

CarnegieMellon

15-826: Multimedia Databases and Data Mining

Lecture #8: Fractals - introduction

C. Faloutsos

CarnegieMellon

Must-read Material

- Christos Faloutsos and Ibrahim Kamel, [*Beyond Uniformity and Independence: Analysis of R-trees Using the Concept of Fractal Dimension*](#), Proc. ACM SIGACT-SIGMOD-SIGART PODS, May 1994, pp. 4-13, Minneapolis, MN.


Recommended Material

optional, but **very** useful:

- Manfred Schroeder *Fractals, Chaos, Power Laws: Minutes from an Infinite Paradise*
W.H. Freeman and Company, 1991
 - Chapter 10: boxcounting method
 - Chapter 1: Sierpinski triangle

Outline

Goal: 'Find **similar / interesting** things'

- Intro to DB
-  • Indexing - similarity search
- Data Mining

Indexing - Detailed outline

- primary key indexing
- secondary key / multi-key indexing
- spatial access methods
 - z-ordering
 - R-trees
 - misc
- fractals
 - intro
 - applications
- text



15-826

Copyright: C. Faloutsos (2019)

5

Intro to fractals - outline


- ➔ • Motivation – 3 problems / case studies
- Definition of fractals and power laws
- Solutions to posed problems
- More examples and tools
- Discussion - putting fractals to work!
- Conclusions – practitioner's guide
- Appendix: gory details - boxcounting plots

15-826

Copyright: C. Faloutsos (2019)

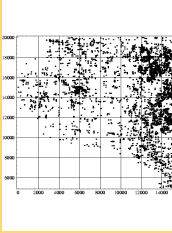
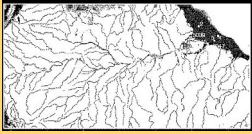
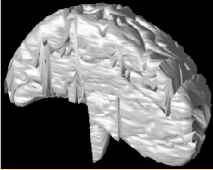
6

CarnegieMellon




Problem

- What patterns are in real k -dim points?

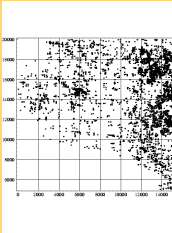

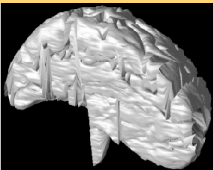
15-826 Copyright: C. Faloutsos (2019) 7

CarnegieMellon



Conclusions

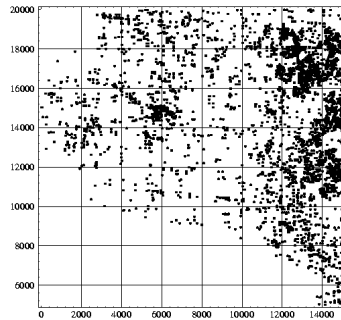
- What patterns are in real k -dim points?
- Self-similarity (= fractals \rightarrow power laws)

15-826 Copyright: C. Faloutsos (2019) 8

CarnegieMellon

Problem #1: GIS - points



Road end-points of
Montgomery county:

- Q1: how many d.a. for an R-tree?
- Q2 : distribution?
 - not uniform
 - not Gaussian
 - no rules??

15-826

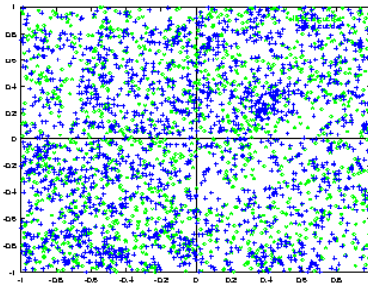
Copyright: C. Faloutsos (2019)

9

CarnegieMellon

Problem #2 - spatial d.m.

Galaxies (Sloan Digital Sky Survey w/ B. Nichol)



- 'spiral' and 'elliptical'
galaxies

(stores and households ...)

- patterns?
- attraction/repulsion?
- how many 'spi' within r
from an 'ell'?

15-826

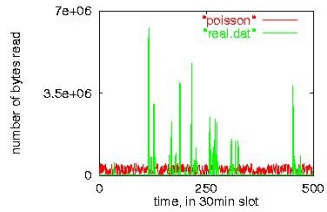
Copyright: C. Faloutsos (2019)

10

CarnegieMellon

Problem #3: traffic

bytes



number of bytes read

time, in 30min slot

time

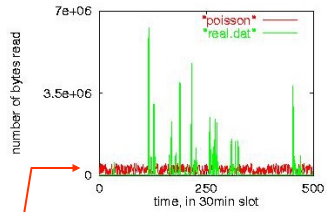
- disk trace (from HP - J. Wilkes); Web traffic - fit a model
- how many explosions to expect?
- queue length distr.?

15-826 Copyright: C. Faloutsos (2019) 11

CarnegieMellon

Problem #3: traffic

bytes



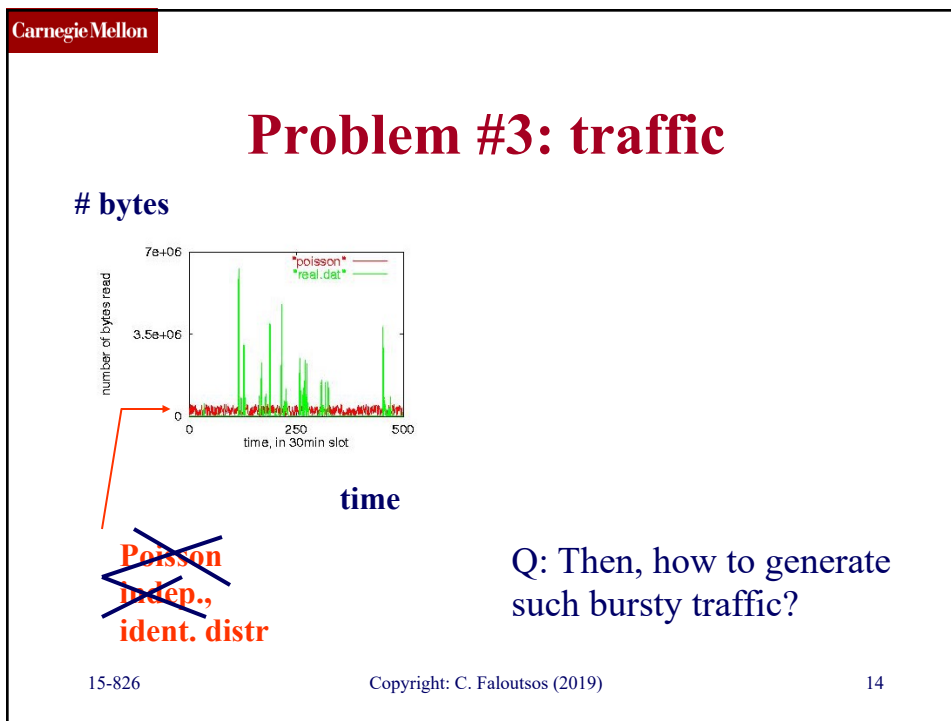
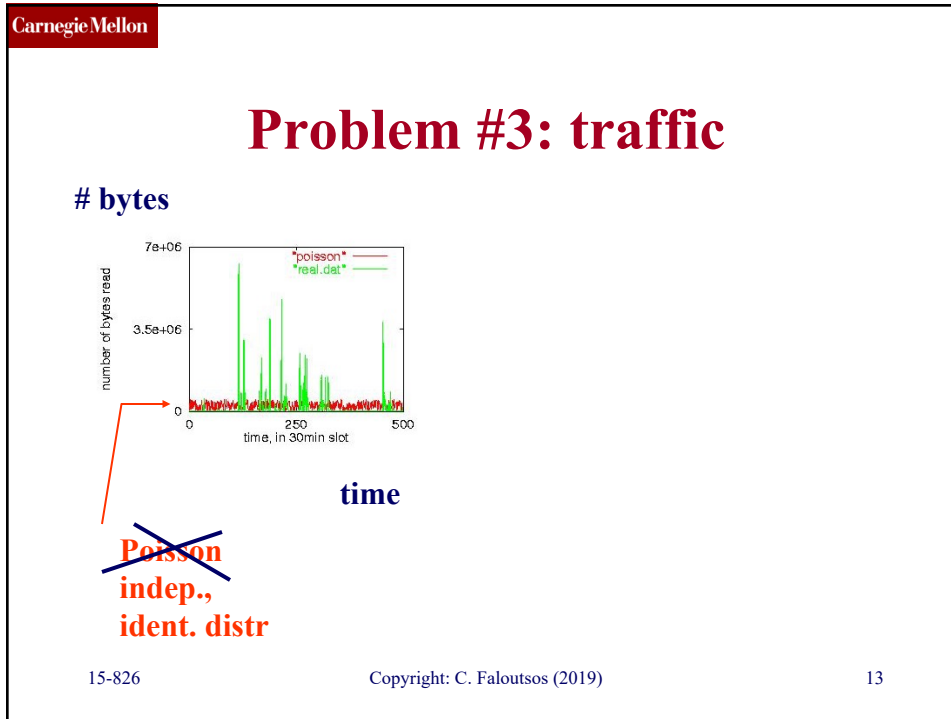
number of bytes read

time, in 30min slot

time

**Poisson
indep.,
ident. distr**


15-826 Copyright: C. Faloutsos (2019) 12



Common answer:

- Fractals / self-similarities / power laws
- Seminal works from Hilbert, Minkowski, Cantor, Mandelbrot, (Hausdorff, Lyapunov, Ken Wilson, ...)

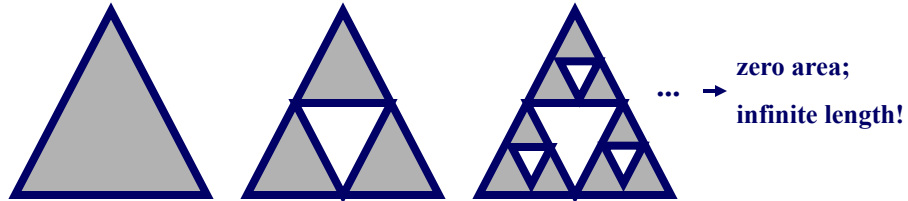
Road map

- Motivation – 3 problems / case studies
-  • Definition of fractals and power laws
- Solutions to posed problems
- More examples and tools
- Discussion - putting fractals to work!
- Conclusions – practitioner's guide
- Appendix: gory details - boxcounting plots

CarnegieMellon

What is a fractal?

= self-similar point set, e.g., Sierpinski triangle:



Dimensionality??

15-826

Copyright: C. Faloutsos (2019)

17

CarnegieMellon

Definitions (cont' d)

- Paradox: Infinite perimeter ; Zero area!
- 'dimensionality' : between 1 and 2
- actually: $\text{Log}(3)/\text{Log}(2) = 1.58\dots$

15-826

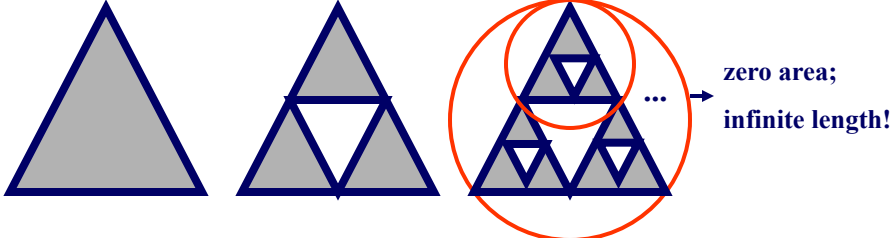
Copyright: C. Faloutsos (2019)

18

CarnegieMellon

Dfn of fd:

ONLY for a perfectly self-similar point set:



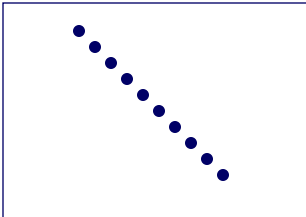
$= \log(n)/\log(f) = \log(3)/\log(2) = 1.58$

15-826 Copyright: C. Faloutsos (2019) 19

CarnegieMellon

Intrinsic ('fractal') dimension

- Q: fractal dimension of a line?
- A: 1 (= $\log(2)/\log(2)$!)

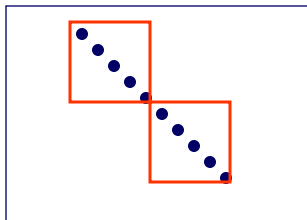


15-826 Copyright: C. Faloutsos (2019) 20

CarnegieMellon

Intrinsic ('fractal') dimension

- Q: fractal dimension of a line?
- A: 1 ($= \log(2)/\log(2)!$)



15-826

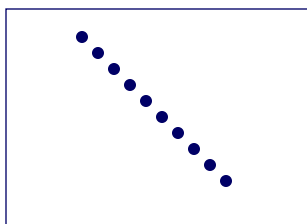
Copyright: C. Faloutsos (2019)

21

CarnegieMellon

Intrinsic ('fractal') dimension

- Q: dfn for a given set of points?



x	y
5	1
4	2
3	3
2	4

15-826

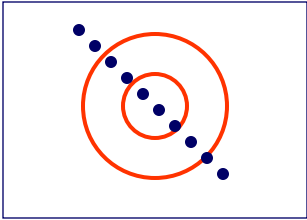
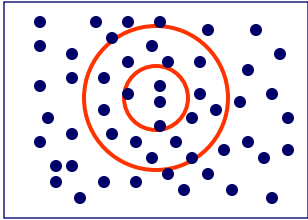
Copyright: C. Faloutsos (2019)

22

CarnegieMellon

Intrinsic ('fractal') dimension

- Q: fractal dimension of a line?
- A: $mn (<= r) \sim r^1$
('power law' : $y=x^a$)
- Q: fd of a plane?
- A: $mn (<= r) \sim r^2$
fd == slope of $(\log(mn) \text{ vs } \log(r))$

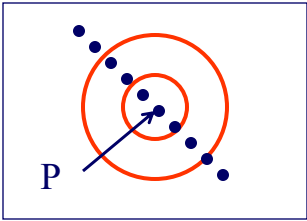
15-826 Copyright: C. Faloutsos (2019) 23

CarnegieMellon

EXPLANATIONS

Intrinsic ('fractal') dimension

- **Local** fractal dimension of point 'P' ?
- A: $mn_p (<= r) \sim r^1$
- If this equation holds for several values of r,
- Then, the **local fractal dimension** of point P:
- Local fd = exp = 1



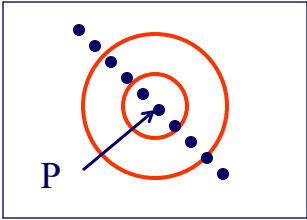
15-826 Copyright: C. Faloutsos (2019) 24

CarnegieMellon

EXPLANATIONS

Intrinsic ('fractal') dimension

- **Local** fractal dimension of point 'A' ?
- A: $nn_P (\leq r) \sim r^1$
- If this is true for all points of the cloud
- Then the exponent is the **global** f.d.
- Or simply the f.d.



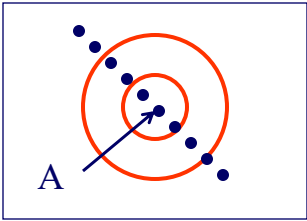
15-826 Copyright: C. Faloutsos (2019) 25

CarnegieMellon

EXPLANATIONS

Intrinsic ('fractal') dimension

- **Global** fractal dimension?
- A: if
- $\sum_{all_P} [nn_P (\leq r)] \sim r^1$
- Then: $exp = global\ f.d.$
- If this is true for all points of the cloud
- Then the exponent is the **global** f.d.
- Or simply the f.d.



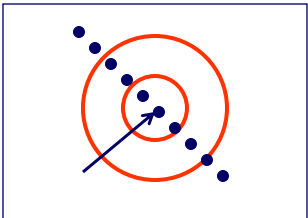
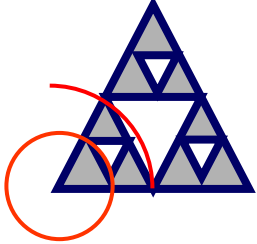
15-826 Copyright: C. Faloutsos (2019) 26

CarnegieMellon

EXPLANATIONS

Intrinsic ('fractal') dimension

- **Local** fractal dimension for sierpinski triangle ?

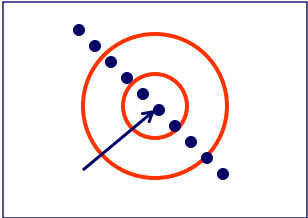
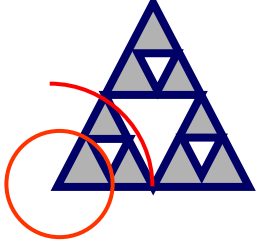
15-826 Copyright: C. Faloutsos (2019) 27

CarnegieMellon

EXPLANATIONS

Intrinsic ('fractal') dimension

- **Local** fractal dimension for sierpinski triangle ?
- 2x radius, 3x points
- $n = r ^ { (\log 3 / \log 2)}$

15-826 Copyright: C. Faloutsos (2019) 28

Intrinsic ('fractal') dimension

- Algorithm, to estimate it?

Notice

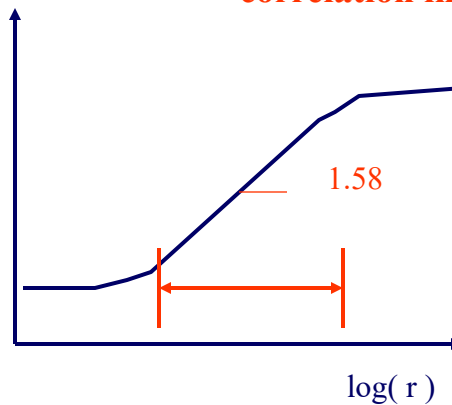
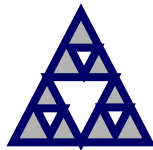
- $Sum_{all_P} [nn_P (<=r)]$ is exactly $tot\#pairs(<=r)$

including 'mirror' pairs

Sierpinsky triangle

== 'correlation integral'

$\log(\#pairs$
within $\leq r$)



Observations:

- Euclidean objects have **integer** fractal dimensions
 - point: 0
 - lines and smooth curves: 1
 - smooth surfaces: 2
- fractal dimension \rightarrow roughness of the periphery

Important properties

- $fd =$ embedding dimension \rightarrow uniform pointset
- a point set may have several fd , depending on scale



CarnegieMellon

Important properties

- fd = embedding dimension \rightarrow uniform pointset
- a point set may have several fd, depending on scale



2-d

15-826

Copyright: C. Faloutsos (2019)

33

CarnegieMellon

Important properties

- fd = embedding dimension \rightarrow uniform pointset
- a point set may have several fd, depending on scale



1-d

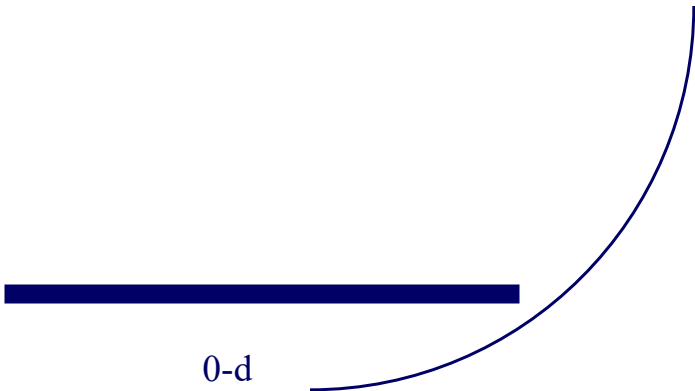
15-826

Copyright: C. Faloutsos (2019)

34

CarnegieMellon

Important properties



0-d

15-826 Copyright: C. Faloutsos (2019) 35

CarnegieMellon

Road map

- Motivation – 3 problems / case studies
- Definition of fractals and power laws
- ➔ • Solutions to posed problems
- More examples and tools
- Discussion - putting fractals to work!
- Conclusions – practitioner's guide
- Appendix: gory details - boxcounting plots

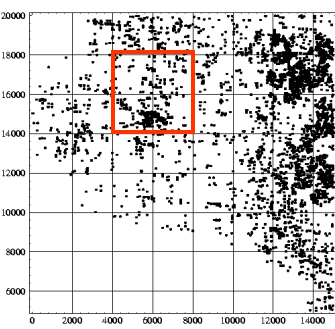
15-826 Copyright: C. Faloutsos (2019) 36

CarnegieMellon

Problem #1: GIS points

Cross-roads of Montgomery county:

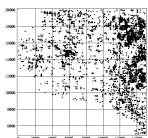
- any rules?



15-826 Copyright: C. Faloutsos (2019) 37

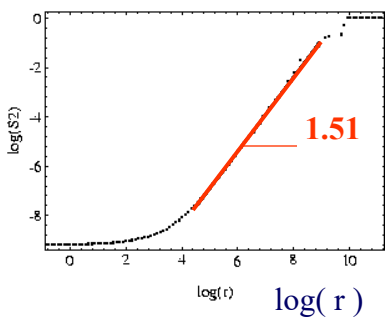
CarnegieMellon

Solution #1



$\log(\#\text{pairs}(\text{within } \leq r))$

SLOPE = 1.51847



$\log(r)$ $\log(r)$

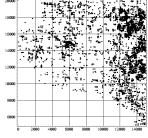
A: self-similarity ->

- \Leftrightarrow fractals
- \Leftrightarrow scale-free
- \Leftrightarrow power-laws
($y=x^a$, $F=C*r^{-2}$)
- $\text{avg}\#\text{neighbors}(\leq r) = r^D$

15-826 Copyright: C. Faloutsos (2019) 38

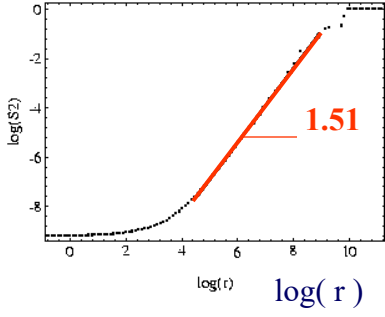
CarnegieMellon

Solution #1



$\log(\#\text{pairs}(\text{within } \leq r))$

SLOPE = 1.51847



$\log(r)$

A: self-similarity

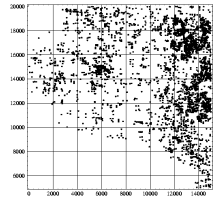
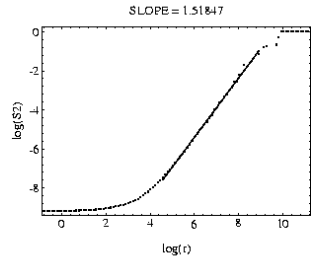
- $\text{avg}\#\text{neighbors}(\leq r) \sim r^{1.51}$

15-826 Copyright: C. Faloutsos (2019) 39

CarnegieMellon

Examples:MG county

- Montgomery County of MD (road end-points)

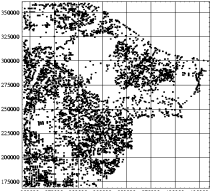
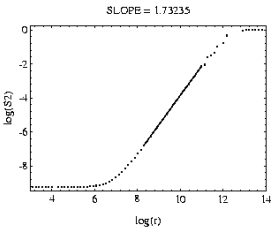
$\log(r)$

15-826 Copyright: C. Faloutsos (2019) 40

CarnegieMellon

Examples:LB county

- Long Beach county of CA (road end-points)

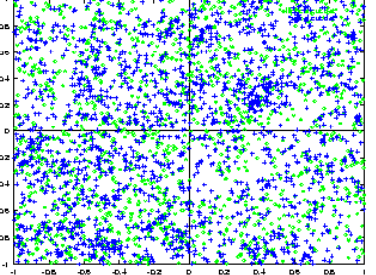
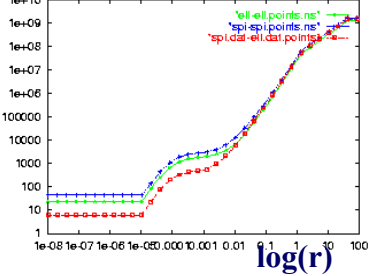



15-826
Copyright: C. Faloutsos (2019)
41

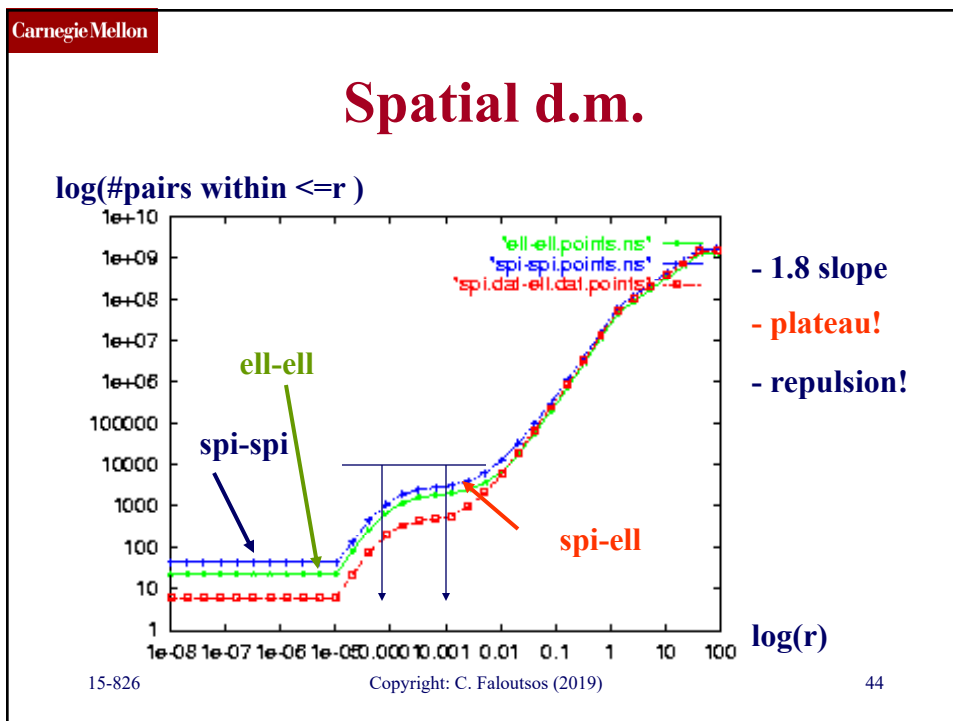
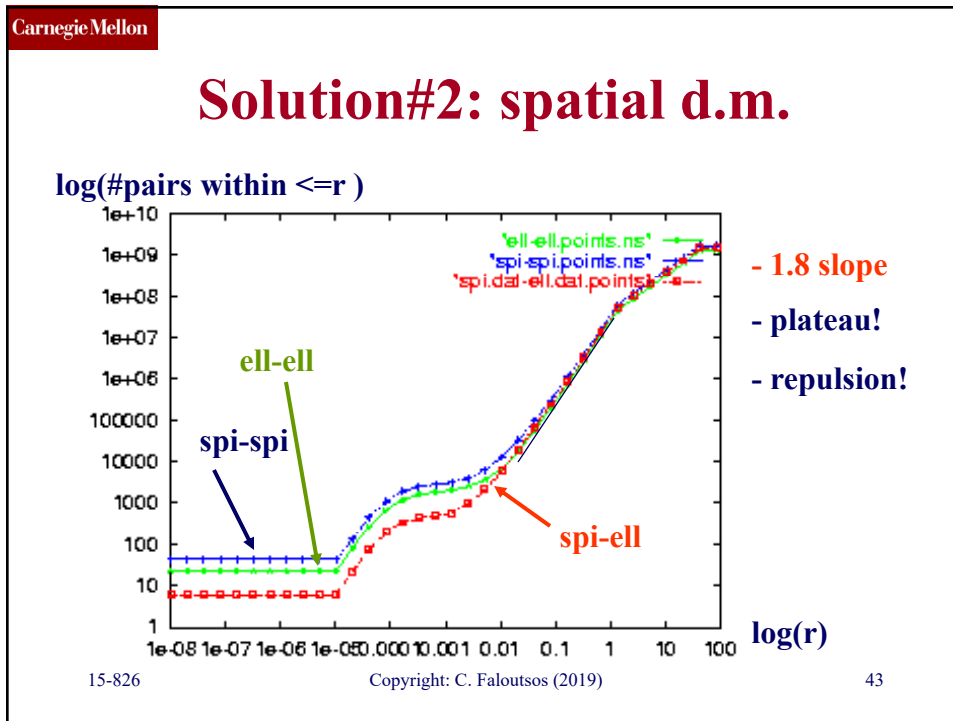
CarnegieMellon

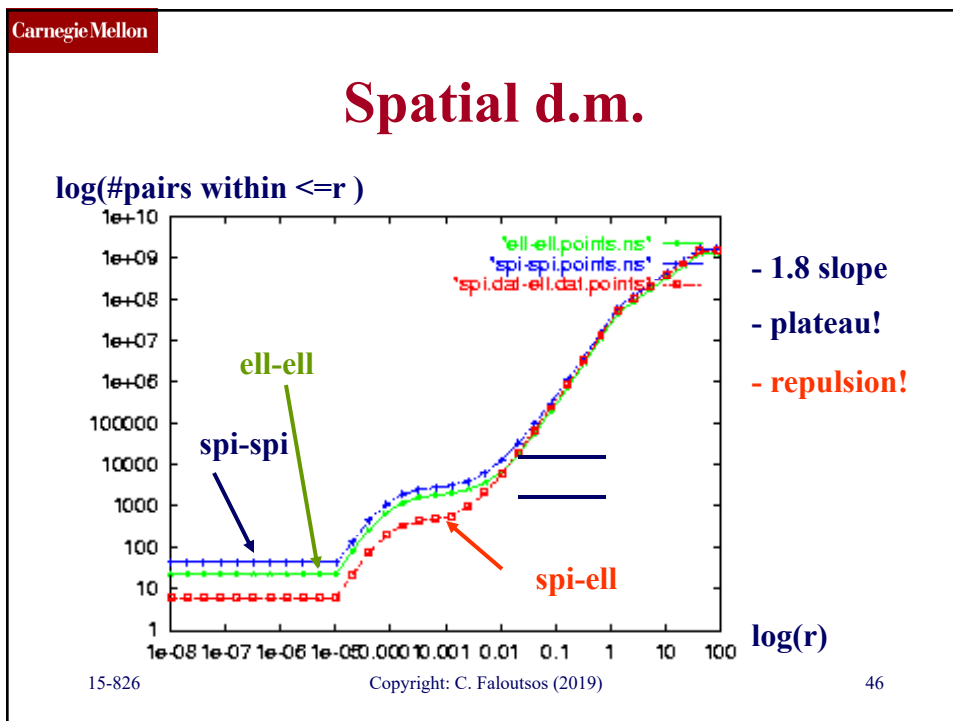
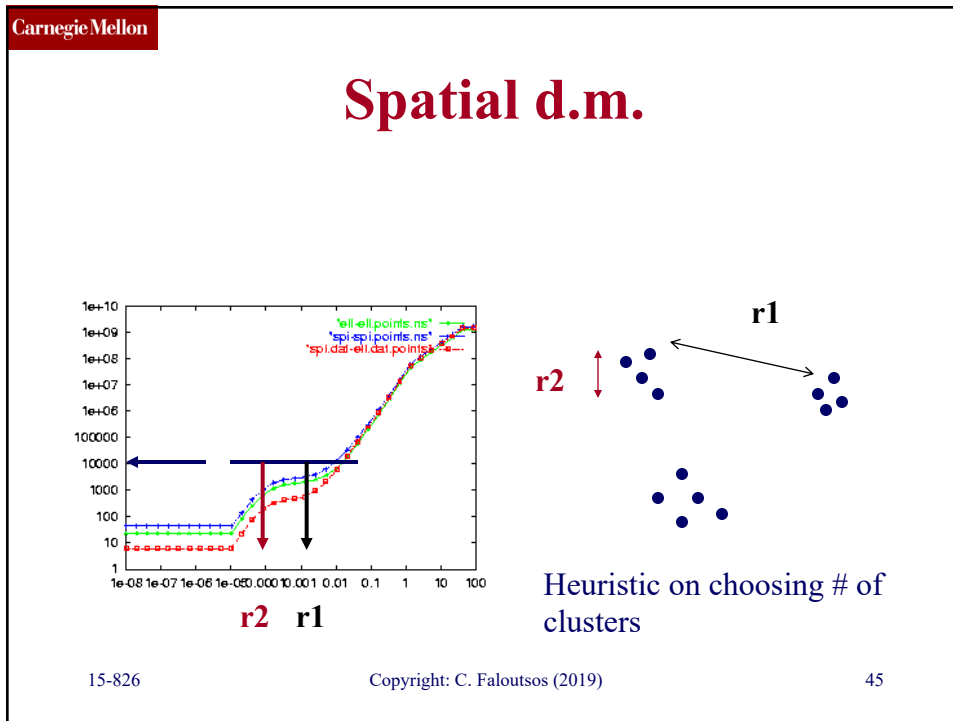
Solution#2: spatial d.m.

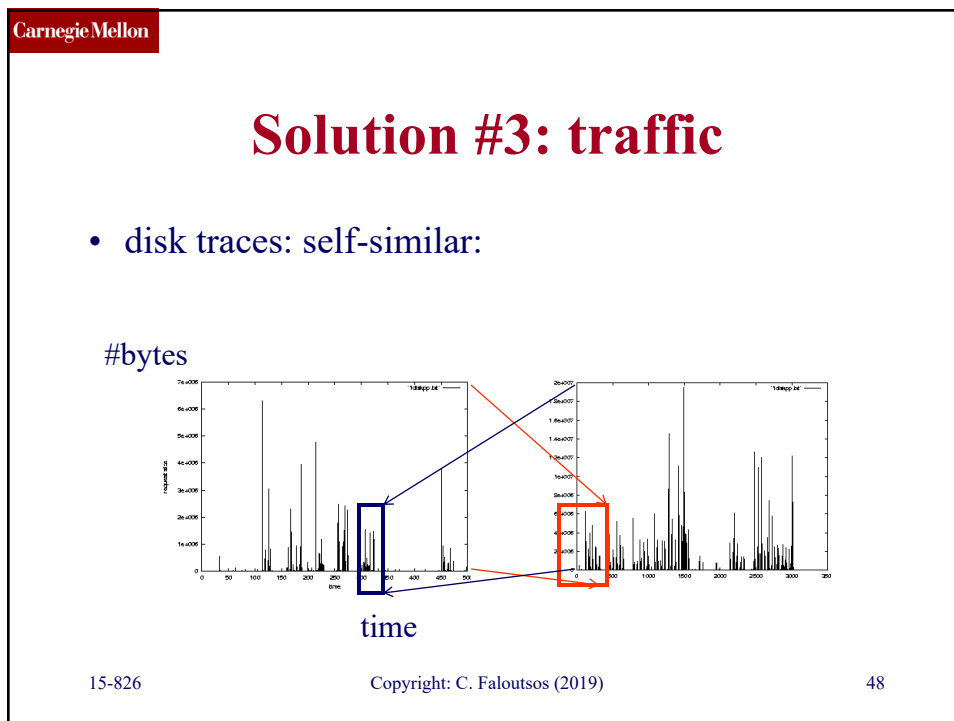
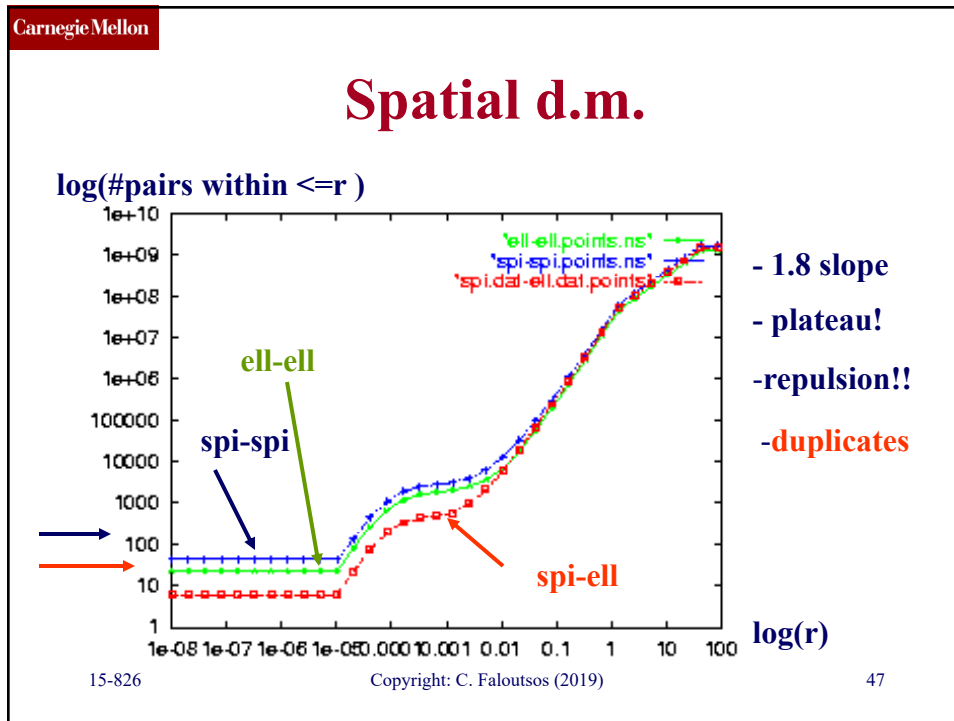
Galaxies (‘BOPS’ plot - [sigmoid2000])

15-826
Copyright: C. Faloutsos (2019)
42







CarnegieMellon

Solution #3: traffic

- disk traces (80-20 'law' = 'multifractal')

#bytes

time

20% ↘ ↙ 80%

15-826 Copyright: C. Faloutsos (2019) 49

CarnegieMellon

80-20 / multifractals

20 ↘ ↙ 80

15-826 Copyright: C. Faloutsos (2019) 50

CarnegieMellon

80-20 / multifractals

- $p ; (1-p)$ in general
- **yes, there are dependencies**

15-826
Copyright: C. Faloutsos (2019)
51

CarnegieMellon

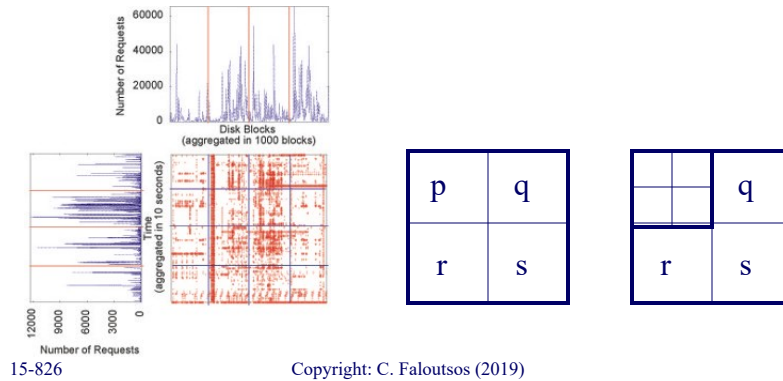
More on 80/20: PQRS

- Part of 'self-* storage' project [Wang+'02]

15-826
Copyright: C. Faloutsos (2019)
52

More on 80/20: PQRS

- Part of 'self-* storage' project [Wang+' 02]



Solution#3: traffic

Clarification:

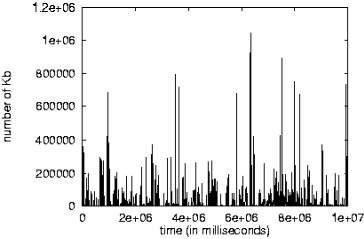
- fractal: a set of points that is self-similar
- multifractal: a probability density function that is self-similar

Many other time-sequences are
bursty/clustered: (such as?)

CarnegieMellon

Example:

- network traffic



<http://repository.cs.vt.edu/lbl-conn-7.tar.Z>

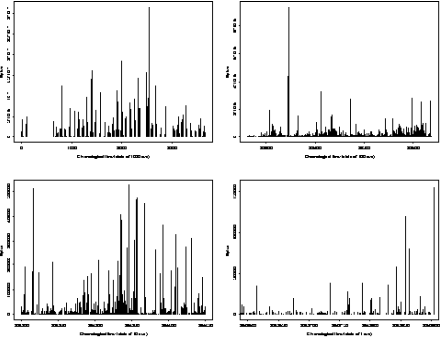
15-826 Copyright: C. Faloutsos (2019) 55

CarnegieMellon

Web traffic

- [Crovella Bestavros, SIGMETRICS' 96]

1000 sec; 100sec
10sec; 1sec



15-826 Copyright: C. Faloutsos (2019) 56

CarnegieMellon

Tape accesses

Tape#1 Tape# N

↔ ↔

★ ★ ★ ★

time

tapes needed, to retrieve n records?

(# days down, due to failures / hurricanes / communication noise...)

15-826
Copyright: C. Faloutsos (2019)
57

CarnegieMellon

Tape accesses

Tape#1 Tape# N

↔ ↔

★ ★ ★ ★

time

tapes retrieved

50-50 = Poisson

qual. records

15-826
Copyright: C. Faloutsos (2019)
58

CarnegieMellon

Road map

- Motivation – 3 problems / case studies
- Definition of fractals and power laws
- Solutions to posed problems
- ➔ • More **tools** and examples
- Discussion - putting fractals to work!
- Conclusions – practitioner's guide
- Appendix: gory details - boxcounting plots

15-826

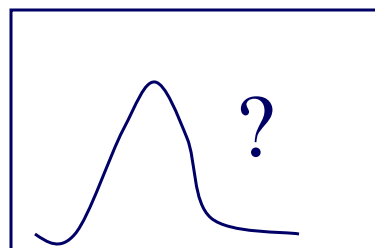
Copyright: C. Faloutsos (2019)

59

CarnegieMellon

A counter-intuitive example

count



avg: 3.3

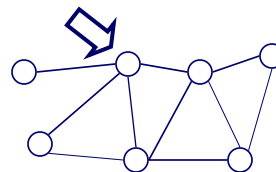
degree

15-826

Copyright: C. Faloutsos (2019)

60

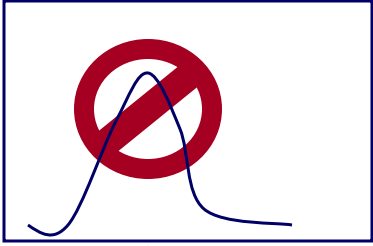
- avg degree is, say 3.3
- pick a node at random – guess its degree, exactly (-> “mode”)



CarnegieMellon

A counter-intuitive example

count



avg: 3.3 degree

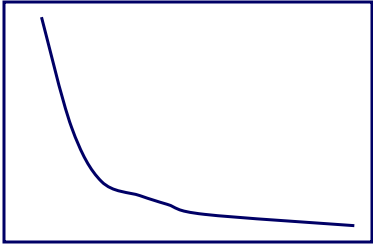
- avg degree is, say 3.3
- pick a node at random – guess its degree, exactly (-> “mode”)
- A: 1!!

15-826 Copyright: C. Faloutsos (2019) 61

CarnegieMellon

A counter-intuitive example

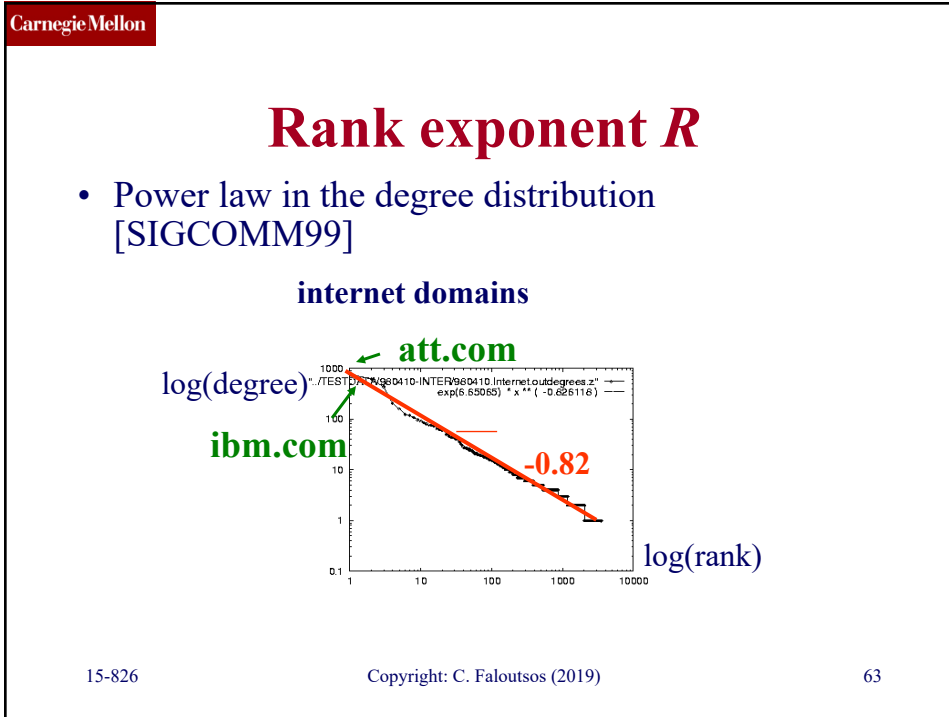
count



avg: 3.3 degree

- avg degree is, say 3.3
- pick a node at random - what is the degree you expect it to have?
- A: 1!!
- A' : very skewed distr.
- Corollary: **the mean is meaningless!**
- (and std -> infinity (!))

15-826 Copyright: C. Faloutsos (2019) 62



CarnegieMellon

A famous power law: Zipf's law

- Q: vocabulary word frequency in a document - any pattern?

15-826 Copyright: C. Faloutsos (2019) 65

CarnegieMellon

A famous power law: Zipf's law

- Bible - rank vs frequency (log-log)

15-826 Copyright: C. Faloutsos (2019) 66

CarnegieMellon

A famous power law: Zipf's law

log(rank)

- Bible - rank vs frequency (log-log)
- similarly, in **many other** languages; for customers and sales volume; city populations etc etc

15-826
Copyright: C. Faloutsos (2019)
67

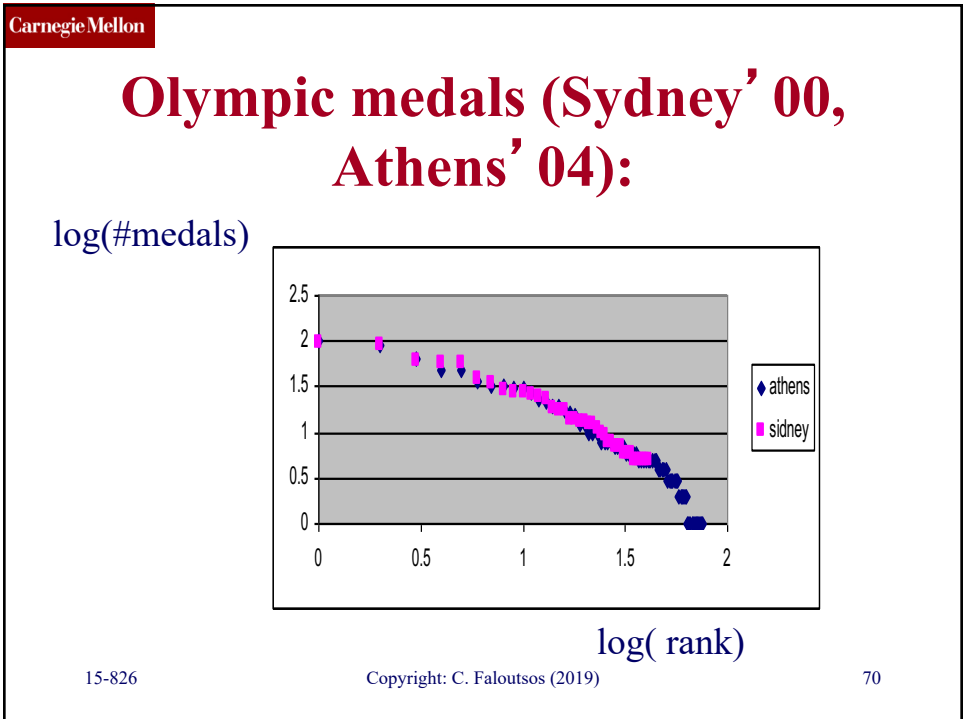
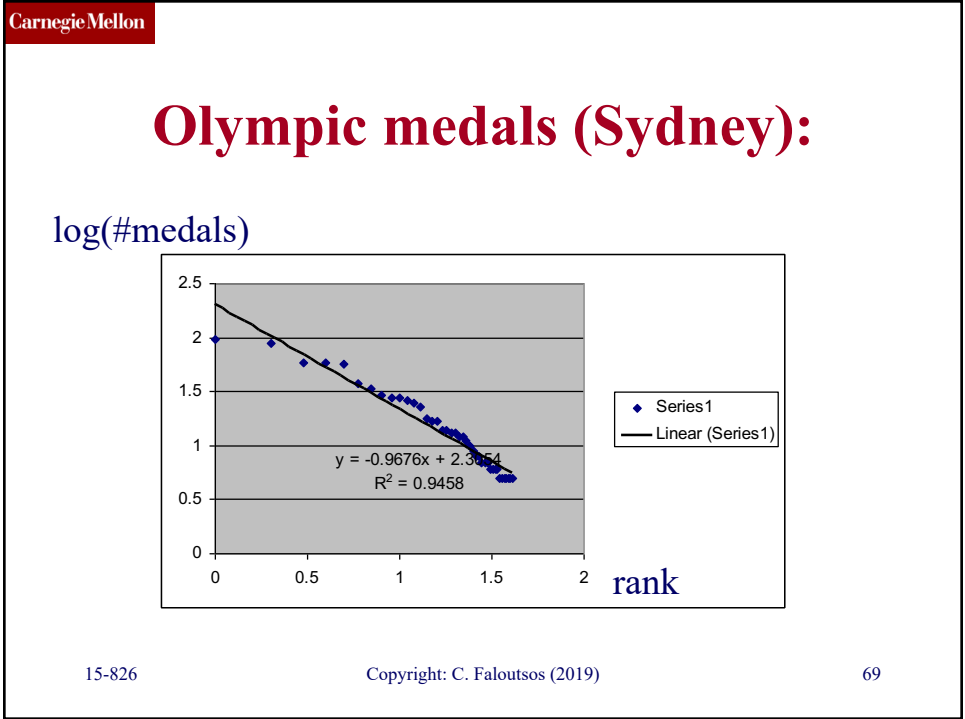
CarnegieMellon

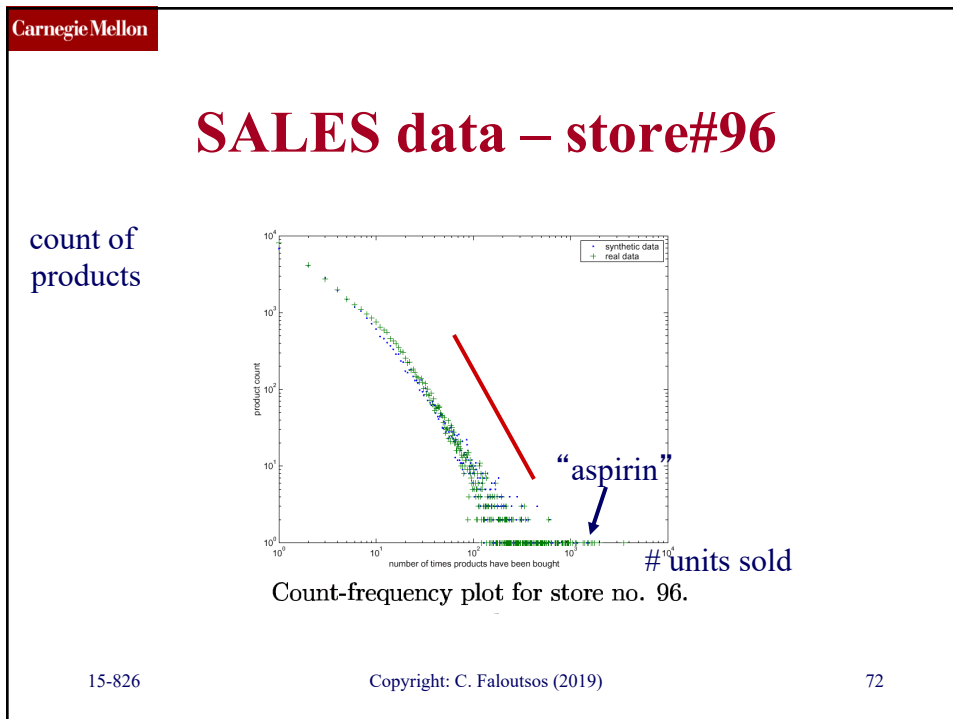
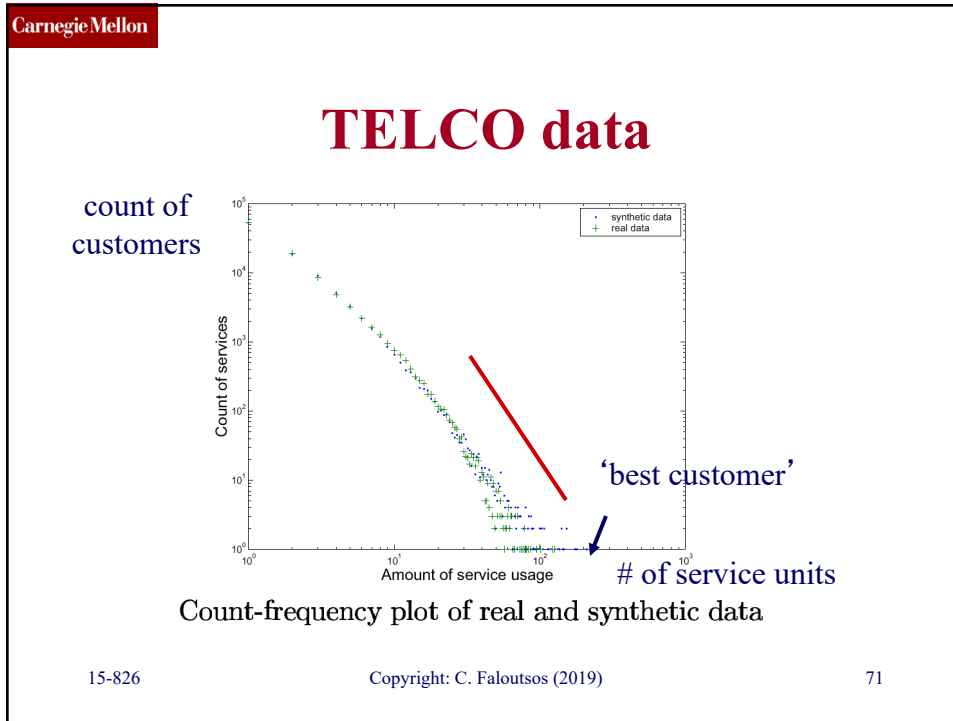
A famous power law: Zipf's law

log(rank)

- Zipf distr:
 $freq = 1 / rank$
- generalized Zipf:
 $freq = 1 / (rank)^a$


15-826
Copyright: C. Faloutsos (2019)
68





CarnegieMellon

More power laws: areas – Korcak's law




Scandinavian lakes
Any pattern?

15-826 Copyright: C. Faloutsos (2019) 73

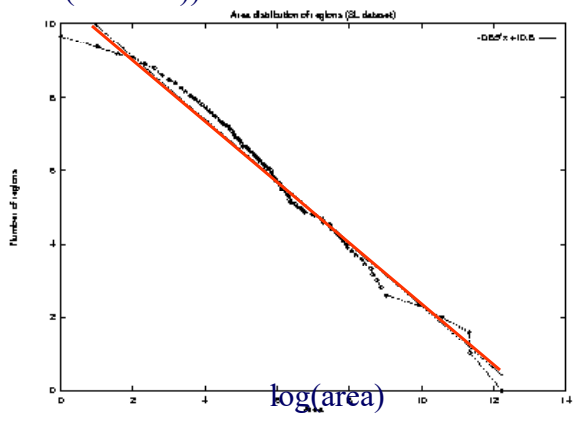
CarnegieMellon

More power laws: areas – Korcak's law

$\log(\text{count}(\geq \text{area}))$



Scandinavian lakes
area vs
complementary
cumulative count
(log-log axes)



Area distribution of regions (SL dataset)

Number of regions

$\log(\text{area})$


-0.62 ± 0.05

15-826 Copyright: C. Faloutsos (2019) 74

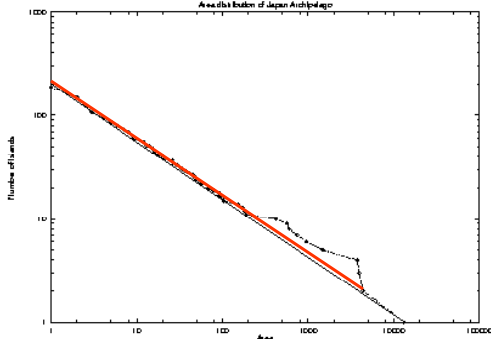
CarnegieMellon

More power laws: Korcak

$\log(\text{count}(\geq \text{area}))$



Japan islands;
area vs cumulative
count (log-log axes)



Number of Islands


Area

$\log(\text{area})$


15-826 Copyright: C. Faloutsos (2019) 75

CarnegieMellon

(Korcak's law: Aegean islands)




Macedonia
Thessaly
Evia & Sporades Islands
North Aegean Islands
Cyclades Islands
Dodecanese Islands
Crete
Peloponnese
Attica
Central Greece
Ionian Islands
Pirios



15-826 Copyright: C. Faloutsos (2019) 76

CarnegieMellon

Korcak's law & "fat fractals"



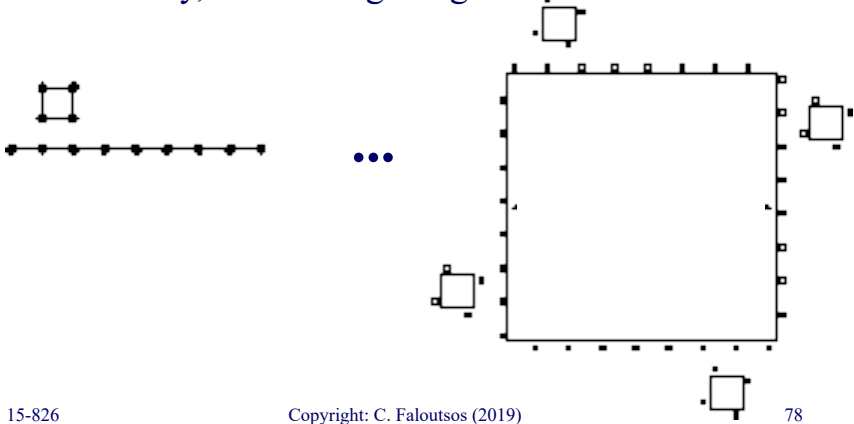
How to generate such regions?

15-826 Copyright: C. Faloutsos (2019) 77

CarnegieMellon

Korcak's law & "fat fractals"

Q: How to generate such regions?
 A: recursively, from a single region



15-826 Copyright: C. Faloutsos (2019) 78

CarnegieMellon

so far we' ve seen:

- concepts:
 - fractals, multifractals and fat fractals
- tools:
 - correlation integral (= pair-count plot)
 - rank/frequency plot (Zipf' s law)
 - CCDF (Korcak' s law)

15-826 Copyright: C. Faloutsos (2019) 79

CarnegieMellon

so far we' ve seen:

- concepts:
 - fractals, multifractals and fat fractals
- tools:
 - correlation integral (= pair-count plot)
 - rank/frequency plot (Zipf' s law)
 - CCDF (Korcak' s law)

same info

15-826 Copyright: C. Faloutsos (2019) 80

CarnegieMellon

Next:

- More examples / applications
- Practitioner's guide
- Box-counting: fast estimation of correlation integral

15-826

Copyright: C. Faloutsos (2019)

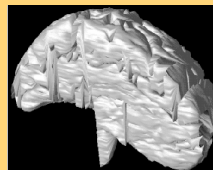
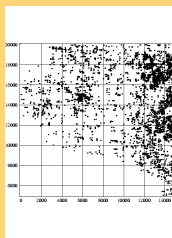
81

CarnegieMellon



Problem

- What patterns are in real k -dim points?



15-826

Copyright: C. Faloutsos (2019)

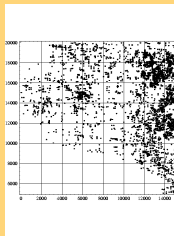
82

CarnegieMellon

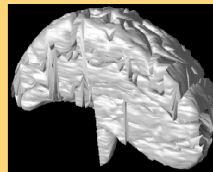


Conclusions

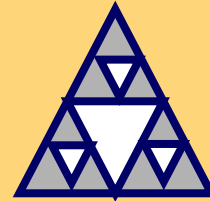
- What patterns are in real k -dim points?
- Self-similarity (= fractals \rightarrow power laws)



15-826



Copyright: C. Faloutsos (2019)



83