

# 15-826: Multimedia Databases and Data Mining

(Project lecture #1)

Lecture #26: Graph mining - patterns

*Christos Faloutsos*

1

## Must-read Material – 1-of-2

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

2

CarnegieMellon

## Must-read Material 2-of-2

- 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, 2019 3

3

CarnegieMellon


## Problem

- Are real graphs random?



15-826 Copyright: C. Faloutsos (2019) 4

4


CarnegieMellon 

## Conclusions

- Are real graphs random?
- NO!
  - Static patterns
    - Small diameters
    - Skewed degree distribution
    - Shrinking diameters
  - Weighted
  - Time-evolving

15-826 Copyright: C. Faloutsos (2019) 5

5

CarnegieMellon 

## Conclusions

- Are real graphs random?
- NO!
  - Static patterns
    - Small diameters
    - Skewed degree distribution
  - Weighted
  - Time-evolving


**• Many power laws – log-logistic**  
**• Take logarithms**

15-826 Copyright: C. Faloutsos (2019) 6

6

CarnegieMellon

## Main outline




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

15-826 (c) C. Faloutsos, 2019 7

7

CarnegieMellon

## Outline




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

15-826 (c) C. Faloutsos, 2019 8

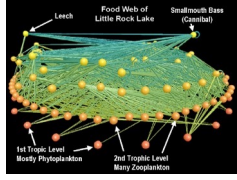
8

CarnegieMellon

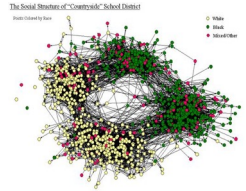
## Graphs - why should we care?



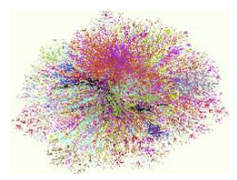
15-826



Food Web  
[Martinez '91]



Friendship Network  
[Moody '01]



Internet Map  
[lumeta.com]

(c) C. Faloutsos, 2019

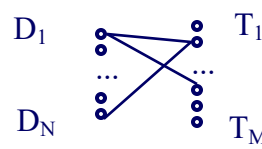
9

9

CarnegieMellon

## Graphs - why should we care?

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



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

(c) C. Faloutsos, 2019

10

10

CarnegieMellon

## Graphs - why should we care?


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

15-826 (c) C. Faloutsos, 2019 11

11

CarnegieMellon

## Outline



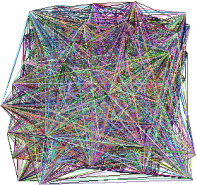
- Introduction – Motivation
- ➔ • Problem: Patterns in graphs
  - Static graphs
  - Weighted graphs
  - Time evolving graphs
- Problem#2: Scalability
- Conclusions

15-826 (c) C. Faloutsos, 2019 12

12

CarnegieMellon

## Problem #1 - network and graph mining



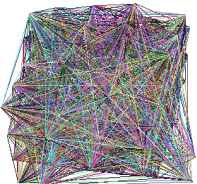
- What does the Internet look like?
- What does FaceBook look like?
- What is 'normal' / 'abnormal' ?
- which patterns/laws hold?

15-826 (c) C. Faloutsos, 2019 13


13

CarnegieMellon

## 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?
  - **anomalies** (rarities) <-> **patterns**

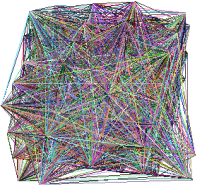


15-826 (c) C. Faloutsos, 2019 14

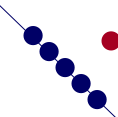

14

CarnegieMellon

## 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?
  - **anomalies** (rarities) <-> **patterns**
  - **Large** datasets reveal patterns/anomalies that may be invisible otherwise...

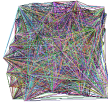



15-826 (c) C. Faloutsos, 2019 15

15

CarnegieMellon

## Graph mining



- Are real graphs random?

15-826 (c) C. Faloutsos, 2019 16

16



## Laws and patterns

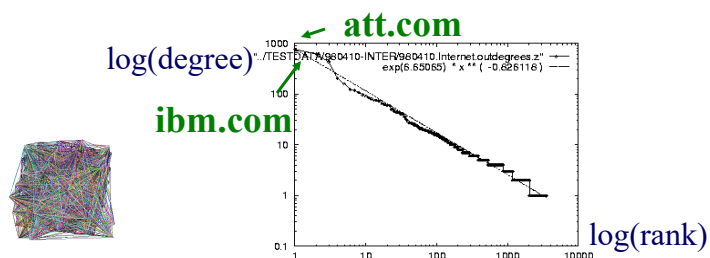
- Are real graphs random?
- A: NO!!
  - Diameter ( ‘6 degrees’ , ‘Kevin Bacon’ )
  - in- and out- degree distributions
  - other (surprising) patterns
- So, let’ s look at the data



## Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

internet domains



CarnegieMellon

## Solution# S.1

- Power law in the degree distribution [SIGCOMM99]

internet domains

15-826 (c) C. Faloutsos, 2019 19

19

CarnegieMellon

## Solution# S.1

- Q: So what?

internet domains

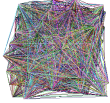
15-826 (c) C. Faloutsos, 2019 20

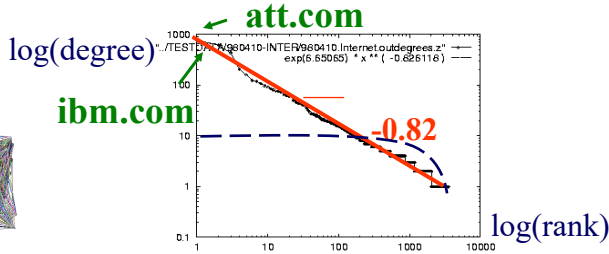
20

CarnegieMellon

## Solution# S.1

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





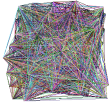
(c) C. Faloutsos, 2019

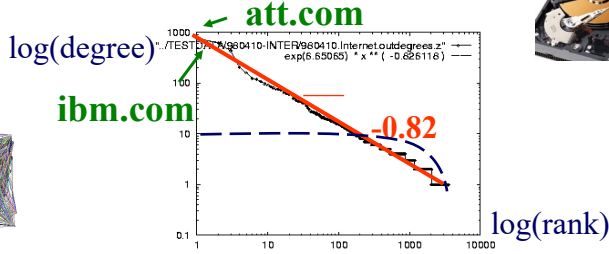
21


CarnegieMellon

## 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**







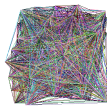
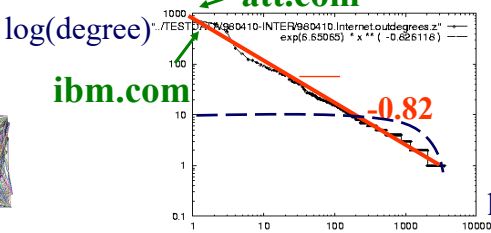

(c) C. Faloutsos, 2019

22

CarnegieMellon

## Solution# S.1

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

15-826 (c) C. Faloutsos, 2019 23

23

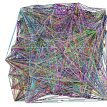
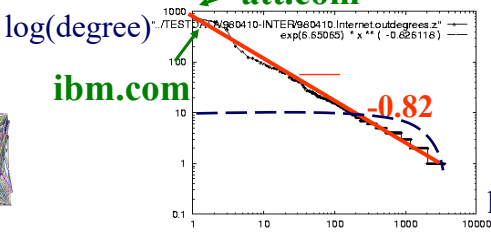

CarnegieMellon

## 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, 2019 DCO @ CMU 24

24

CarnegieMellon

## Gaussian trap

### Solution# S.1

- Q: So what?
- A1: # of two-step-away inter

Such patterns ->  
New algorithms

?)  $\sim 10M^2$   
↓  
 $\sim 0.8PB \rightarrow$   
a data center(!)

15-826 (c) C. Faloutsos, 2019 25

25

CarnegieMellon

## Observation – big-data:

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

1B

↑

1B

←

15-826 (c) C. Faloutsos, 2019 26

26

CarnegieMellon

## Observation – big-data:

- $O(N^2)$  algorithms are ~intractable -  $N=1B$
- $N^2$  seconds = ~~31B~~<sup>31M</sup> years
- 1,000 machines

1B

15-826

(c) C. Faloutsos, 2019

27

27

CarnegieMellon

## Observation – big-data:

- $O(N^2)$  algorithms are ~intractable -  $N=1B$
- $N^2$  seconds = ~~31B~~<sup>31K</sup> years
- 1M machines

1B

15-826

(c) C. Faloutsos, 2019

28

28

CarnegieMellon

### Observation – big-data:

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

15-826 (c) C. Faloutsos, 2019 29

29

CarnegieMellon

### Observation – big-data:

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

**And parallelism might not help**

- $N^2$  seconds = ~~31B~~<sup>3</sup> years
- 10B machines ~ \$10Trillion

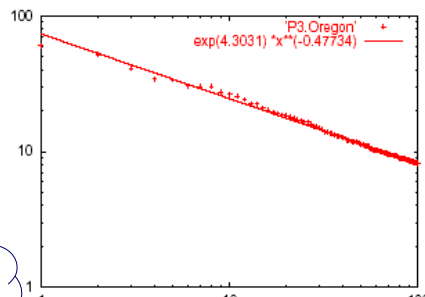
15-826 (c) C. Faloutsos, 2019 30

30

CarnegieMellon

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

Eigenvalue



Rank of decreasing eigenvalue

Exponent = slope

$E = -0.48$

May 2001

$Ax = \lambda x$

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

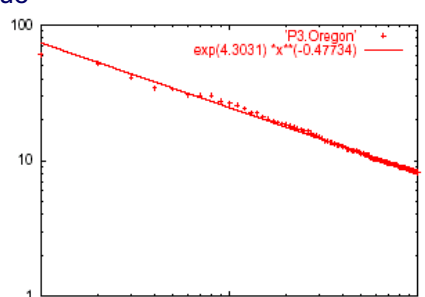
15-826
(c) C. Faloutsos, 2019
31

31

CarnegieMellon

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

Eigenvalue



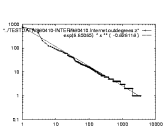
Rank of decreasing eigenvalue

Exponent = slope

$E = -0.48$

May 2001

- [Mihail, Papadimitriou '02]: slope is  $\frac{1}{2}$  of rank exponent



15-826
(c) C. Faloutsos, 2019
32

32



CarnegieMellon

**But:**

How about graphs from other domains?

15-826 (c) C. Faloutsos, 2019 33

33

CarnegieMellon

**More power laws:**

- web hit counts [w/ A. Montgomery]

Count  
(log scale)

Web Site Traffic

Count of Websites

Number of Visits Websites Receive

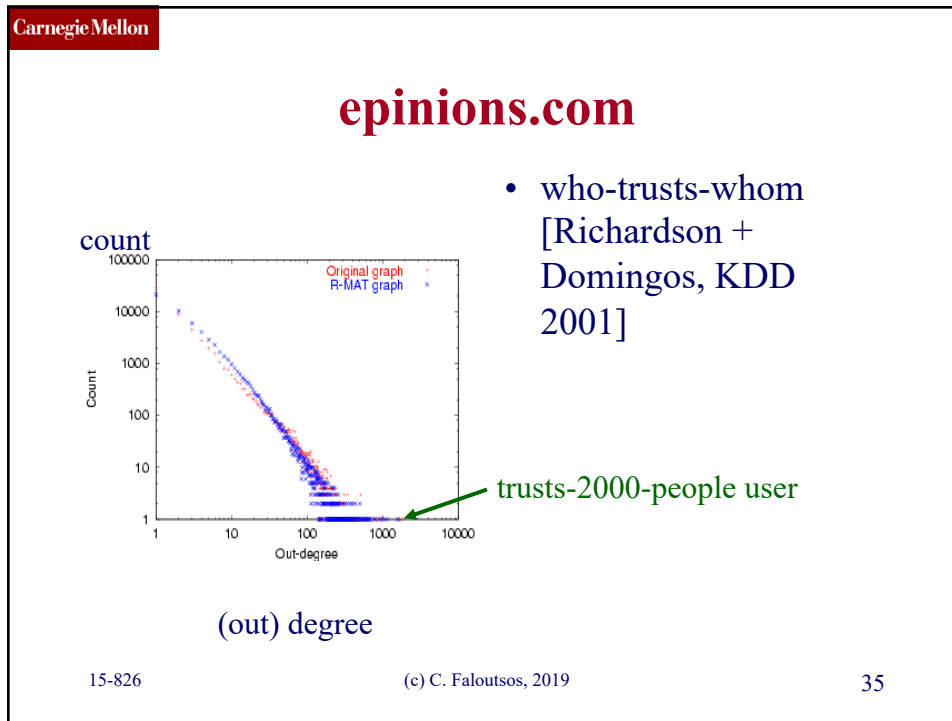
users

sites

in-degree (log scale)

15-826 (c) C. Faloutsos, 2019 34

34



35

CarnegieMellon

## 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, 2019 36

36

CarnegieMellon

## List of Static Patterns



- ✓ • S.1 degree
- ✓ • S.2 eigenvalues
- S.3 small diameter
- S.4/5 Triangle laws
- (S.6) NLCC non-largest conn. components
- (S.7) eigen plots
- (S.8) radius plot

} In textbook


15-826 (c) C. Faloutsos, 2019 37

37

CarnegieMellon

## S.3 small diameters

- Small diameter ( $\sim$  constant!) –
  - six degrees of separation / ‘Kevin Bacon’
  - small worlds [Watts and Strogatz]




15-826 (c) C. Faloutsos, 2019 38

38

CarnegieMellon

## List of Static Patterns



- ✓ • S.1 degree
- ✓ • S.2 eigenvalues
- ✓ • S.3 small diameter
- S.4/5 Triangle laws
- (S.6) NLCC non-largest conn. components
- (S.7) eigen plots
- (S.8) radius plot

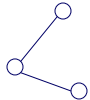
} In textbook

15-826 (c) C. Faloutsos, 2019 39

39

CarnegieMellon

## Solution# S.4: Triangle 'Laws'



- Real social networks have a lot of triangles

15-826 (c) C. Faloutsos, 2019 40

40

CarnegieMellon

## Solution# S.4: Triangle 'Laws'




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

15-826
(c) C. Faloutsos, 2019
41

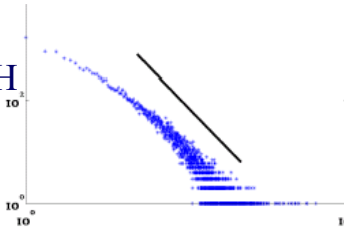
41

CarnegieMellon

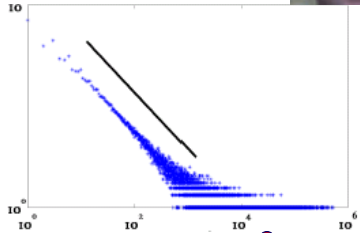
## Triangle Law: #S.4 [Tsourakakis ICDM 2008]



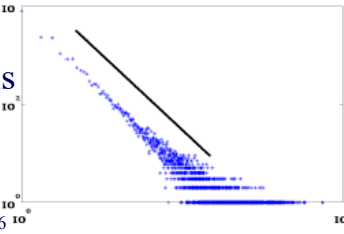
HEP-TH




ASN



Epinions

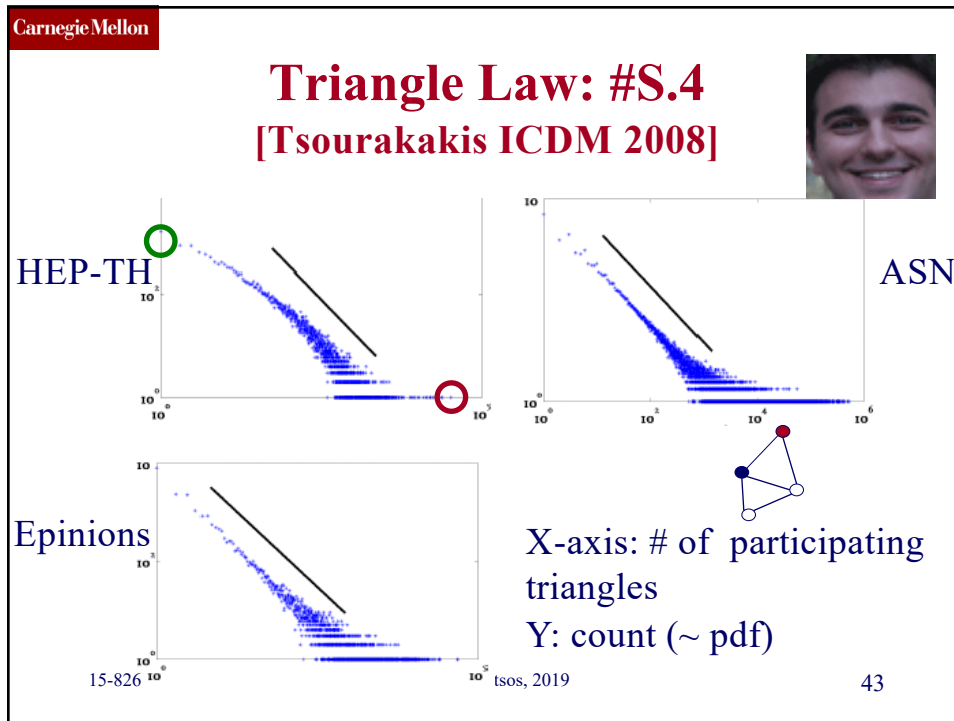




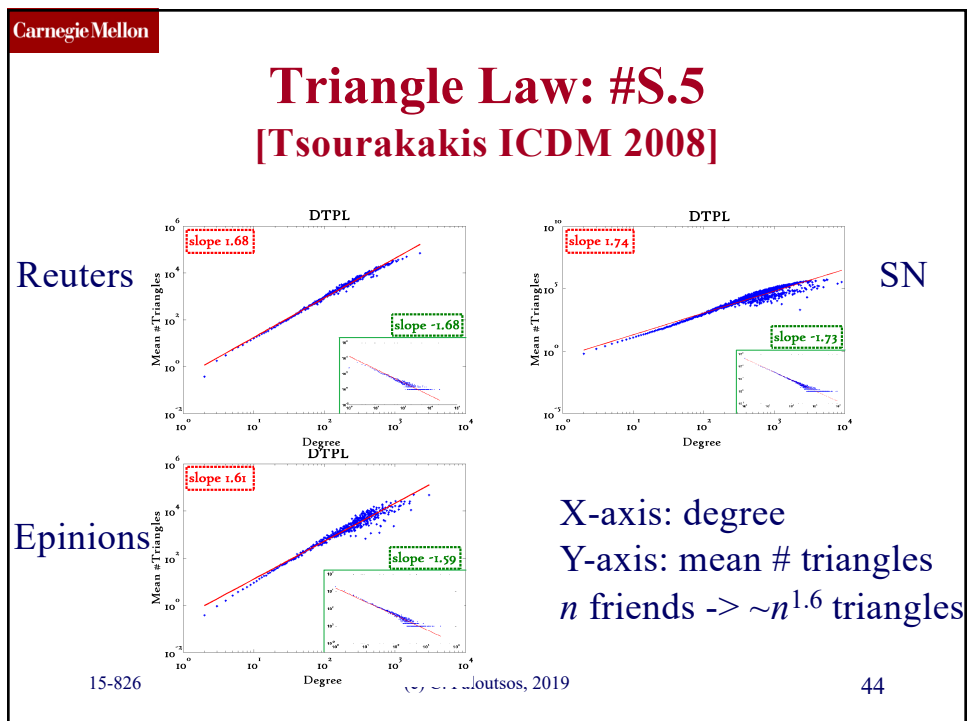
X-axis: # of participating triangles  
Y: count (~ pdf)

15-826
tsos, 2019
42


42



43



44

CarnegieMellon 

## 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, 2019 45

45

CarnegieMellon 

## Triangle Law: Computations

[Tsourakakis ICDM 2008]

But: triangles are expensive to compute  
(3-way join; several approx. algos)

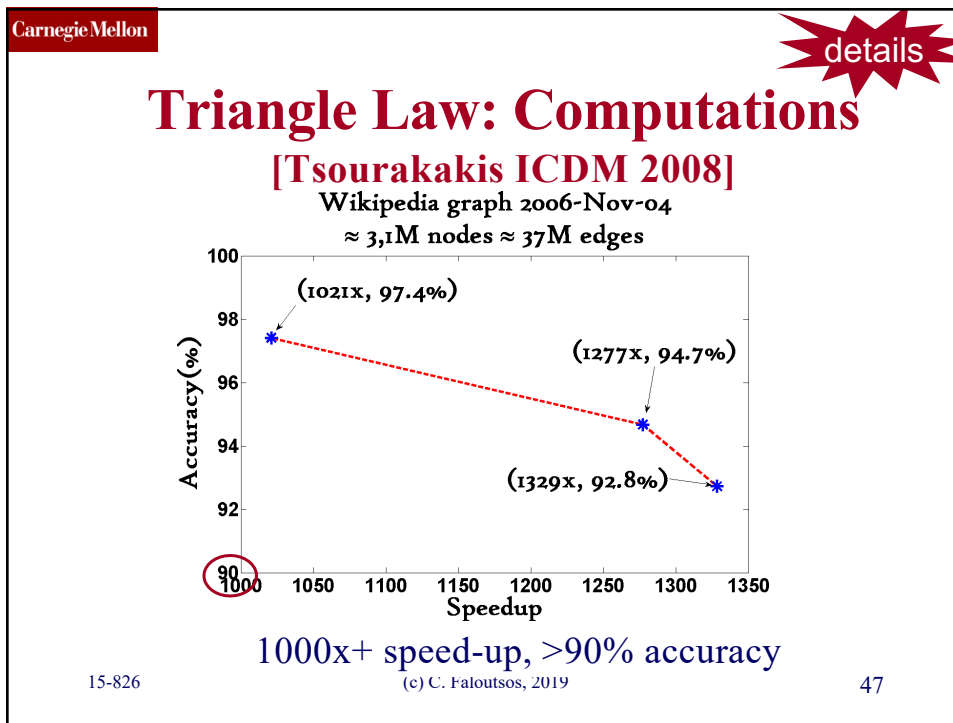
Q: Can we do that quickly?

A: Yes!

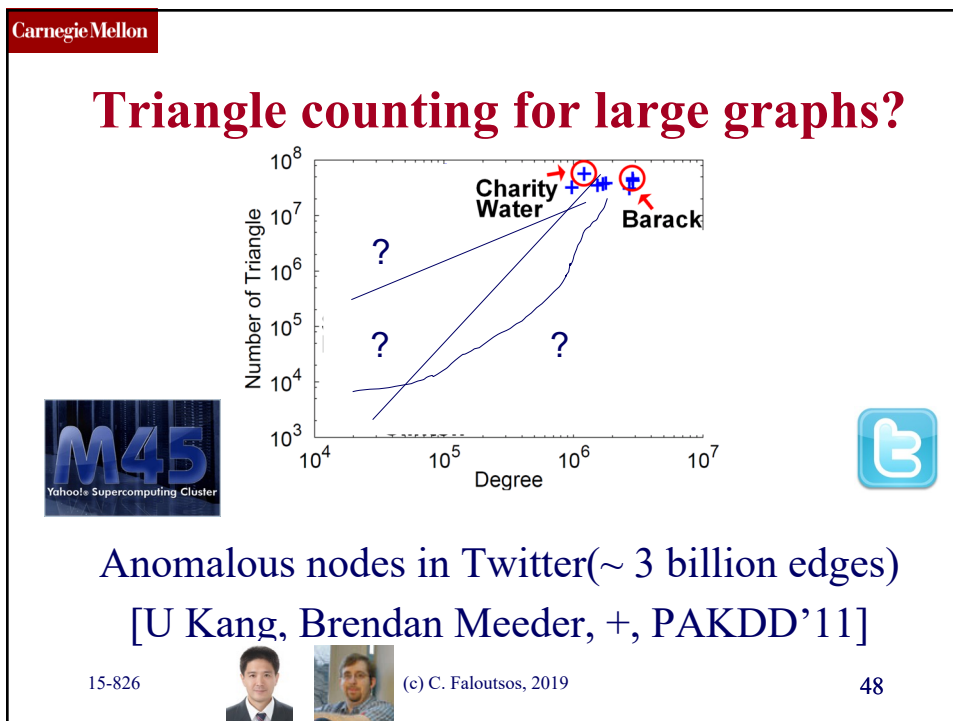
**#triangles =  $1/6 \text{ Sum } (\lambda_i^3)$**   
(and, because of skewness (S2),  
we only need the top few eigenvalues!)

15-826 (c) C. Faloutsos, 2019 46

46



47



48



CarnegieMellon

## Triangle counting for large graphs?

Number of Triangle

Degree

Sarah Palin

Hillary Clinton

John McCain

Barack Obama

Charity Water

Omitted

Twitter +

M45  
Yahoo! Supercomputing Cluster

Twitter

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

15-826 (c) C. Faloutsos, 2019 49

49

CarnegieMellon

## Triangle counting for large graphs?

Number of Triangle

Degree

Sarah Palin

Hillary Clinton

John McCain

Barack Obama

Charity Water

Omitted

Twitter +

M45  
Yahoo! Supercomputing Cluster

Twitter

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

15-826 (c) C. Faloutsos, 2019 50

50

CarnegieMellon

## Triangle counting for large graphs?

Number of Triangle

Degree

Twitter +

Adult Advertiser

Charity Water

Barack Obama

John McCain

Omitted

Sarah Palin

Hillary Clinton

15-826 (c) C. Faloutsos, 2019 51

51

CarnegieMellon

## Triangle counting for large graphs?

Number of Triangle

Degree

Twitter +

Adult Advertiser

Charity Water

Barack Obama

John McCain

Omitted

Sarah Palin


Hillary Clinton

15-826 (c) C. Faloutsos, 2019 52

52

CarnegieMellon

## List of Static Patterns



- ✓ • S.1 degree
- ✓ • S.2 eigenvalues
- ✓ • S.3 small diameter
- ✓ • S.4/5 Triangle laws
- (S.6) NLCC non-largest conn. components
- (S.7) eigen plots
- (S.8) radius plot

In textbook

15-826 (c) C. Faloutsos, 2019 53

53

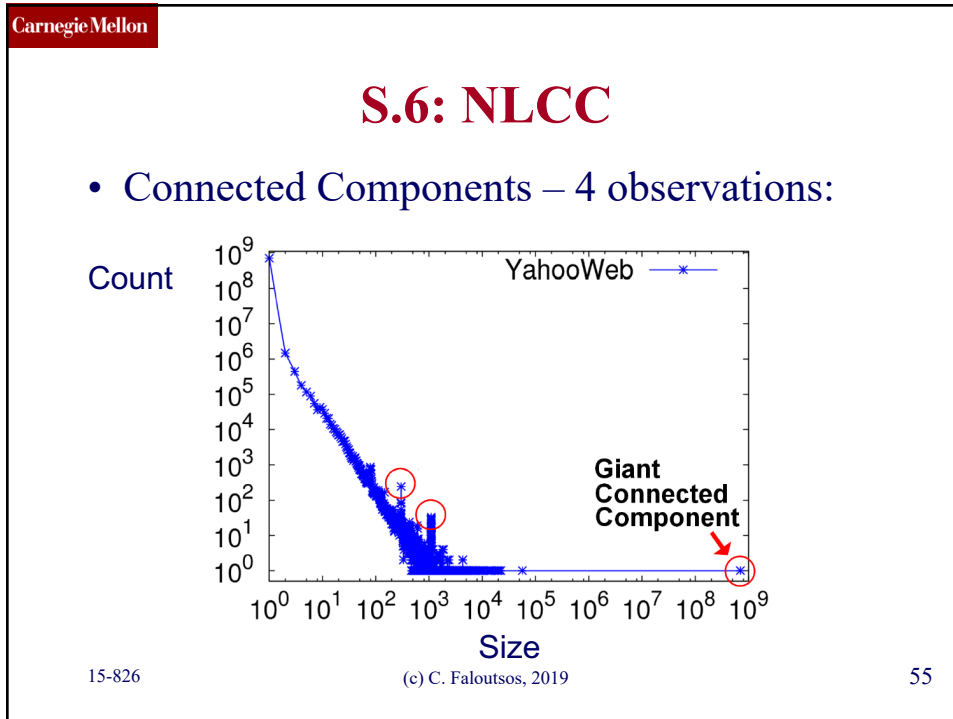
CarnegieMellon

## Generalized Iterated Matrix Vector Multiplication (GIMV)

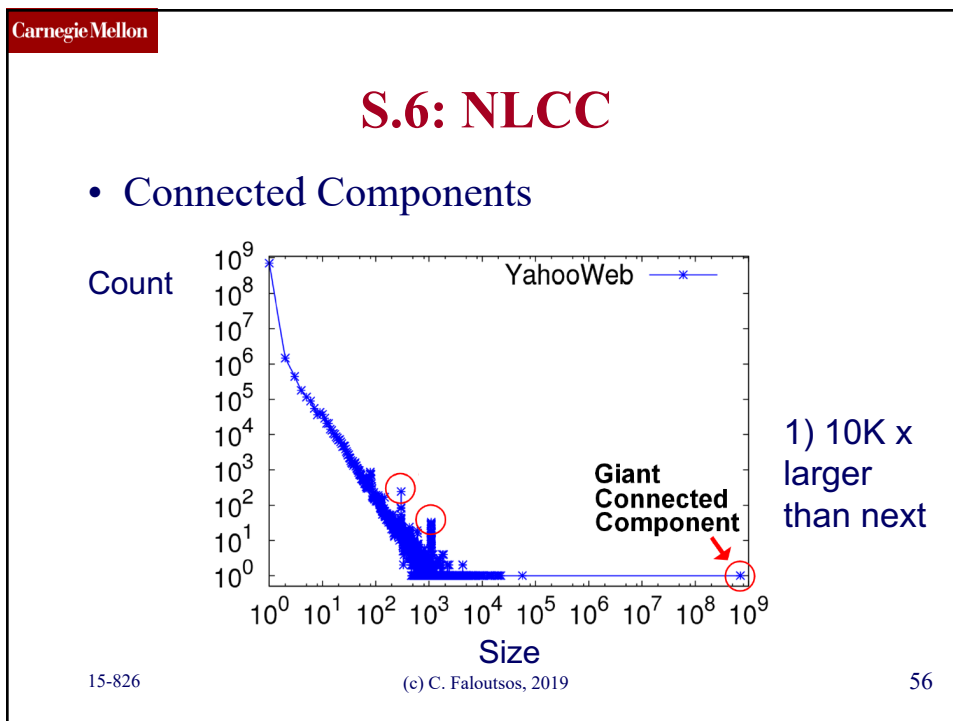
[PEGASUS: A Peta-Scale Graph Mining System - Implementation and Observations.](#)  
 U Kang, Charalampos E. Tsourakakis,  
 and Christos Faloutsos.  
 (ICDM) 2009, Miami, Florida, USA.  
 Best Application Paper (runner-up).

15-826 (c) C. Faloutsos, 2019 54

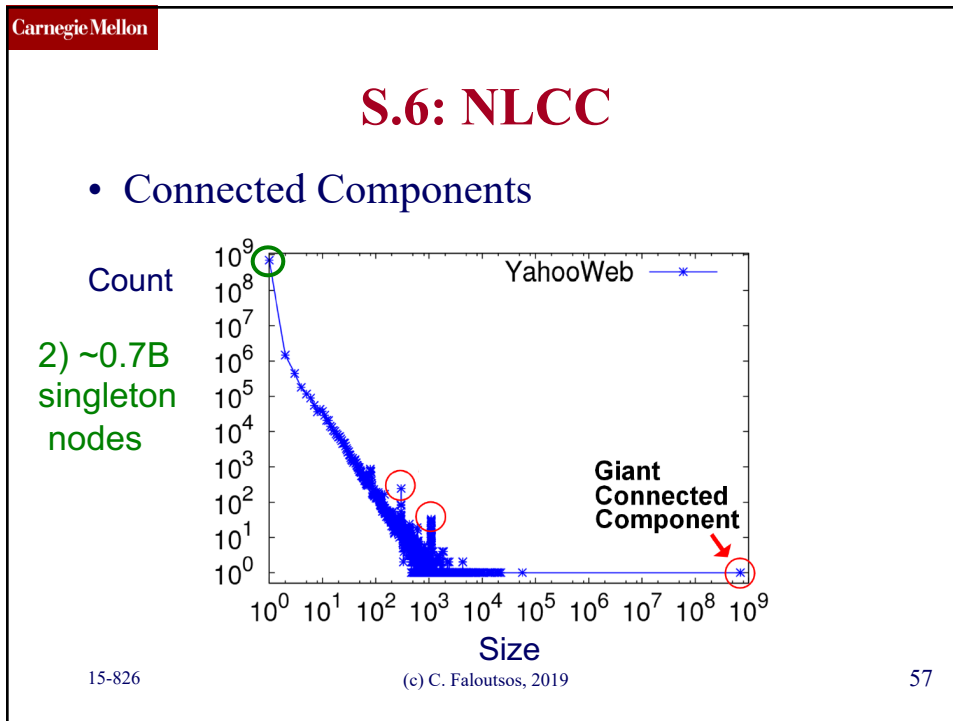
54



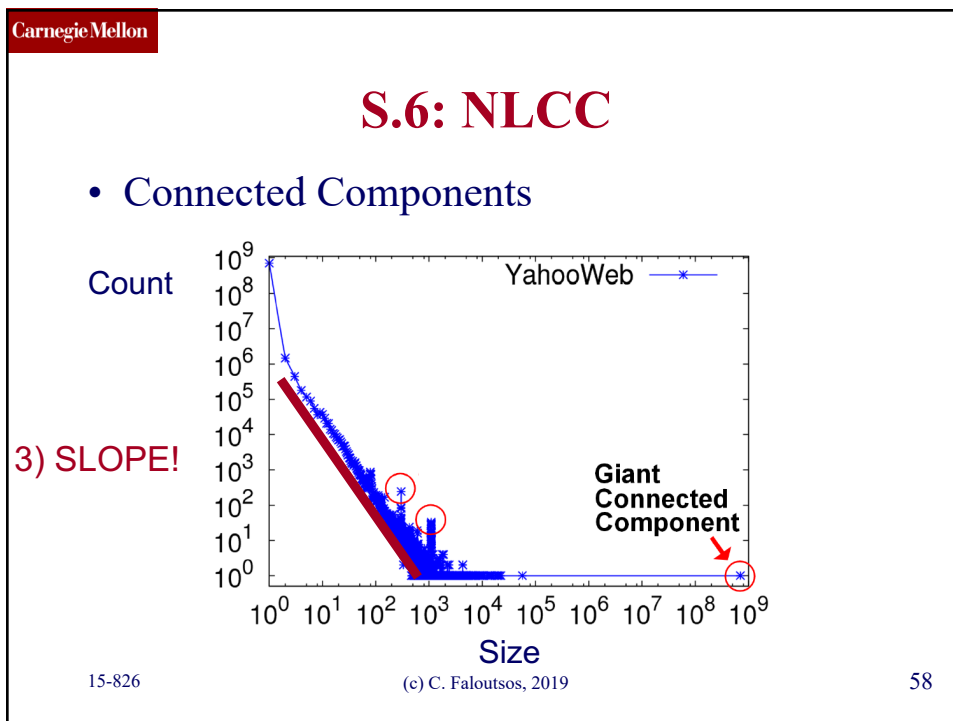
55



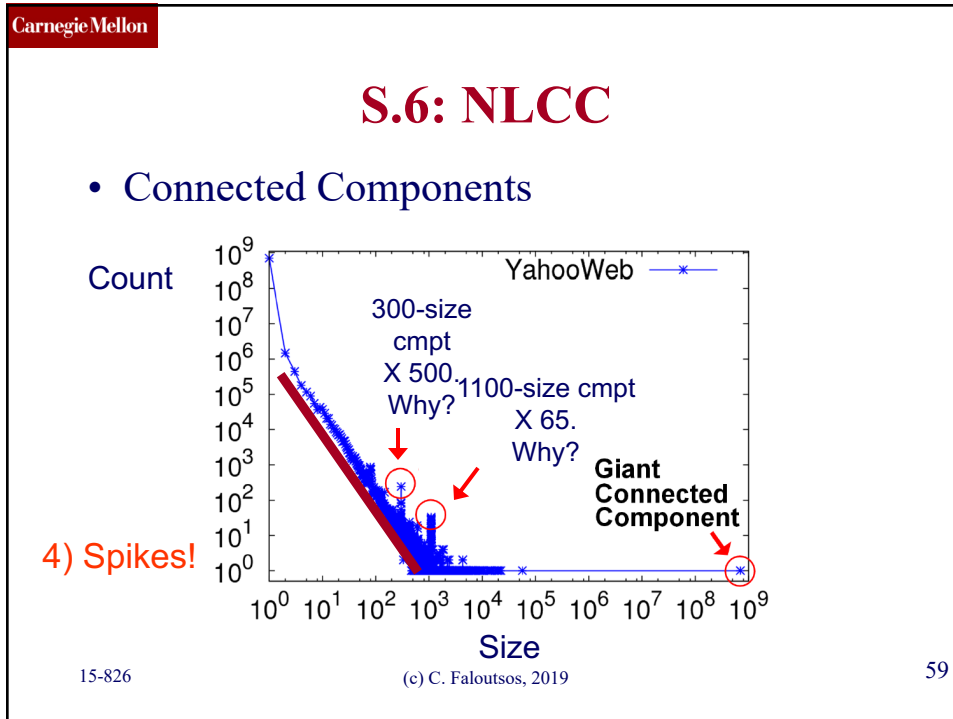
56



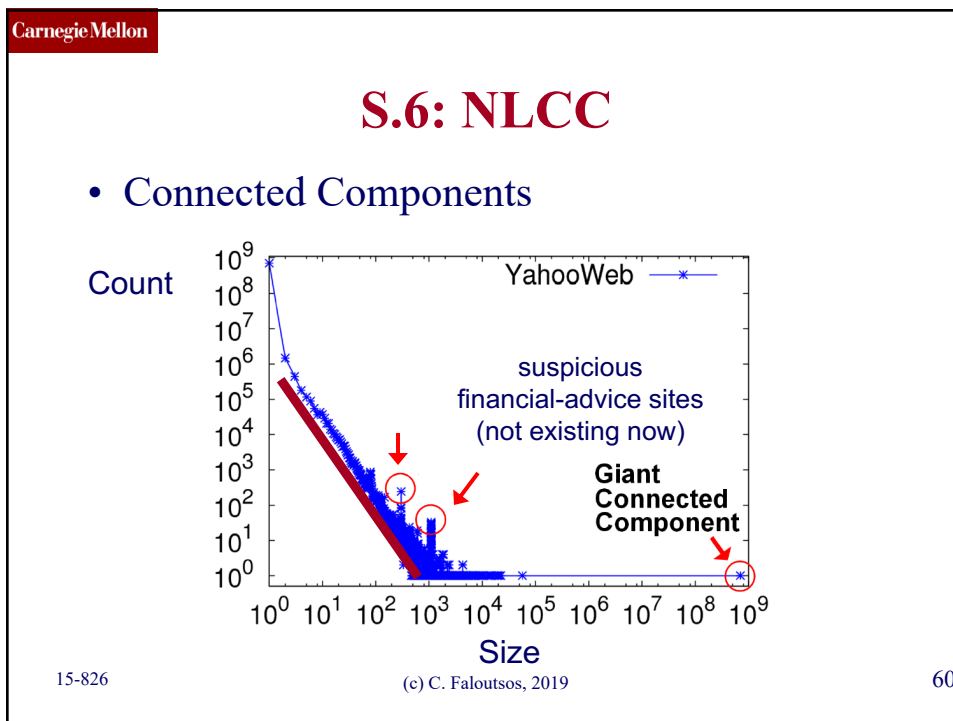
57



58



59



60

CarnegieMellon

## S.6: persists over time

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

2003  
Unstable Slope  
Giant Connected Component

2004  
Slope = -2.75  
Giant Connected Component

Stable tail slope after the gelling point

2005  
Slope = -2.75  
Giant Connected Component

2006  
Slope = -2.75  
Giant Connected Component

15-826 (c) C. Faloutsos, 2019 61

61

CarnegieMellon

## List of Static Patterns

- ✓ • S.1 degree
- ✓ • S.2 eigenvalues
- ✓ • S.3 small diameter
- ✓ • S.4/5 Triangle laws
- ✓ • (S.6) NLCC non-largest conn. components
- (S.7) eigen plots
- (S.8) radius plot

} In textbook

15-826 (c) C. Faloutsos, 2019 62

62

CarnegieMellon

## 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!

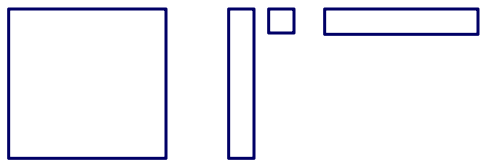
15-826 (c) C. Faloutsos, 2019 63

63

CarnegieMellon

## EigenSpokes

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

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


15-826 (c) C. Faloutsos, 2019 64

64



CarnegieMellon

## EigenSpokes

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

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

15-826 65  
(c) C. Faloutsos, 2019

65

CarnegieMellon

## EigenSpokes

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

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

15-826 66  
(c) C. Faloutsos, 2019

66

CarnegieMellon

## EigenSpokes

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

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

$\vec{u}_1$     $\vec{u}_i$

15-826 67

(c) C. Faloutsos, 2019

67

CarnegieMellon

## EigenSpokes

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

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

$\vec{u}_1$     $\vec{u}_i$

15-826 68

(c) C. Faloutsos, 2019

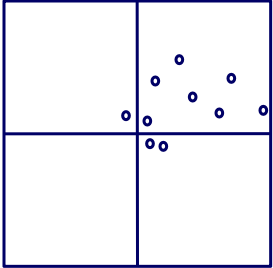
68

CarnegieMellon

## EigenSpokes

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

2<sup>nd</sup> Principal component  
 $u_2$



$u_1$

1<sup>st</sup> Principal component

15-826 (c) C. Faloutsos, 2019 69

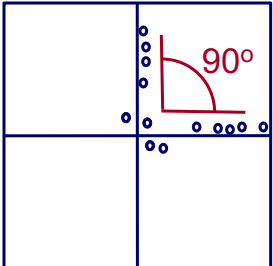
69

CarnegieMellon

## EigenSpokes

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

$u_2$



$u_1$

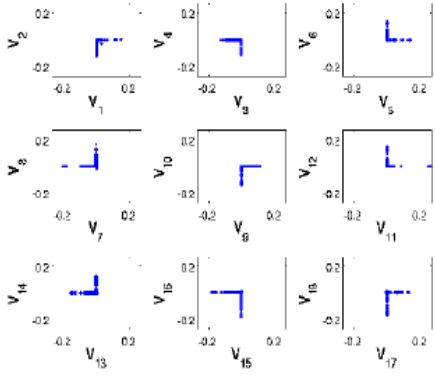
15-826 (c) C. Faloutsos, 2019 70

70

CarnegieMellon

## EigenSpokes - pervasiveness

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



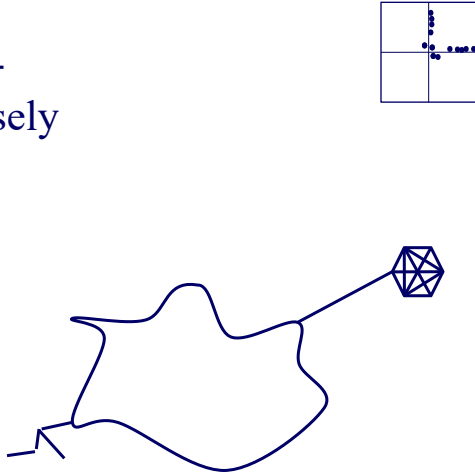
15-826 (c) C. Faloutsos, 2019 71

71

CarnegieMellon

## EigenSpokes - explanation

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



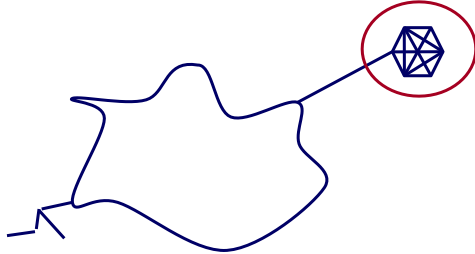
15-826 (c) C. Faloutsos, 2019 72

72

CarnegieMellon

## EigenSpokes - explanation

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



The diagram shows a large, irregularly shaped network with a jagged boundary. A single, dense, triangular-like clique is highlighted with a red circle. A line connects this clique to the rest of the network. In the top right corner, there is a small 2x2 grid with a red circle around the top-right cell.

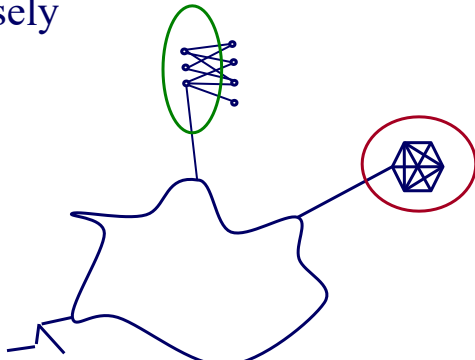
15-826 (c) C. Faloutsos, 2019 73

73

CarnegieMellon

## EigenSpokes - explanation

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



The diagram shows a large, irregularly shaped network with a jagged boundary. Two distinct cliques are highlighted: one with a green circle and one with a red circle. A line connects these two cliques to the rest of the network. In the top right corner, there is a small 2x2 grid with a green circle around the top-left cell and a red circle around the top-right cell.

15-826 (c) C. Faloutsos, 2019 74

74

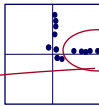
CarnegieMellon

## EigenSpokes - explanation

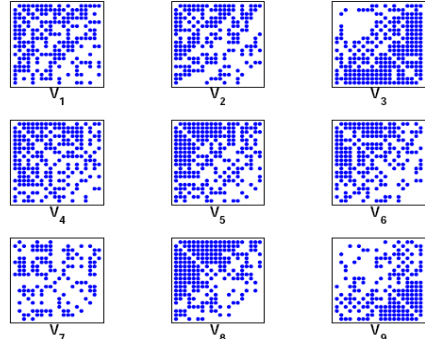
Near-cliques, or near-bipartite-cores, loosely connected

So what?

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



spy plot of top 20 nodes



15-826 (c) C. Faloutsos, 2019 75

75

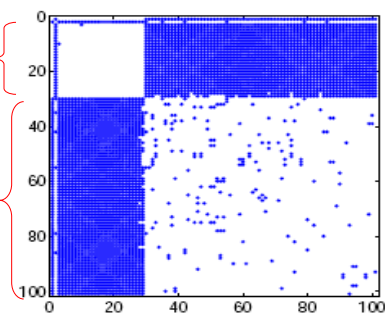
CarnegieMellon

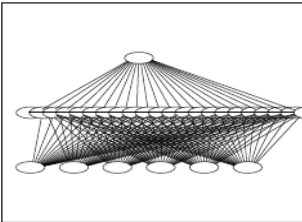
## Bipartite Communities!

patents from same inventor(s)

‘cut-and-paste’ bibliography!

magnified bipartite community





Useful for fraud detection!

15-826 (c) C. Faloutsos, 2019 76

76

CarnegieMellon

## Bipartite Communities!

IP – port scanners

victims


Useful for fraud detection!

15-826 (c) C. Faloutsos, 2019 77

77

CarnegieMellon

## List of Static Patterns




- ✓ • S.1 degree
- ✓ • S.2 eigenvalues
- ✓ • S.3 small diameter
- ✓ • S.4/5 Triangle laws
- ✓ • (S.6) NLCC non-largest conn. components
- ✓ • (S.7) eigen plots
- (S.8) radius plot

In textbook

15-826 (c) C. Faloutsos, 2019 78

78

CarnegieMellon

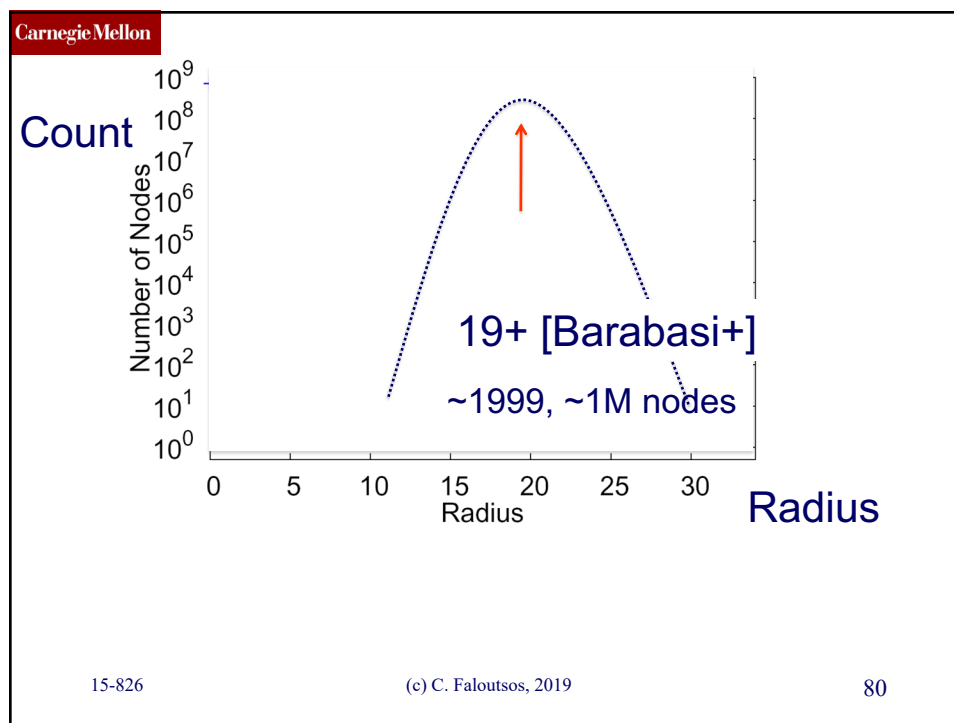


## HADI for diameter estimation

- *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

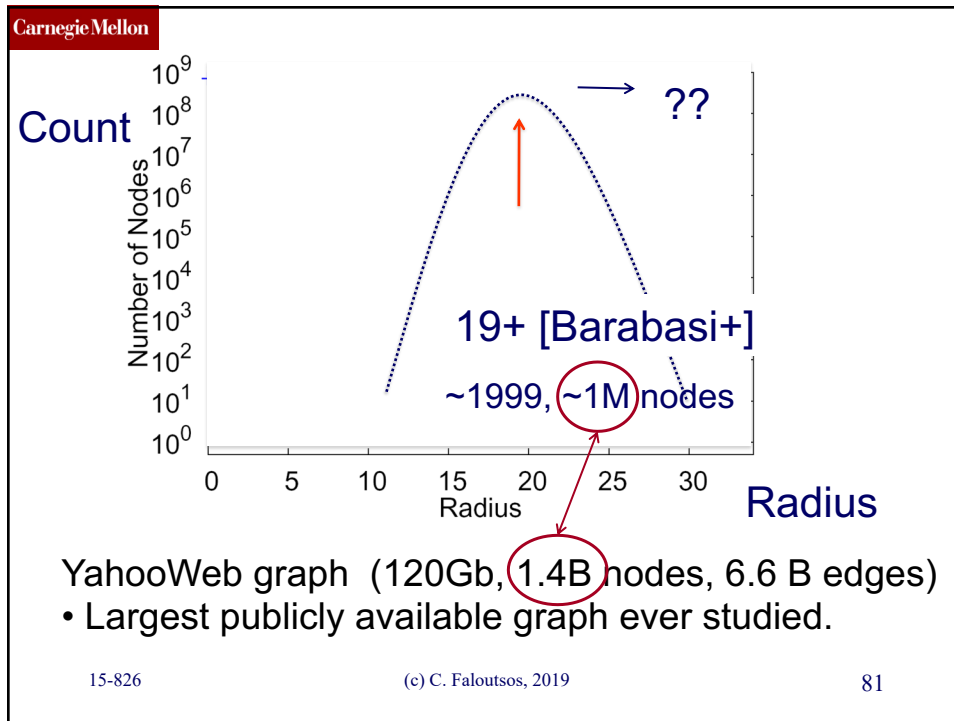
15-826 (c) C. Faloutsos, 2019 79

79

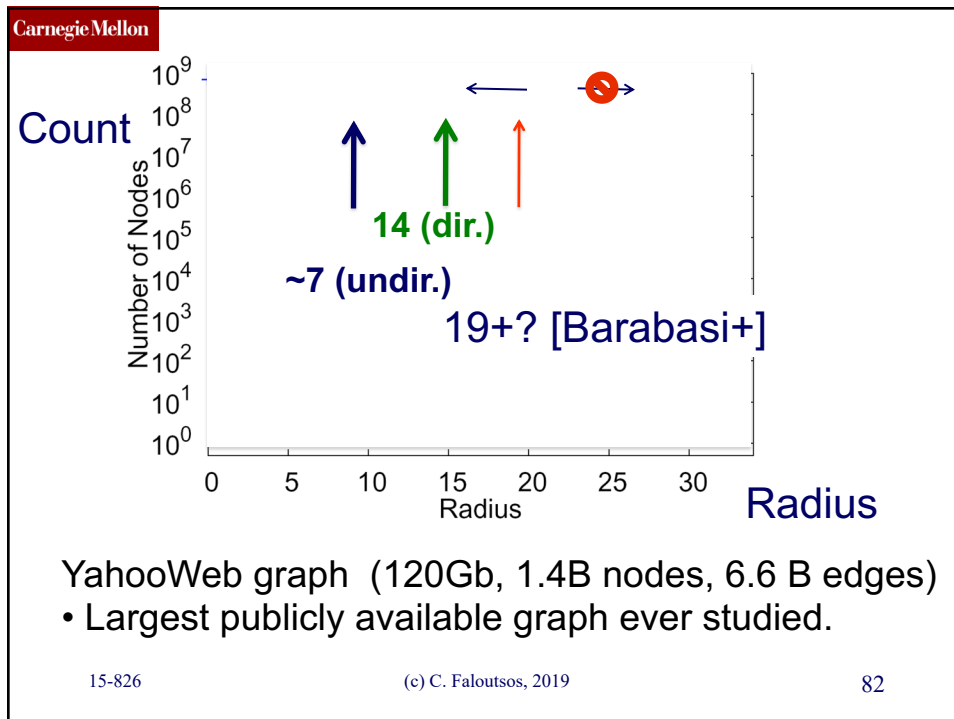


80

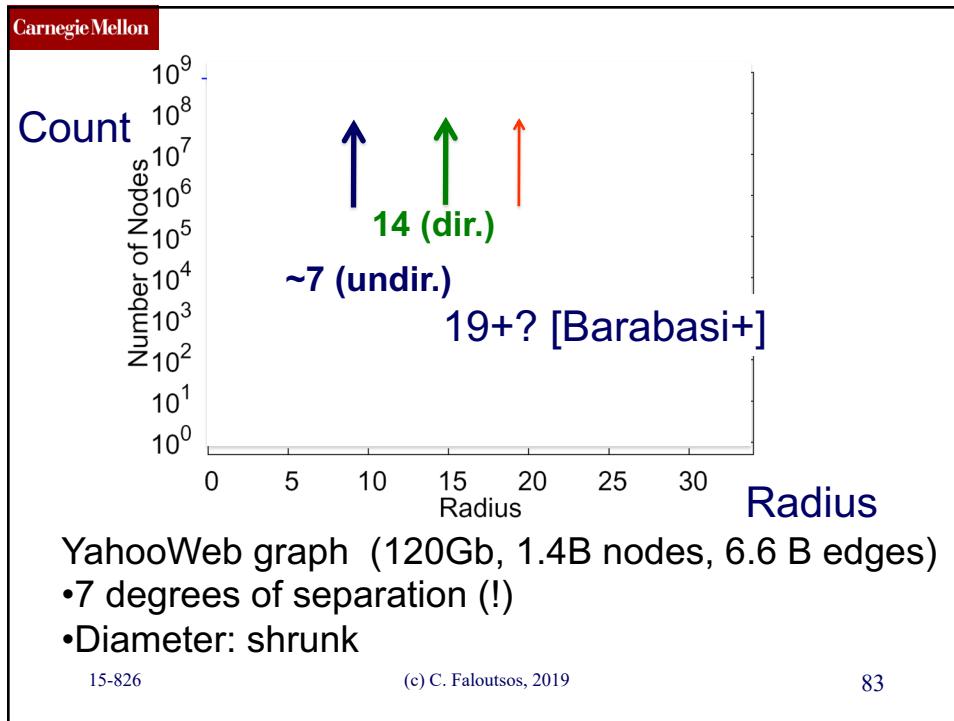




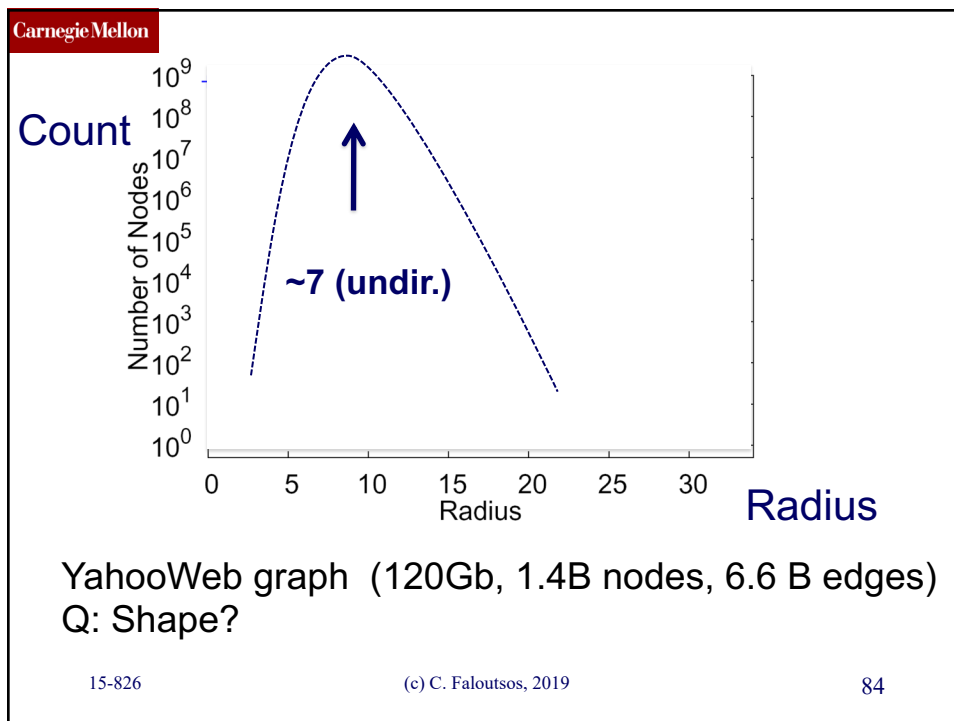
81



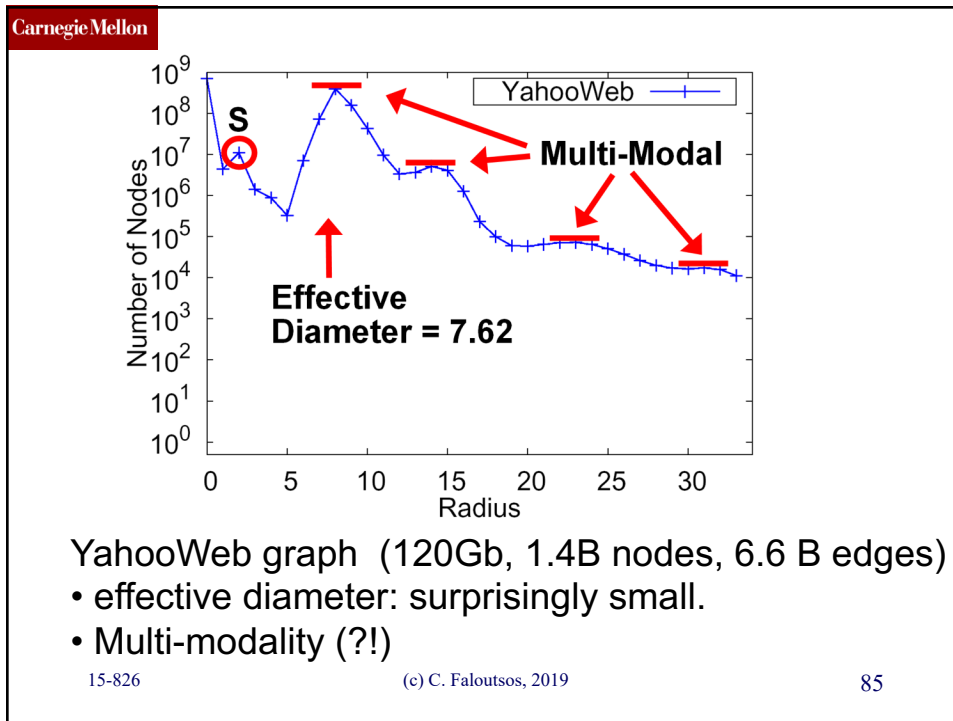
82



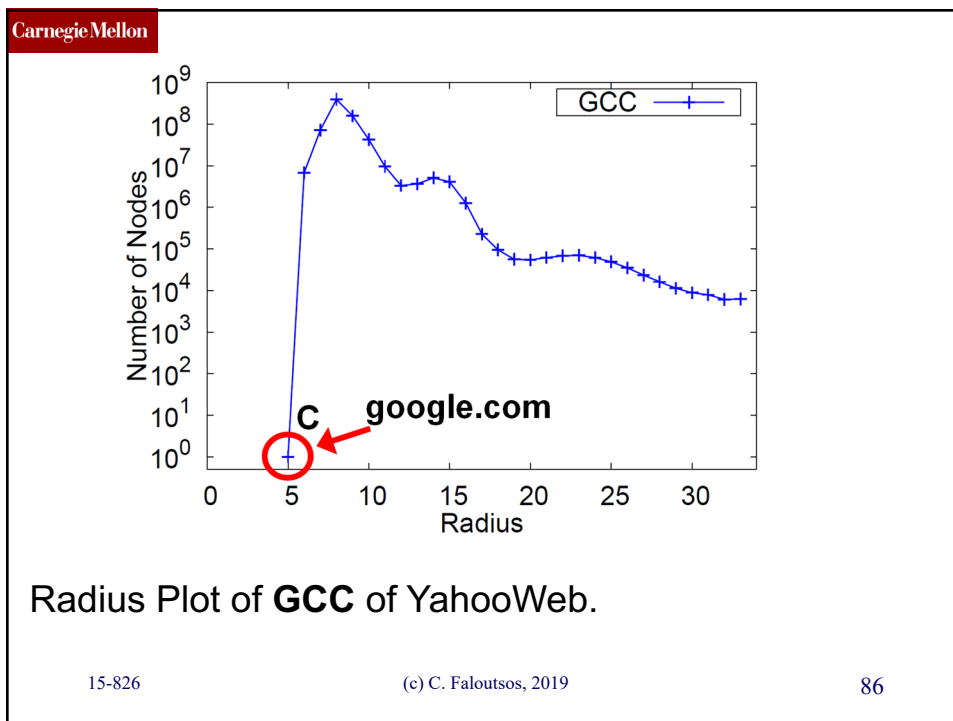
83



84



85



86

**CarnegieMellon**

Number of Nodes

Radius

YahooWeb

Effective Diameter = 7.62

Multi-Modal

S

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, 2019 87

87

**CarnegieMellon**

Number of Nodes

Radius

YahooWeb

Effective Diameter = 7.62

Multi-Modal

S

Conjecture:

EN

DE

BR

~7

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, 2019 88

88

**CarnegieMellon**

Conjecture:

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, 2019 89

89

**CarnegieMellon**

## List of Static Patterns

- ✓ • S.1 degree
- ✓ • S.2 eigenvalues
- ✓ • S.3 small diameter
- ✓ • S.4/5 Triangle laws
- ✓ • (S.6) NLCC non-largest conn. components
- ✓ • (S.7) eigen plots
- ✓ • (S.8) radius plot
- Other observations / patterns?

In textbook

15-826 (c) C. Faloutsos, 2019 90

90

CarnegieMellon

Attention:  
Phase I

## List of Static Patterns

- ✓ • S.1 degree
- ✓ • S.2 eigenvalues
- ✓ • S.3 small diameter
- ✓ • S.4/5 Triangle laws
- ✓ • (S.6) NLCC non-largest conn. components
- ✓ • (S.7) eigen plots
- ✓ • (S.8) radius plot
- Other observations / patterns?

In textbook

15-826 (c) C. Faloutsos, 2019 91

91

CarnegieMellon

## Any other 'laws' ?

Yes!

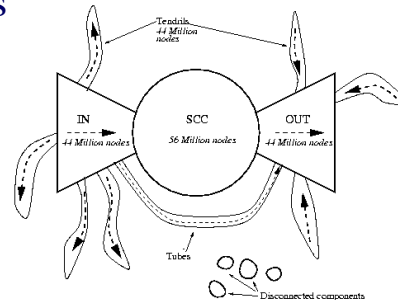
- Small diameter (~ constant!) –
  - six degrees of separation / 'Kevin Bacon'
  - small worlds [Watts and Strogatz]

15-826 (c) C. Faloutsos, 2019 92

92

## Any other 'laws' ?

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



15-826

(c) C. Faloutsos, 2019

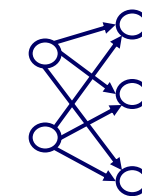
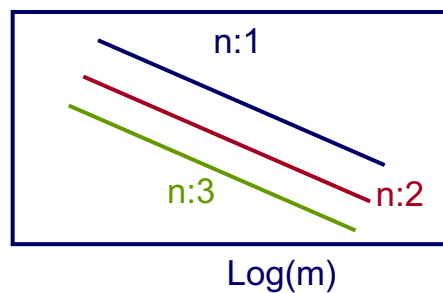
93

93

## 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, 2019

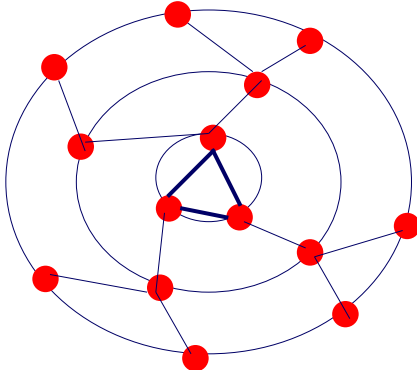
94

94

CarnegieMellon

## Any other 'laws' ?

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




15-826 (c) C. Faloutsos, 2019 95

95

CarnegieMellon

## Outline



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

15-826 (c) C. Faloutsos, 2019 96


96



CarnegieMellon

## 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, 2019 97

97

CarnegieMellon

## Observation W.1: Fortification

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

15-826 (c) C. Faloutsos, 2019 98

98

CarnegieMellon

## Observation W.1: Fortification

**More donors,  
more \$ ?**

15-826 (c) C. Faloutsos, 2019 99

99

CarnegieMellon

## 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  
(\$)

**Orgs-Candidates**


Edges (# donors)

(c) C. Faloutsos, 2019 100

100

CarnegieMellon

## Outline



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



15-826 (c) C. Faloutsos, 2019 101

101

CarnegieMellon

## Problem: Time evolution


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



15-826 (c) C. Faloutsos, 2019 102

102

CarnegieMellon



## List of Dynamic Patterns

- D.1 diameter
- D.2 densification
- D.3 gelling point
- D.4 NLCC over time
- D.5 Eigenvalue over time
- D.6 Popularity over time
- D.7 phonecall duration

In textbook

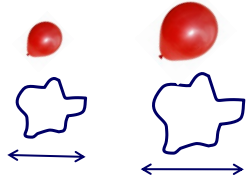
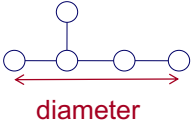
15-826 (c) C. Faloutsos, 2019 103

103

CarnegieMellon

## D.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?

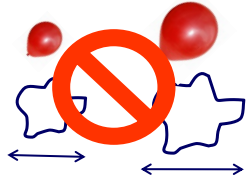
15-826 (c) C. Faloutsos, 2019 104

104

CarnegieMellon

## D.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



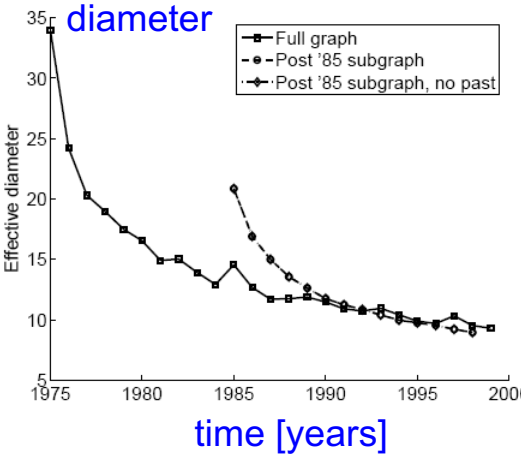
15-826 (c) C. Faloutsos, 2019 105

105

CarnegieMellon

## D.1 Diameter – “Patents”

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




Effective diameter

time [years]

15-826 (c) C. Faloutsos, 2019 106

106

CarnegieMellon



## List of Dynamic Patterns

- ✓ • D.1 diameter
- D.2 densification
- D.3 gelling point
- D.4 NLCC over time
- D.5 Eigenvalue over time
- D.6 Popularity over time
- D.7 phonecall duration

In textbook

15-826 (c) C. Faloutsos, 2019 107

107

CarnegieMellon

## D.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, 2019 108

108

CarnegieMellon

## D.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, 2019 109

109

CarnegieMellon


## D.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, 2019 110

110

CarnegieMellon



## List of Dynamic Patterns

- ✓ • D.1 diameter
- ✓ • D.2 densification
- D.3 gelling point
- D.4 NLCC over time
- D.5 Eigenvalue over time
- D.6 Popularity over time
- D.7 phonecall duration

In textbook

15-826 (c) C. Faloutsos, 2019 111

111

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) C. Faloutsos, 2019 112

112

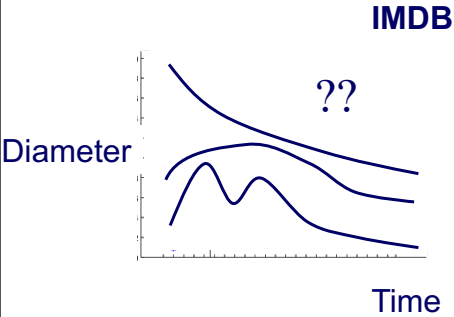


CarnegieMellon

## D.3 Gelling Point

- Diameter, over time

IMDB



Time

15-826 (c) C. Faloutsos, 2019 113

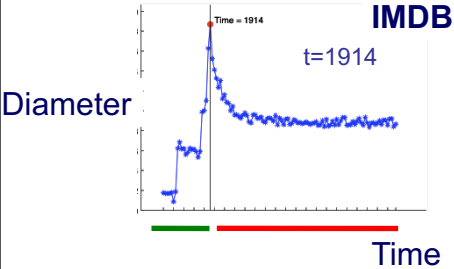
113

CarnegieMellon

## D.3 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

15-826 (c) C. Faloutsos, 2019 114

114

CarnegieMellon

## D.3 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, 2019

115

115

CarnegieMellon

## List of Dynamic Patterns

- ✓ • D.1 diameter
- ✓ • D.2 densification
- ✓ • D.3 gelling point
- D.4 NLCC over time
- D.5 Eigenvalue over time
- D.6 Popularity over time
- D.7 phonecall duration

In textbook

15-826

(c) C. Faloutsos, 2019

116

116

## Observation D.4: 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, 2019

117

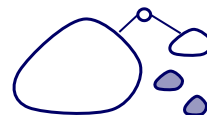
117

## Observation D.4: 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, 2019

118

118

## Observation D.4: 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

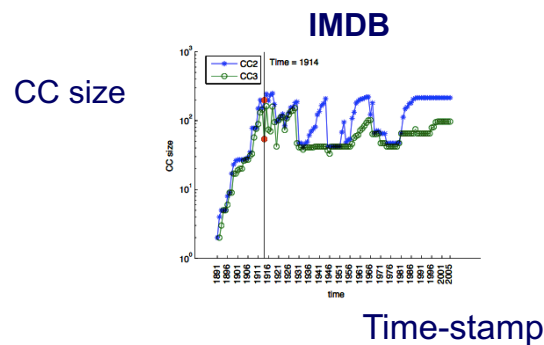
(c) C. Faloutsos, 2019

119

119

## Observation D.4: NLCC behavior

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




15-826

(c) C. Faloutsos, 2019

120

120

CarnegieMellon



## List of Dynamic Patterns

- ✓ • D.1 diameter
- ✓ • D.2 densification
- ✓ • D.3 gelling point
- ✓ • D.4 NLCC over time
- ~~D.5 Eigenvalue over time~~
- D.6 Popularity over time
- D.7 phonecall duration

In textbook

15-826 (c) C. Faloutsos, 2019 121

121

CarnegieMellon

## Timing for Blogs

*Cascading Behavior in Large Blog Graphs:  
Patterns and a model*

Jure Leskovec, Mary McGlohon, Christos Faloutsos, Natalie Glance, Matthew Hurst  
SDM'07

15-826 (c) C. Faloutsos, 2019 122

122

CarnegieMellon

## D.6 : popularity over time

# in links

lag: days after post

Post popularity drops-off – exponentially?

@t + lag      @t

15-826 (c) C. Faloutsos, 2019 123

123

CarnegieMellon

## D.6 : popularity over time

# in links  
(log)

days after post  
(log)

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

15-826 (c) C. Faloutsos, 2019 124

124

CarnegieMellon

## D.6 : 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, 2019 125

125

CarnegieMellon


## -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.** 126

126

CarnegieMellon



## List of Dynamic Patterns

- ✓ • D.1 diameter
- ✓ • D.2 densification
- ✓ • D.3 gelling point
- ✓ • D.4 NLCC over time
- ~~D.5 Eigenvalue over time~~
- ✓ • D.6 Popularity over time
- D.7 phonecall duration

In textbook


15-826 (c) C. Faloutsos, 2019 127

127

CarnegieMellon

## D.7: duration of phonecalls

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



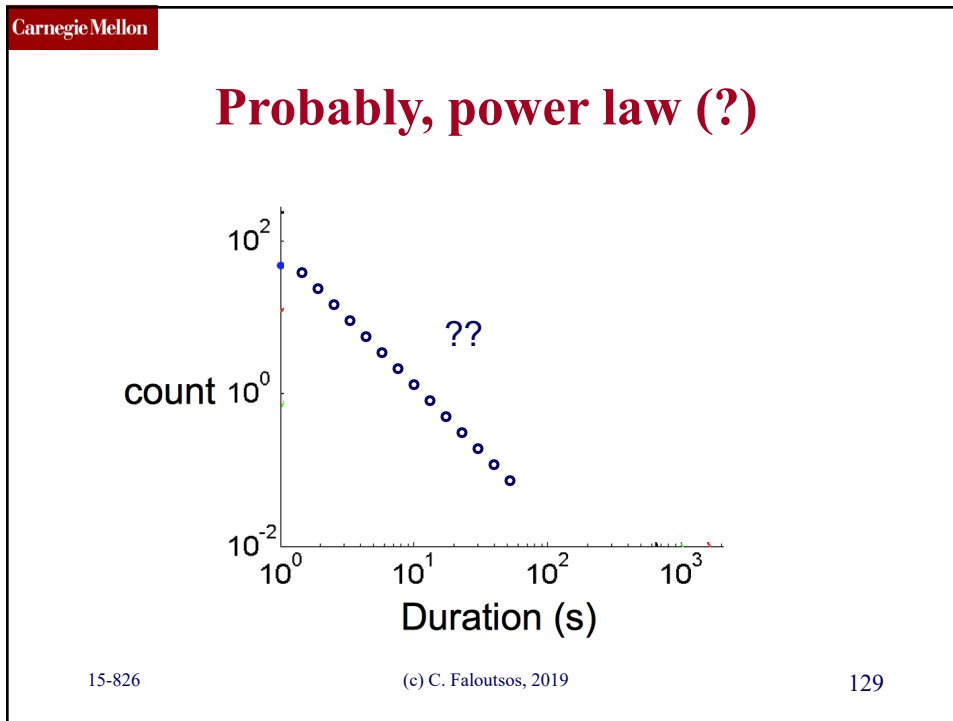
Pedro O. S. Vaz de Melo, Leman Akoglu, Christos Faloutsos, Antonio A. F. Loureiro

PKDD 2010

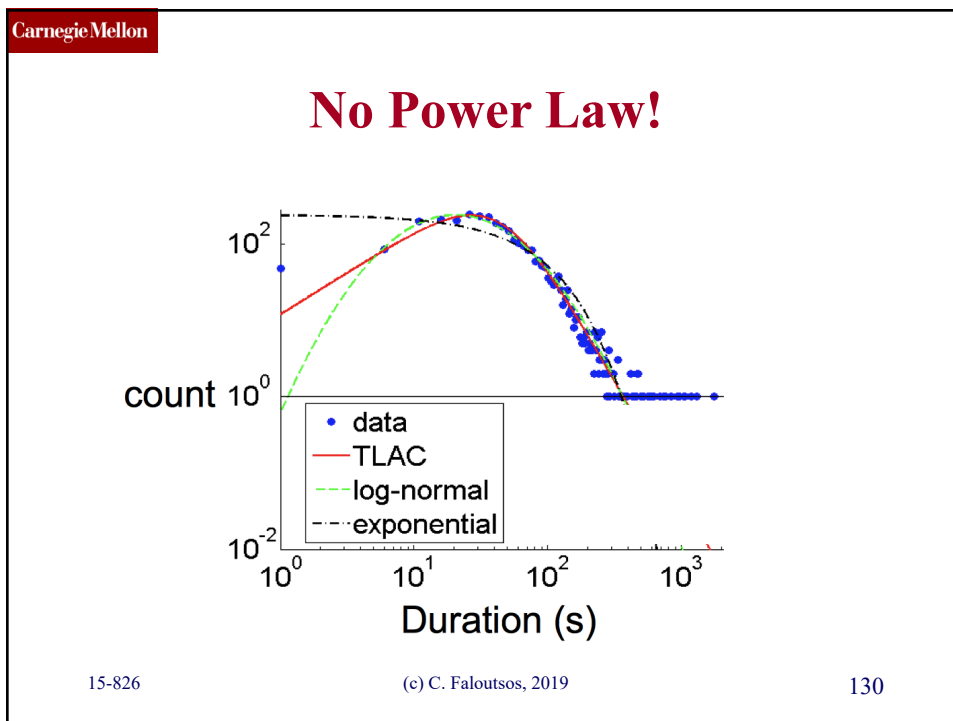
15-826 (c) C. Faloutsos, 2019 128

128





129



130

**CarnegieMellon**

## ‘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, 2019 131

131

**CarnegieMellon**

## 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, 2019 132

132

CarnegieMellon

## Log-logistic distribution

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

count

Duration (s)

OR(t)

duration (s)

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

15-826 (c) C. Faloutsos, 2019 133

133

CarnegieMellon

## Log-logistic distribution

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

OR(t)

Duration ( t )

15-826 134

134

CarnegieMellon

## Log-logistic distribution

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

$x \rightarrow \ln(x)$

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

15-826

CDF(x) =  $1/(1+1/x)$

(c) C. Faloutsos, 2019

135

135

CarnegieMellon

## Log-logistic distribution

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

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

15-826

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

(c) C. Faloutsos, 2019

136

136

CarnegieMellon

Attention:  
Phase 1

## Log-logistic distribution

Nice 1 page description: section II of

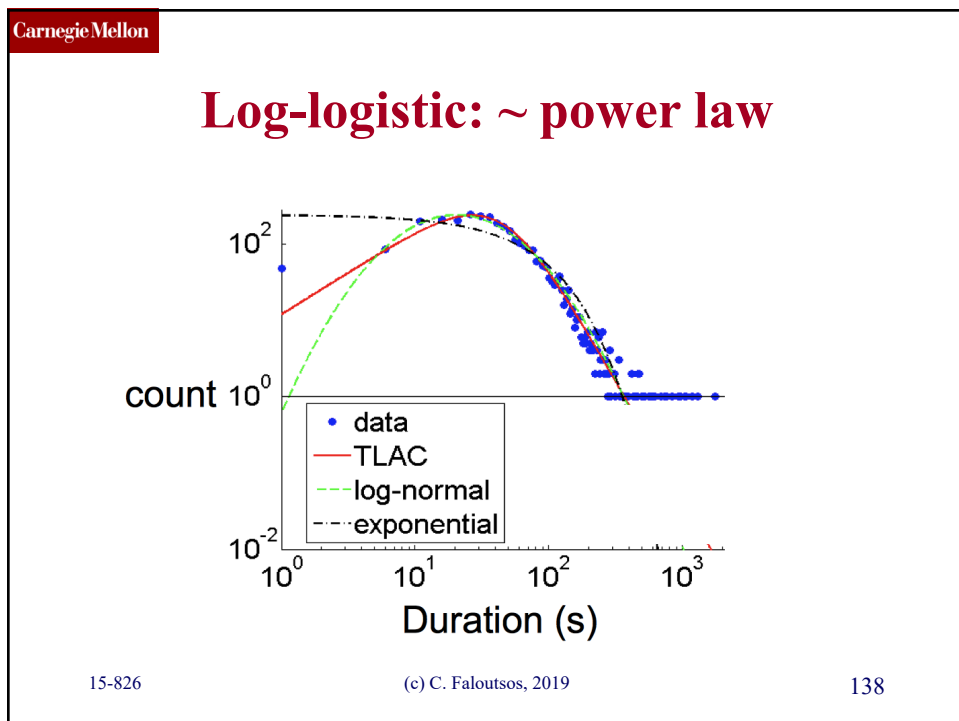
Pravallika Devineni, Danai Koutra, Michalis Faloutsos, and Christos Faloutsos.

[\*If walls could talk: Patterns and anomalies in Facebook wallposts.\*](#)

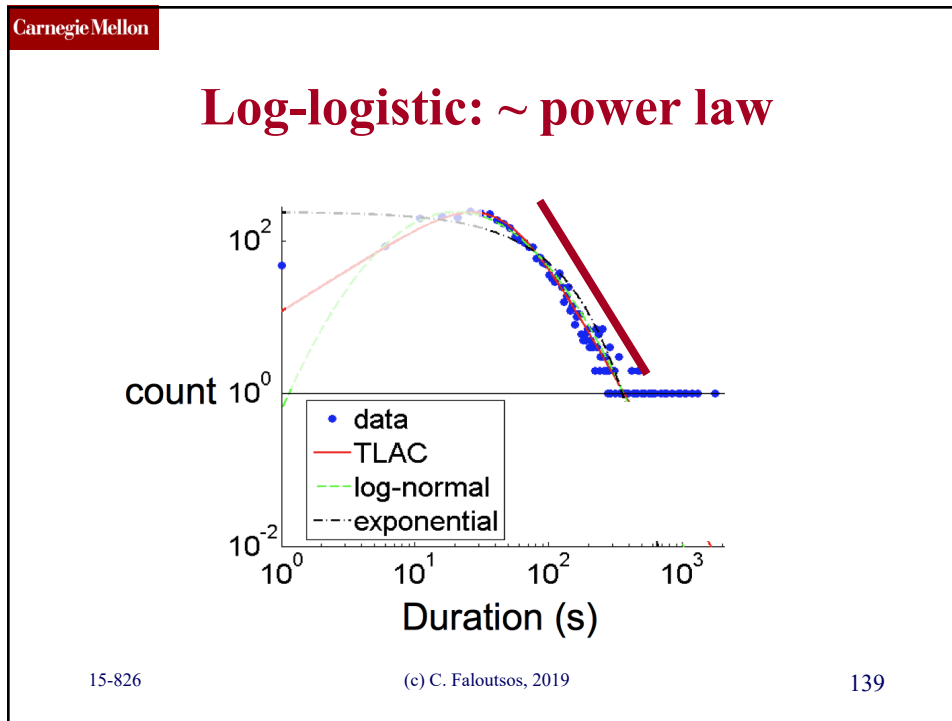
*ASONAM 2015, pp 367-374.*

15-826 (c) C. Faloutsos, 2019 137

137



138



139

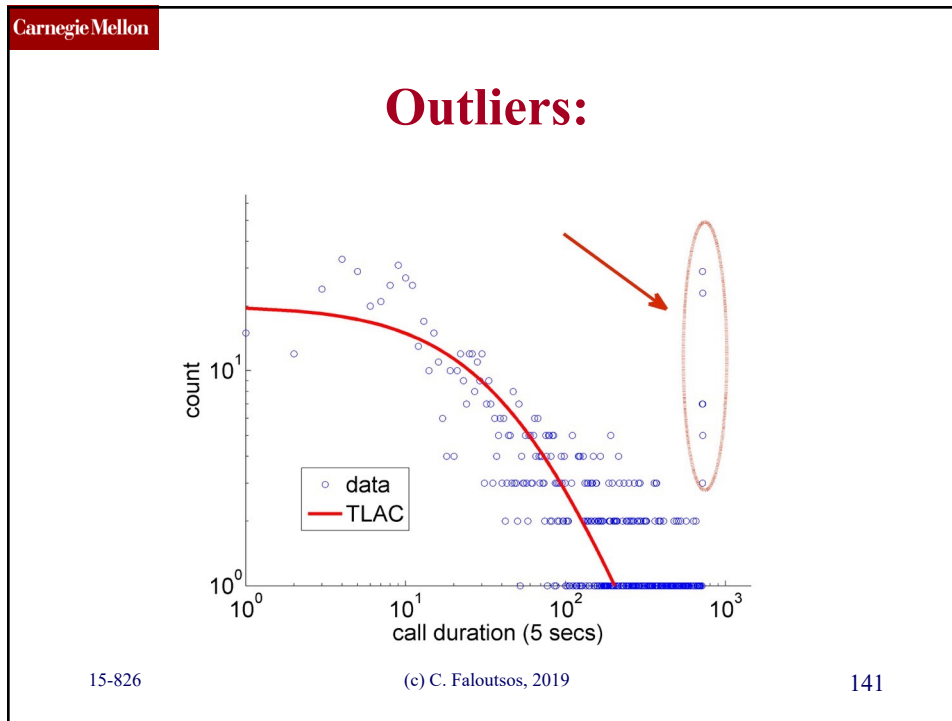
CarnegieMellon

## 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, 2019 140

140



141


CarnegieMellon 

## Conclusions

- Are real graphs random?
- NO!
  - Static patterns
    - Small diameters
    - Skewed degree distribution
    - Shrinking diameters
  - Weighted
  - Time-evolving

15-826 Copyright: C. Faloutsos (2019) 142

142

CarnegieMellon 

## Conclusions

- Are real graphs random?
- NO!
  - Static patterns
    - Small diameter
    - Skewed
  - Evolving
    - Many power laws – log-logistic
    - Take logarithms

15-826 Copyright: C. Faloutsos (2019) 143

143

CarnegieMellon

## Next lecture:

- Anomaly detection tools (OddBall, etc)

15-826 Copyright: C. Faloutsos (2019) 144

144