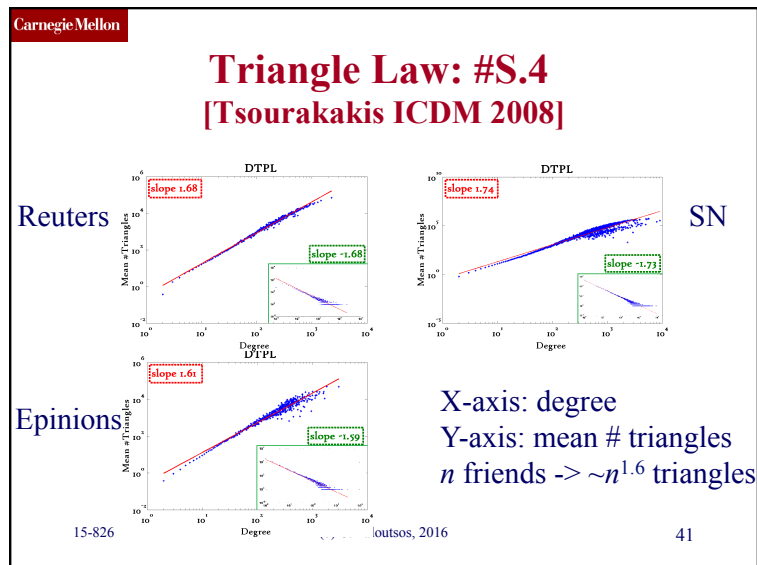
X-axis: # of participating triangles
Y: count (\sim pdf)



15-826

tsos, 2016

40



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?

15-826

(c) C. Faloutsos, 2016

42

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?

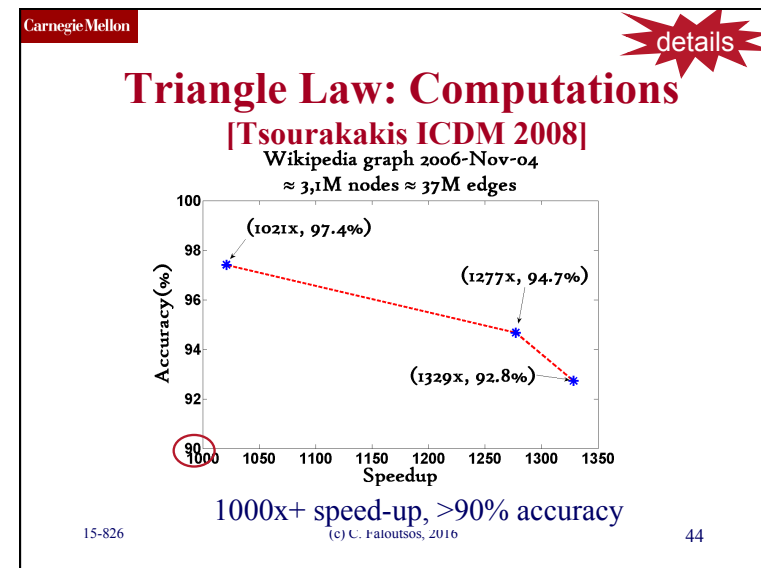
A: Yes!

#triangles = $1/6 \sum (\lambda_i^3)$
(and, because of skewness (S2),
we only need the top few eigenvalues!)

15-826

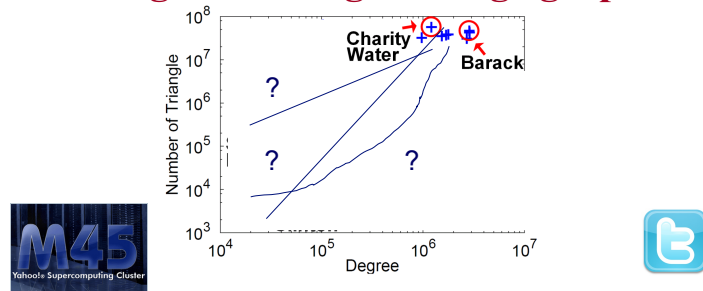
(c) C. Faloutsos, 2016

43



Carnegie Mellon

Triangle counting for large graphs?



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

15-826

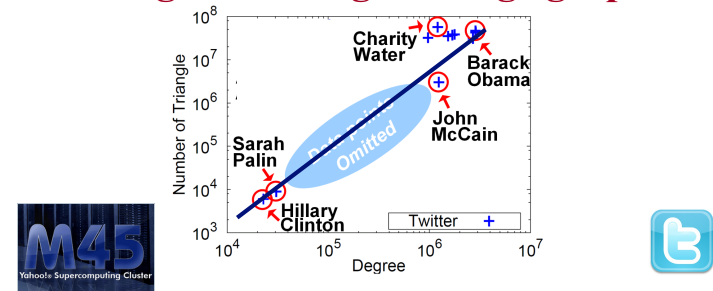


(c) C. Faloutsos, 2016

45

Carnegie Mellon

Triangle counting for large graphs?



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

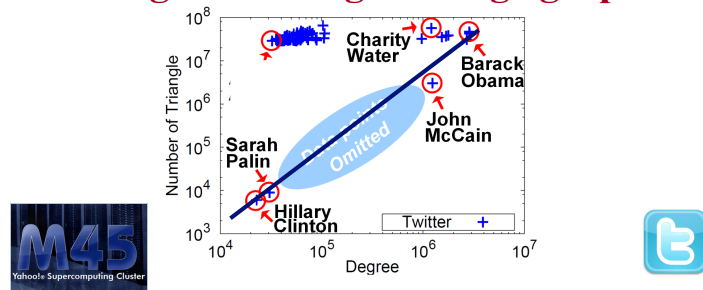
15-826

(c) C. Faloutsos, 2016

46

Carnegie Mellon

Triangle counting for large graphs?



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

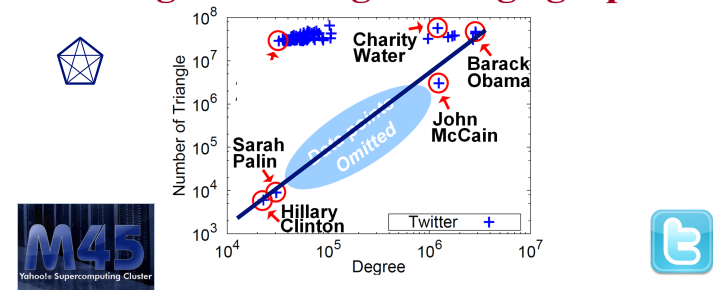
15-826

(c) C. Faloutsos, 2016

47

Carnegie Mellon

Triangle counting for large graphs?



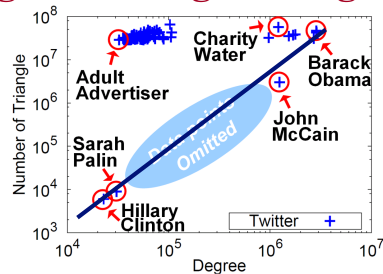
Anomalous nodes in Twitter (~ 3 billion edges)
[U Kang, Brendan Meeder, +, PAKDD'11]

15-826

(c) C. Faloutsos, 2016

48

Triangle counting for large graphs?



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

15-826

(c) C. Faloutsos, 2016

49

Any other 'laws'?

Yes!

15-826

(c) C. Faloutsos, 2016

50

Any other 'laws'?

Yes!

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

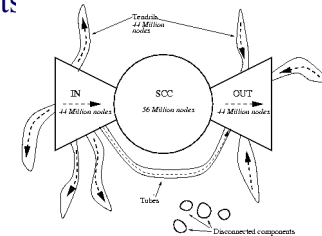
15-826

(c) C. Faloutsos, 2016

51

Any other 'laws'?

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



15-826

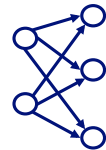
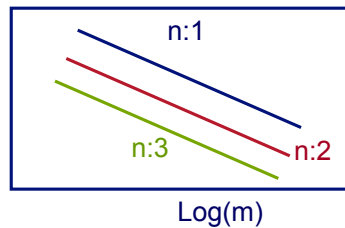
(c) C. Faloutsos, 2016

52

Any other 'laws' ?

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

Log(count)



2:3 core
(m:n core)

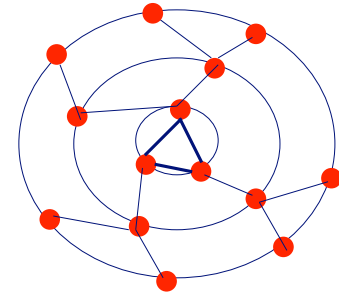
15-826

(c) C. Faloutsos, 2016

53

Any other 'laws' ?

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



15-826

(c) C. Faloutsos, 2016

54

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

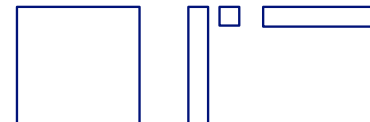
(c) C. Faloutsos, 2016

55

EigenSpokes

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

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



15-826

(c) C. Faloutsos, 2016

56

Carnegie Mellon

EigenSpokes details

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

$A = U\Sigma U^T$

15-826 (c) C. Faloutsos, 2016 57

Carnegie Mellon

EigenSpokes details

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

$A = U\Sigma U^T$

15-826 (c) C. Faloutsos, 2016 58

Carnegie Mellon

EigenSpokes details

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

$A = U\Sigma U^T$

15-826 (c) C. Faloutsos, 2016 59

Carnegie Mellon

EigenSpokes details

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

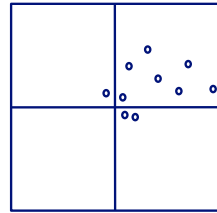
$A = U\Sigma U^T$

15-826 (c) C. Faloutsos, 2016 60

EigenSpokes

- EE plot:
- Scatter plot of scores of u_1 vs u_2
- One would expect
 - Many points @ origin
 - A few scattered ~randomly

2nd Principal component
 u_2



u_1
1st Principal component

15-826

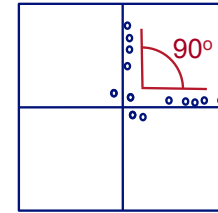
(c) C. Faloutsos, 2016

61

EigenSpokes

- EE plot:
- Scatter plot of scores of u_1 vs u_2
- One would expect
 - Many points @ origin
 - A few scattered ~randomly

u_2



u_1

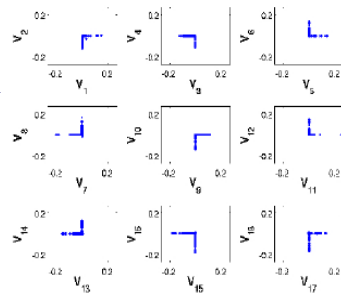
15-826

(c) C. Faloutsos, 2016

62

EigenSpokes - pervasiveness

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



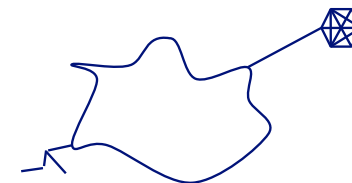
15-826

(c) C. Faloutsos, 2016

63

EigenSpokes - explanation

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



15-826

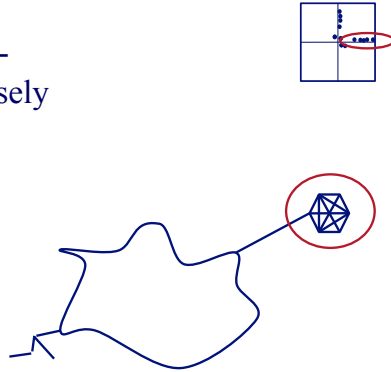
(c) C. Faloutsos, 2016

64

Carnegie Mellon

EigenSpokes - explanation

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



15-826

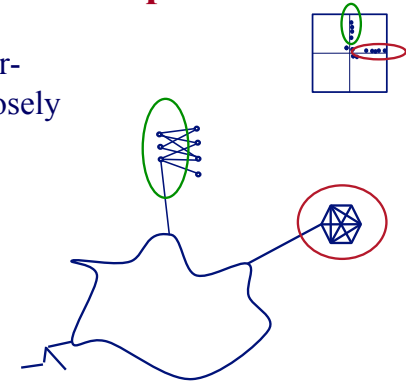
(c) C. Faloutsos, 2016

65

Carnegie Mellon

EigenSpokes - explanation

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



15-826

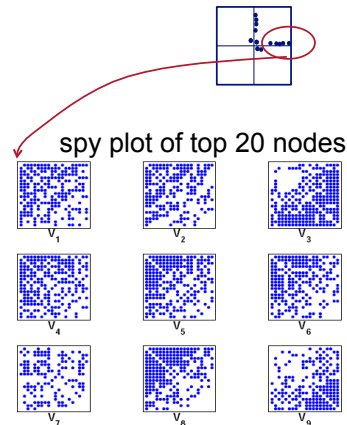
(c) C. Faloutsos, 2016

66

Carnegie Mellon

EigenSpokes - explanation

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



So what?

- Extract nodes with high scores
- high connectivity
- Good “communities”

15-826

(c) C. Faloutsos, 2016

67

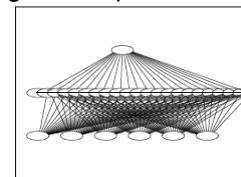
Carnegie Mellon

Bipartite Communities!

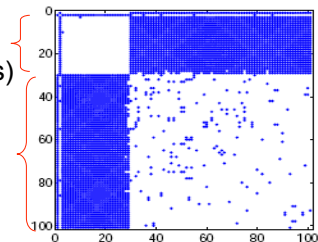
patents from same inventor(s)

‘cut-and-paste’ bibliography!

magnified bipartite community



Useful for fraud detection!

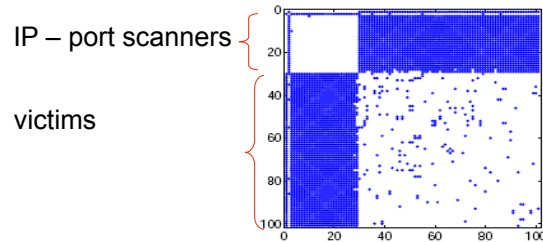


15-826

(c) C. Faloutsos, 2016

68

Bipartite Communities!



Useful for fraud detection!

15-826

(c) C. Faloutsos, 2016

69

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, 2016

70

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, 2016

71

Observation W.1: Fortification

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

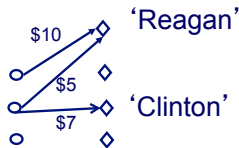
15-826

(c) C. Faloutsos, 2016

72

Observation W.1: Fortification

More donors,
more \$?



15-826

(c) C. Faloutsos, 2016

73

Observation W.1: fortification: Snapshot Power Law

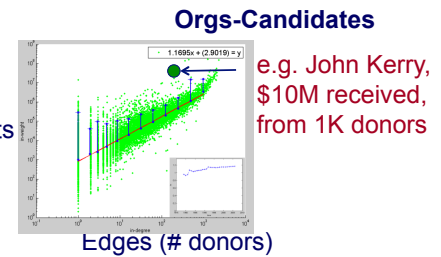
- Weight: super-linear on in-degree
- exponent 'iw': $1.01 < iw < 1.26$

More donors,
even more \$



15-826

In-weights
(\$)



(c) C. Faloutsos, 2016

74

Outline



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

15-826

(c) C. Faloutsos, 2016

75

Problem: Time evolution

- with Jure Leskovec (CMU -> Stanford)
- and Jon Kleinberg (Cornell – sabb. @ CMU)



15-826

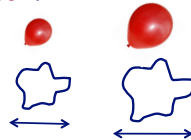
(c) C. Faloutsos, 2016

76

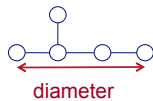
T.1 Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:

- [diameter $\sim O(N^{1/3})$]
- diameter $\sim O(\log N)$
- diameter $\sim O(\log \log N)$



- What is happening in real data?



15-826

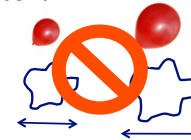
(c) C. Faloutsos, 2016

77

T.1 Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:

- [diameter $\sim O(N^{1/3})$]
- diameter $\sim O(\log N)$
- diameter $\sim O(\log \log N)$



- What is happening in real data?
- Diameter **shrinks** over time

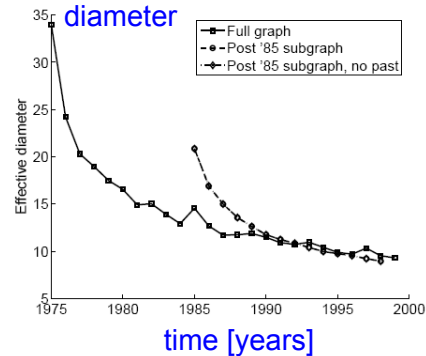
15-826

(c) C. Faloutsos, 2016

78

T.1 Diameter – “Patents”

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



15-826

(c) C. Faloutsos, 2016

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

$$E(t+1) = ? * E(t)$$

15-826

(c) C. Faloutsos, 2016

80

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
 - $E(t+1) = ? * E(t)$
- A: over-doubled!
 - But obeying the “Densification Power Law”

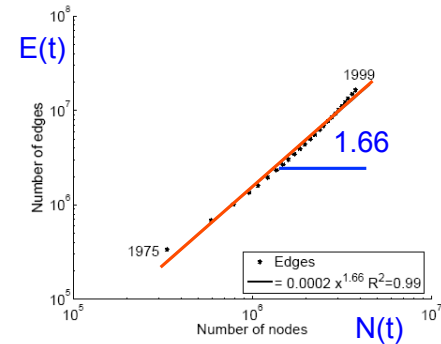
15-826

(c) C. Faloutsos, 2016

81

T.2 Densification – Patent Citations

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



15-826

(c) C. Faloutsos, 2016

82

Outline



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

15-826

(c) C. Faloutsos, 2016

83

More on Time-evolving graphs

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

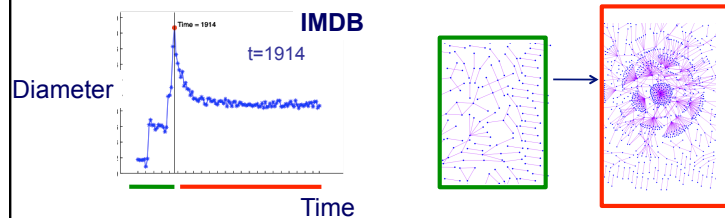
15-826

(c) C. Faloutsos, 2016

84

[Gelling Point]

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



15-826

(c) C. Faloutsos, 2016

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?



15-826

(c) C. Faloutsos, 2016

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 shrink?
- or stabilize?



15-826

(c) C. Faloutsos, 2016

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?

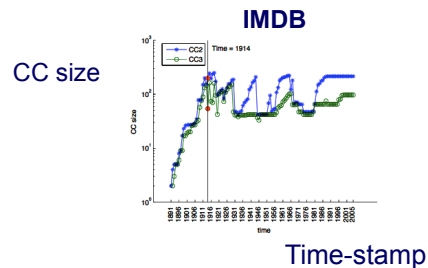
15-826

(c) C. Faloutsos, 2016

88

Observation T.3: NLCC behavior

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



15-826

(c) C. Faloutsos, 2016

89

Timing for Blogs

- with Mary McGlohon (CMU->Google)
 - Jure Leskovec (CMU->Stanford)
 - Natalie Glance (now at Google)
 - Mat Hurst (now at MSR)
- [SDM' 07]

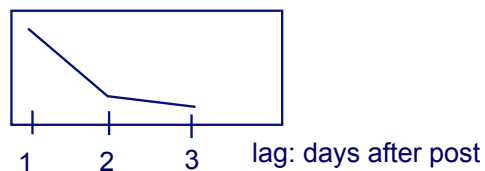
15-826

(c) C. Faloutsos, 2016

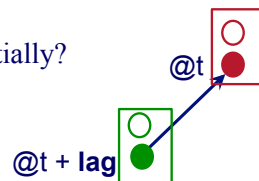
90

T.4 : popularity over time

in links



Post popularity drops-off – exponentially?

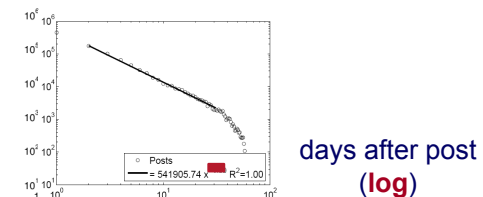


15-826

(c) C. Faloutsos, 2016

91

T.4 : popularity over time

in links
(log)

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

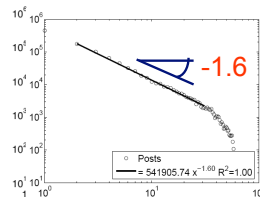
15-826

(c) C. Faloutsos, 2016

92

T.4 : popularity over time

in links
(log)

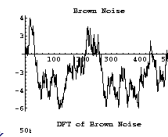


days after post
(log)

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) C. Faloutsos, 2016

93

-1.5 slope

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

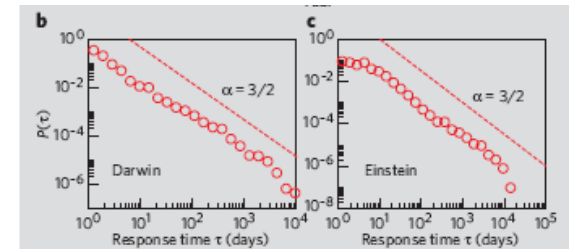


Figure 1 | The correspondence patterns of Darwin and Einstein.

1

94

T.5: duration of phonecalls

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



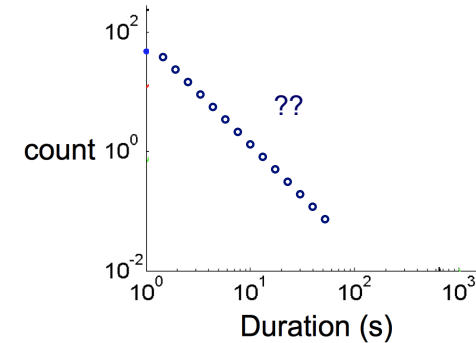
Pedro O. S. Vaz de Melo, Leman
Akoglu, Christos Faloutsos, Antonio
A. F. Loureiro
PKDD 2010

15-826

(c) C. Faloutsos, 2016

95

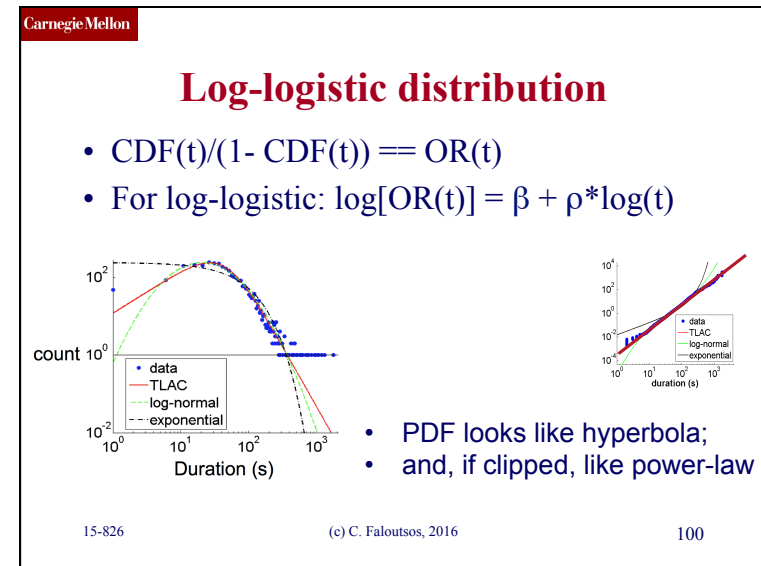
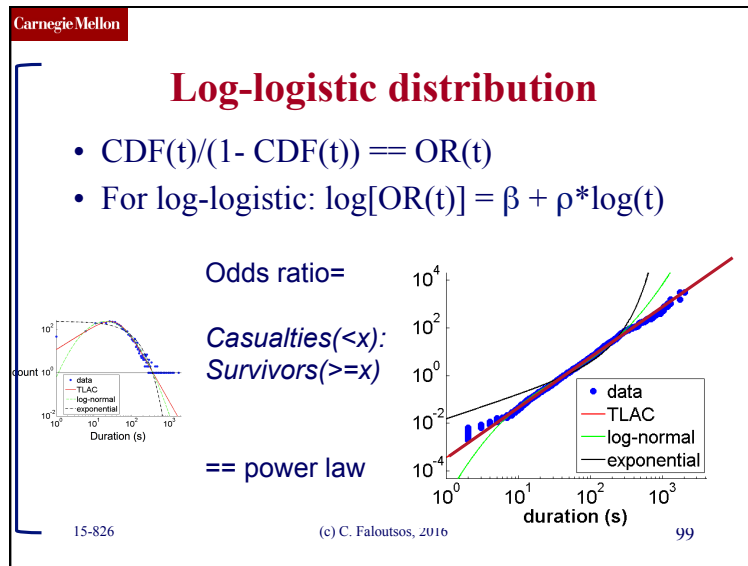
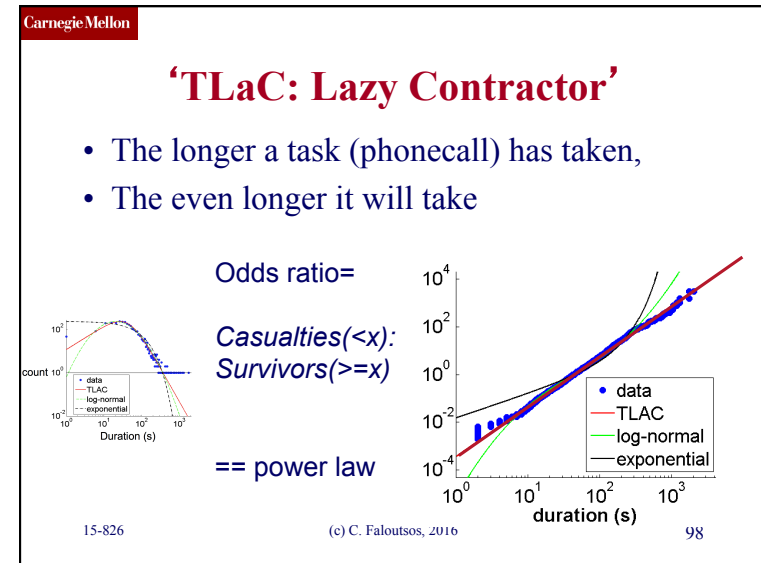
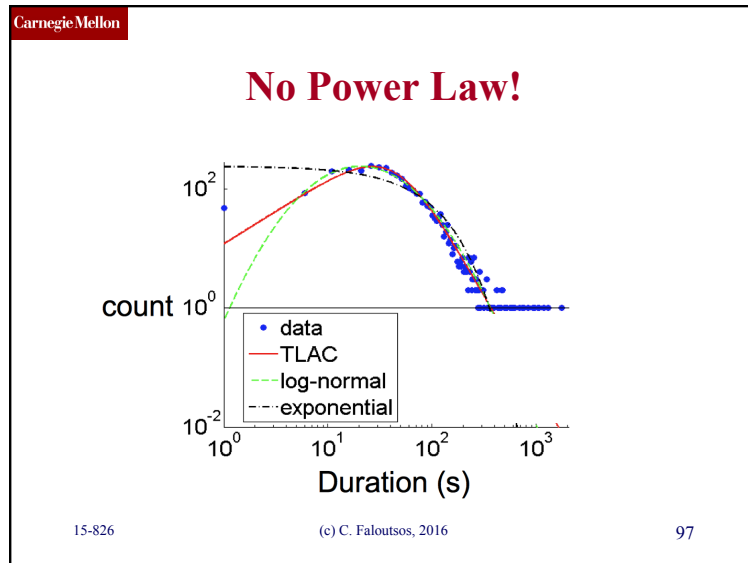
Probably, power law (?)



15-826

(c) C. Faloutsos, 2016

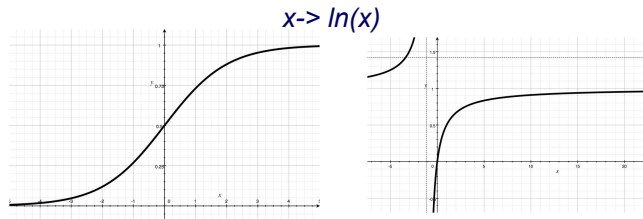
96



Carnegie Mellon

Log-logistic distribution

- Logistic distribution: CDF \rightarrow sigmoid
- LOG-Logistic distribution:



$$\text{CDF}(x) = 1/(1+\exp(-x))$$

15-826

(c) C. Faloutsos, 2016

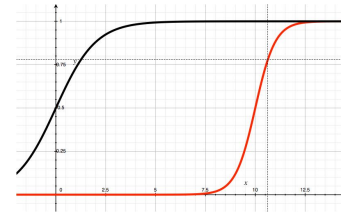
101

$$\text{CDF}(x) = 1/(1+1/x)$$

Carnegie Mellon

Log-logistic distribution

- Logistic distribution: CDF \rightarrow sigmoid
- LOG-Logistic distribution:



$$\text{CDF}(x) = 1/(1+\exp(-(x-m)/s)) \quad \text{CDF}(x) = 1/(1+\exp(-(\ln(x)-m)/s))$$

15-826

(c) C. Faloutsos, 2016

102

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)

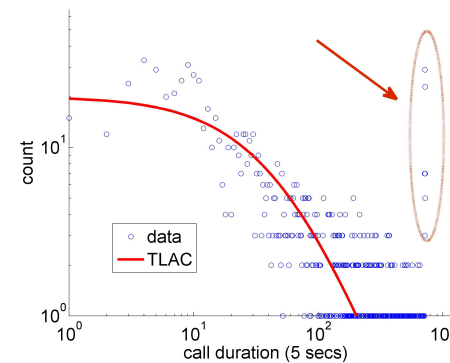
15-826

(c) C. Faloutsos, 2016

103

Carnegie Mellon

Outliers:



15-826

(c) C. Faloutsos, 2016

104

Outline



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

15-826

(c) C. Faloutsos, 2016

105

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, 2016

106

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) C. Faloutsos, 2016

107

HADI for diameter estimation

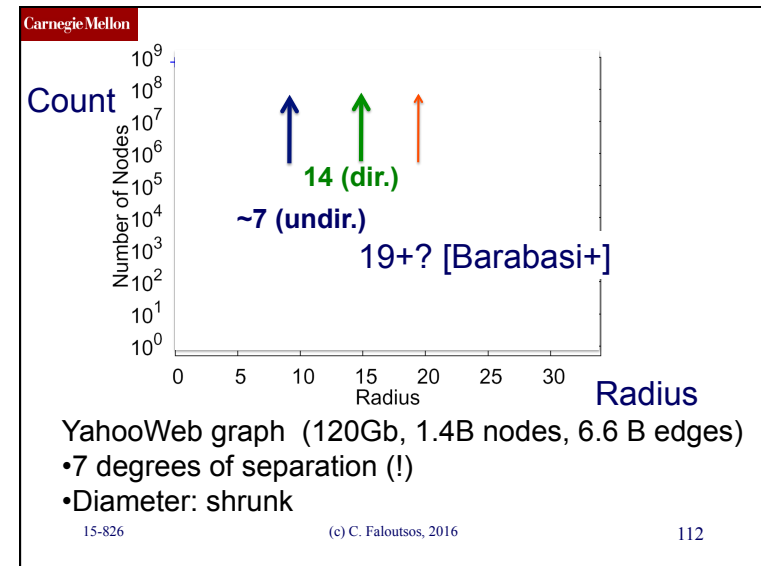
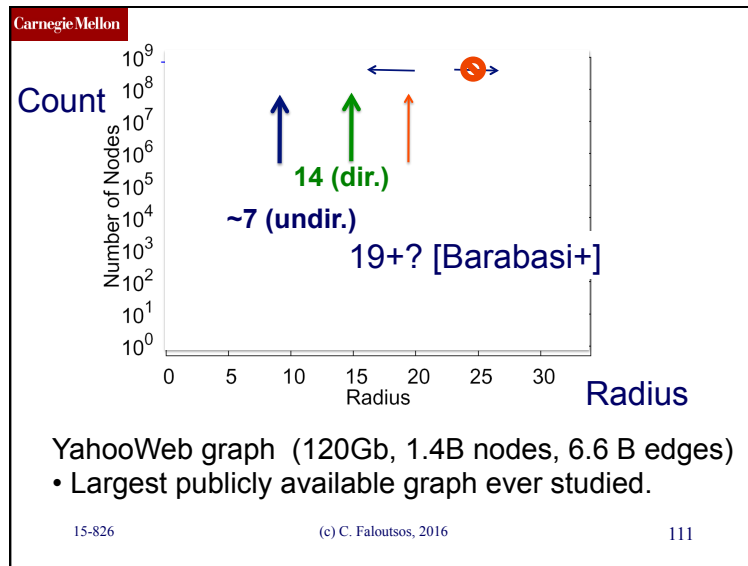
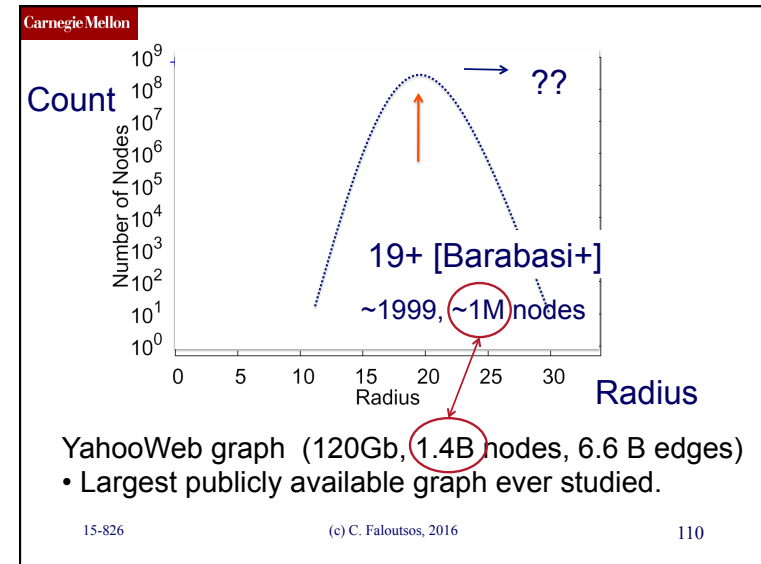
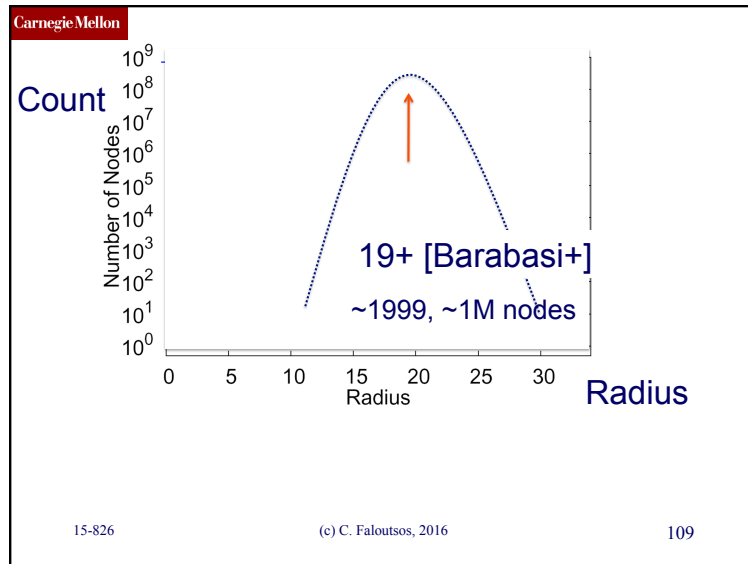


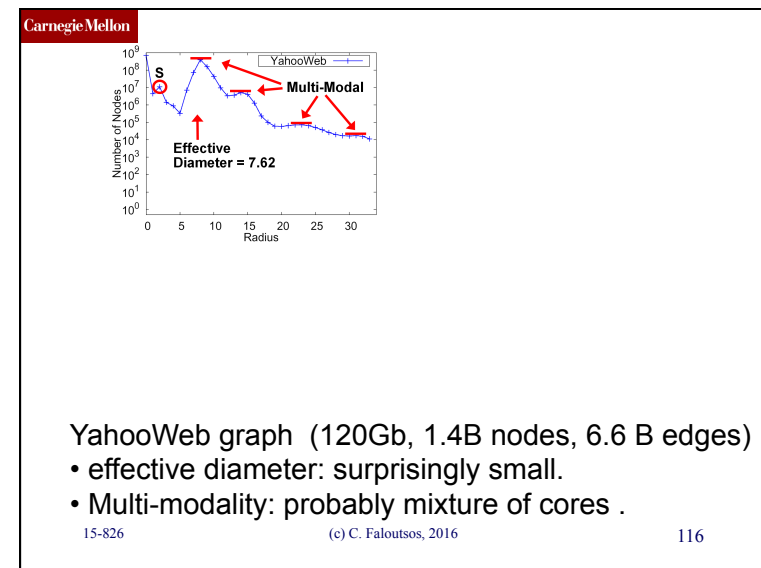
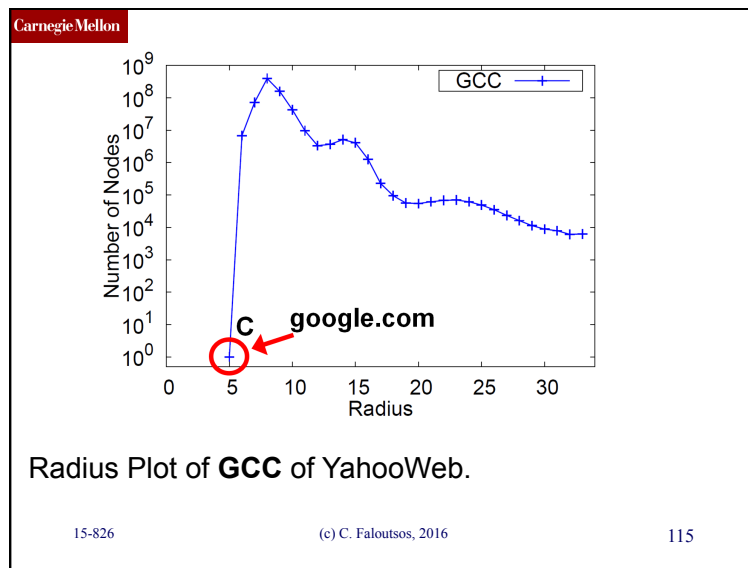
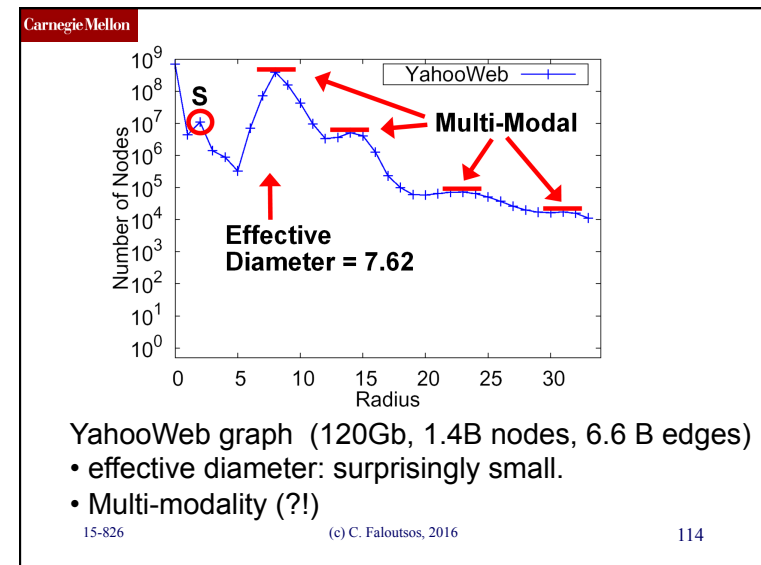
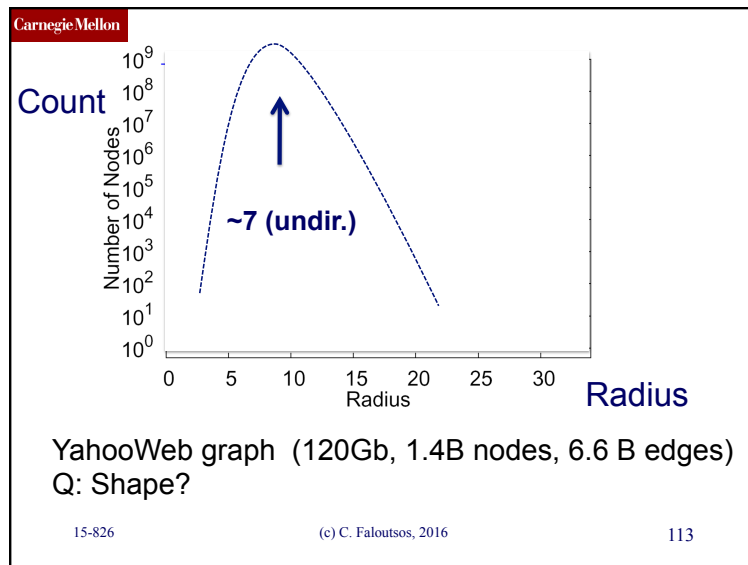
- *Radius Plots for Mining Tera-byte Scale Graphs* **U Kang**, Charalampos Tsourakakis, Ana Paula Appel, Christos Faloutsos, Jure Leskovec, SDM’10
- Naively: diameter needs $O(N^2)$ space and up to $O(N^3)$ time – **prohibitive** ($N \sim 1B$)
- Our HADI: linear on E ($\sim 10B$)
 - Near-linear scalability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

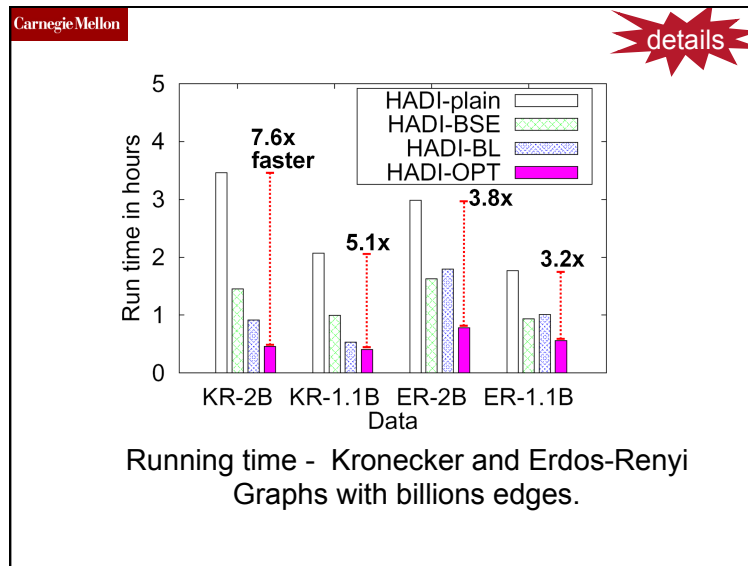
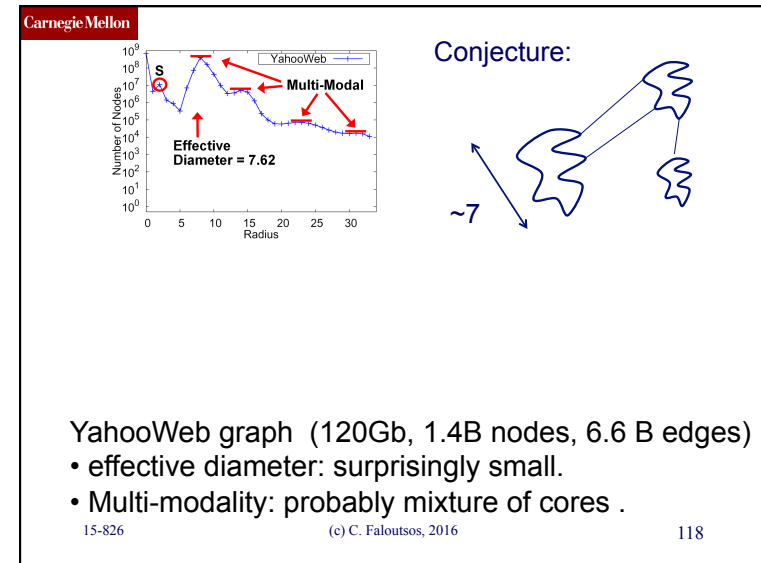
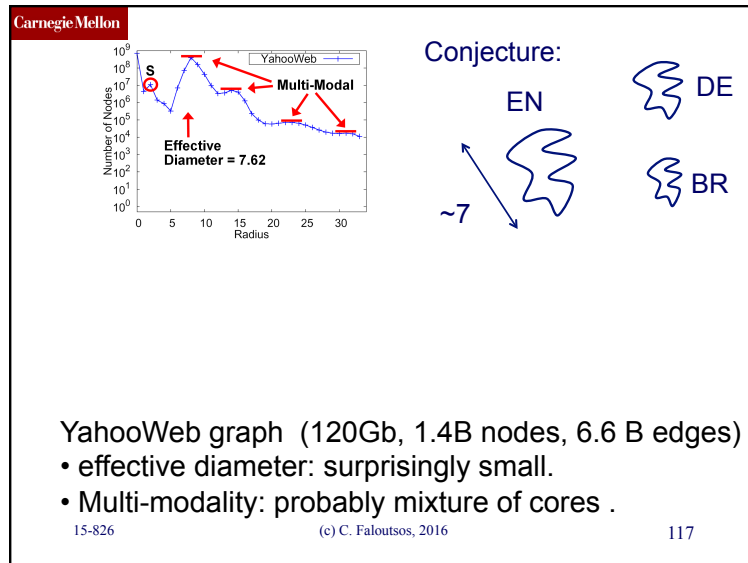
15-826

(c) C. Faloutsos, 2016

108







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) C. Faloutsos, 2016 120

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, 2016

121

Generalized Iterated Matrix Vector Multiplication (GIMV)

details

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

Matrix – vector
Multiplication
(iterated)

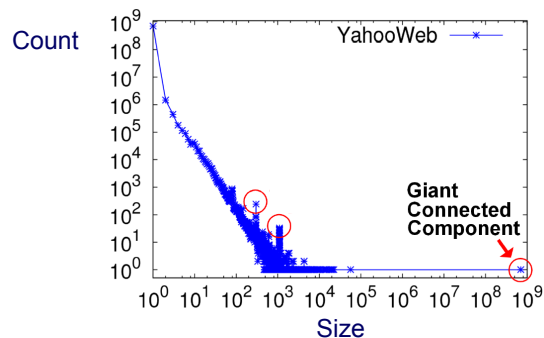
15-826

(c) C. Faloutsos, 2016

122

Example: GIM-V At Work

- Connected Components – 4 observations:



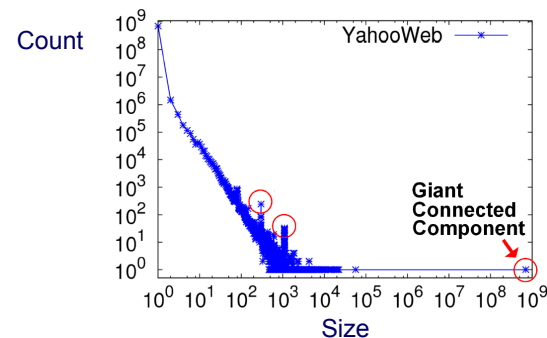
15-826

(c) C. Faloutsos, 2016

123

Example: GIM-V At Work

- Connected Components



1) 10K x
larger
than next

15-826

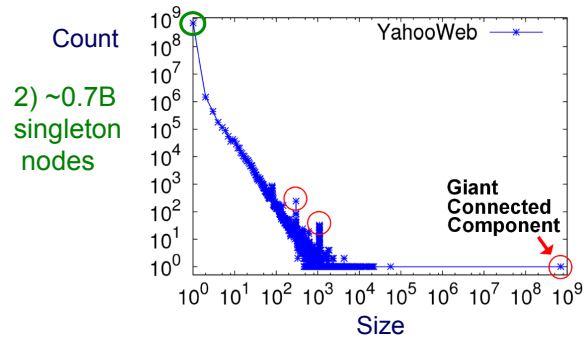
(c) C. Faloutsos, 2016

124

Carnegie Mellon

Example: GIM-V At Work

- Connected Components



15-826

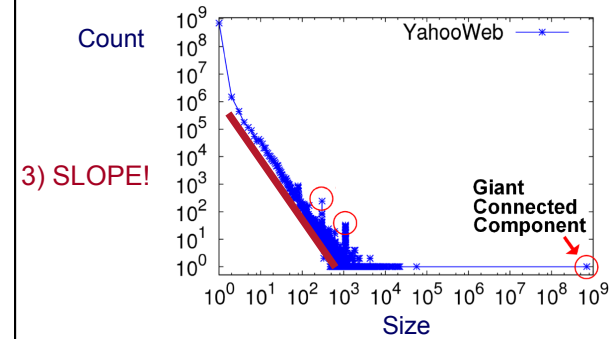
(c) C. Faloutsos, 2016

125

Carnegie Mellon

Example: GIM-V At Work

- Connected Components



15-826

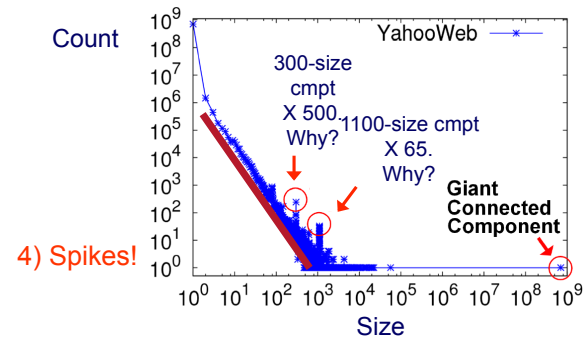
(c) C. Faloutsos, 2016

126

Carnegie Mellon

Example: GIM-V At Work

- Connected Components



4) Spikes!

15-826

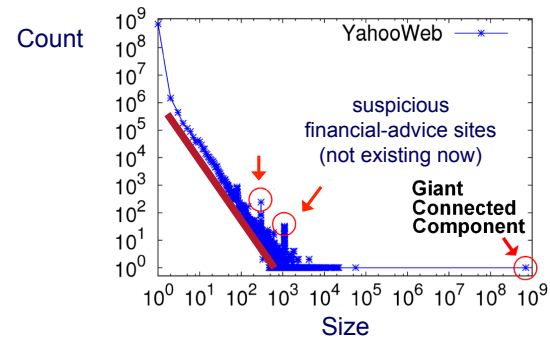
(c) C. Faloutsos, 2016

127

Carnegie Mellon

Example: GIM-V At Work

- Connected Components



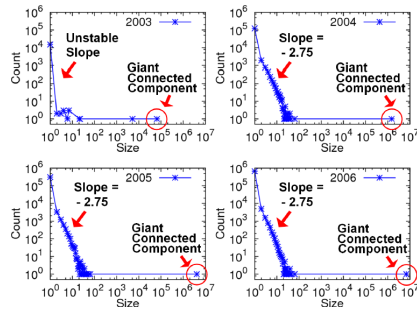
15-826

(c) C. Faloutsos, 2016

128

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, 2016

129

Outline



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

15-826

(c) C. Faloutsos, 2016

130

OVERALL CONCLUSIONS – low level:

- Several new **patterns** (fortification, shrinking diameter, triangle-laws, conn. components, etc)
- Log-logistic distribution: ubiquitous
- New **tools**:
 - anomaly detection (OddBall), belief propagation, immunization
- **Scalability**: PEGASUS / hadoop

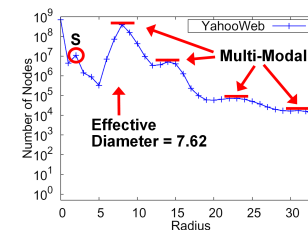
15-826

(c) C. Faloutsos, 2016

131

OVERALL CONCLUSIONS – high level

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



15-826

(c) C. Faloutsos, 2016

132

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, 2016

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, 2016

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, 2016

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, 2016

136

References

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

15-826

(c) C. Faloutsos, 2016

137

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, 2016

138

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, 2016

139

(Project info)

www.cs.cmu.edu/~pegasus



Chau,
Polo



Akoglu,
Leman

Koutra,
Danae



Kang, U

Prakash,
Aditya



McGlohon,
Mary



Tong,
Hanghang

15-826

(c) C. Faloutsos, 2016

140