



15-721 DB Sys. Design & Impl.

Fractals

Christos Faloutsos

www.cs.cmu.edu/~christos



Roadmap

- 1) Roots: System R and Ingres
- 2) Implementation: buffering, indexing, q-opt
- <...>
- 7) Data Analysis - data mining
 - data cubes / data warehousing / OLAP
 - Association Rules
 - SVD
- Fractals
 - Streams, time sequences, wavelets
- 8) Benchmarks
- 9) vision statements

15-721

C. Faloutsos

2



Citations

[pods94] Christos Faloutsos and Ibrahim Kamel, *Beyond Uniformity and Independence: Analysis of R-trees Using the Concept of Fractal Dimension*, PODS, Minneapolis, MN, May 24-26, 1994, pp. 4-13
<http://www.cs.cmu.edu/~christos/PUBLICATIONS.OLDER/pods94.ps.gz>

Alberto Belussi and Christos Faloutsos *Estimating the Selectivity of Spatial Queries Using the 'Correlation' Fractal Dimension* VLDB, Sept. 1995, Zurich, Switzerland, pp. 299-310
<http://www.cs.cmu.edu/~christos/PUBLICATIONS.OLDER/vldb95.ps.gz>

15-721

C. Faloutsos

3



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

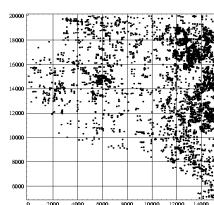
15-721

C. Faloutsos

4



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??

15-721

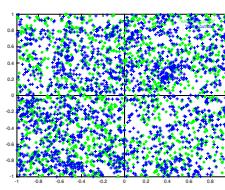
C. Faloutsos

5



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'?

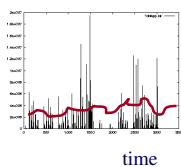
15-721

C. Faloutsos

6

Problem #3: traffic

- disk trace (from HP - J. Wilkes); Web traffic - fit a model
#bytes



Poisson

- how many
explosions to
expect?
- queue length
distr.?

15-721

C. Faloutsos

7

Common answer:

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

15-721

C. Faloutsos

8

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

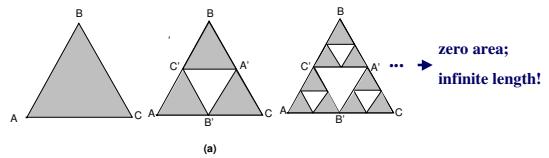
15-721

C. Faloutsos

9

What is a fractal?

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



15-721

C. Faloutsos

10

Definitions (cont'd)

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

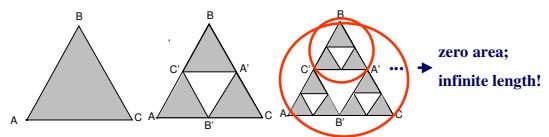
15-721

C. Faloutsos

11

Dfn of fd:

ONLY for a perfectly self-similar point set:



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

15-721

C. Faloutsos

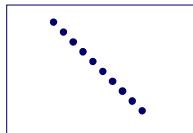
12



CMU SCS

Intrinsic ('fractal') dimension

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



15-721

C. Faloutsos

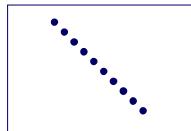
13



CMU SCS

Intrinsic ('fractal') dimension

- Q: dfn for a given set of points?



15-721

C. Faloutsos

14

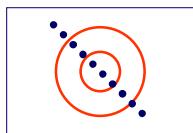
x	y
5	1
4	2
3	3
2	4



CMU SCS

Intrinsic ('fractal') dimension

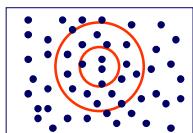
- Q: fractal dimension of a line?
- A: $nn (\leq r) \sim r^1$
('power law': $y=x^a$)
- Q: fd of a plane?
- A: $nn (\leq r) \sim r^2$
fd == slope of $(\log(nn) \text{ vs } \log(r))$



15-721

C. Faloutsos

15



CMU SCS

Intrinsic ('fractal') dimension

- Algorithm, to estimate it?
- Notice
- avg $nn(\leq r)$ is exactly $\text{tot}\#\text{pairs}(\leq r) / (N)$

15-721

C. Faloutsos

16

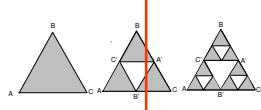


CMU SCS

Sierpinsky triangle

== 'correlation integral'

log(#pairs
within $\leq r$)



15-721

C. Faloutsos

17



CMU SCS

Observations:

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



15-721

C. Faloutsos

18

Important properties

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

15-721

C. Faloutsos

19

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

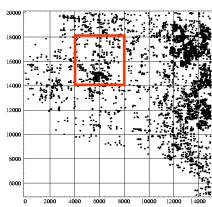
15-721

C. Faloutsos

20

Problem #1: GIS points

Cross-roads of Montgomery county:
• any rules?

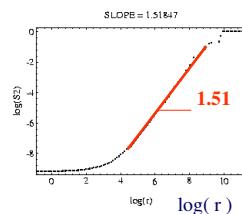


15-721

C. Faloutsos

21

Solution #1

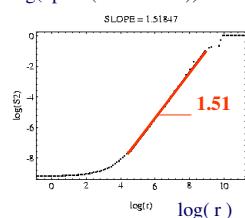
 $\log(\#\text{pairs}(\text{within } \leq r))$ 

15-721

C. Faloutsos

22

Solution #1

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

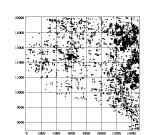
15-721

C. Faloutsos

23

Examples: MG county

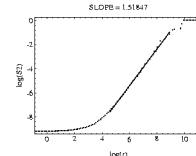
- Montgomery County of MD (road endpoints)



15-721

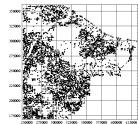
C. Faloutsos

24

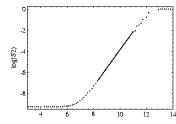


Examples:LB county

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



15-721



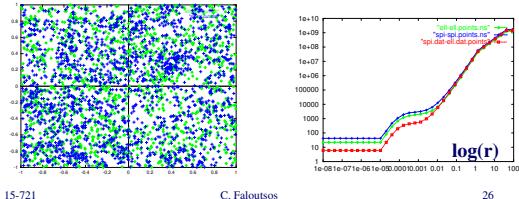
C. Faloutsos

25

Solution#2: spatial d.m.

Galaxies ('BOPS' plot - [sigmod2000])

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

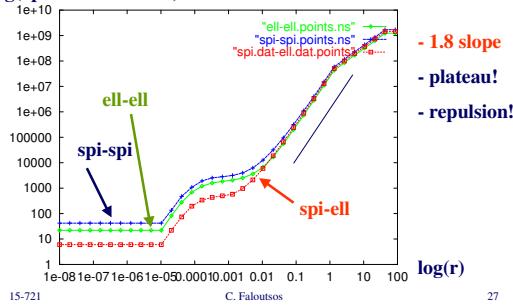


15-721 C. Faloutsos

26

Solution#2: spatial d.m.

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



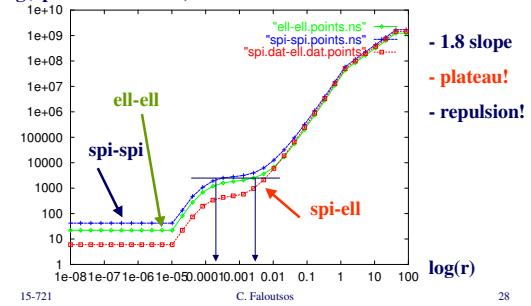
15-721

C. Faloutsos

27

spatial d.m.

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

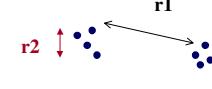
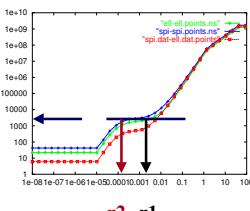


15-721

C. Faloutsos

28

spatial d.m.



Heuristic on choosing # of clusters

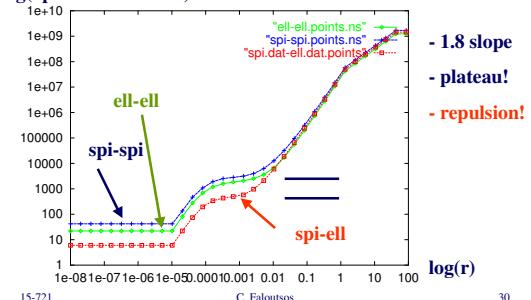
15-721

C. Faloutsos

29

spatial d.m.

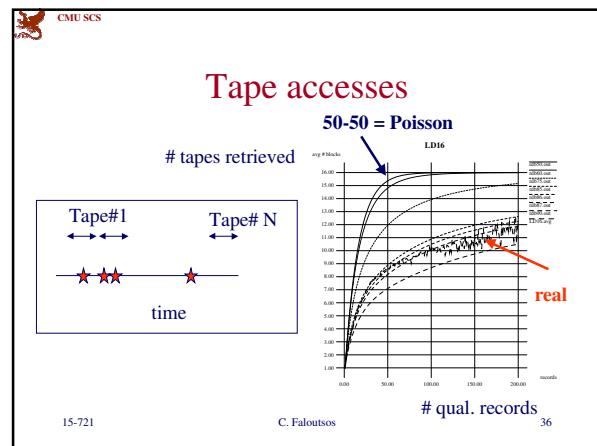
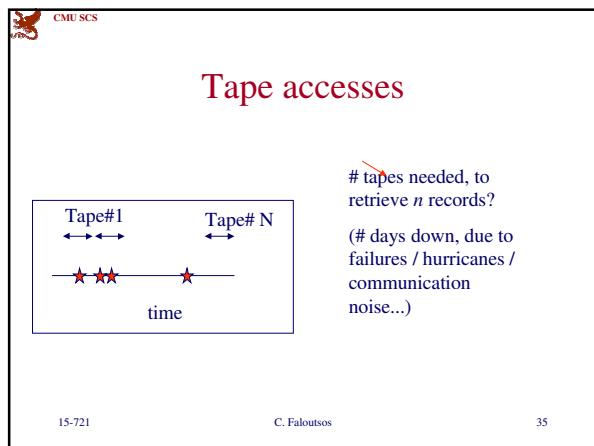
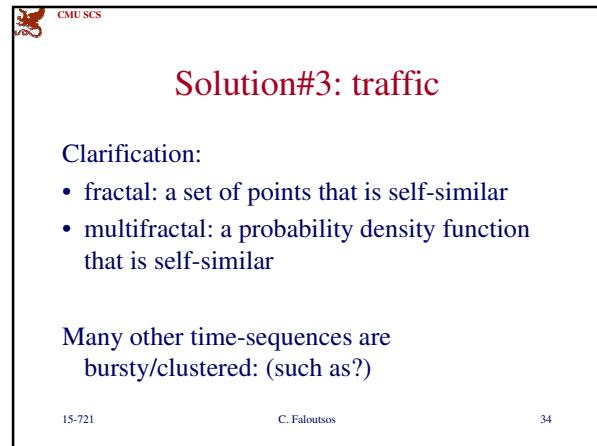
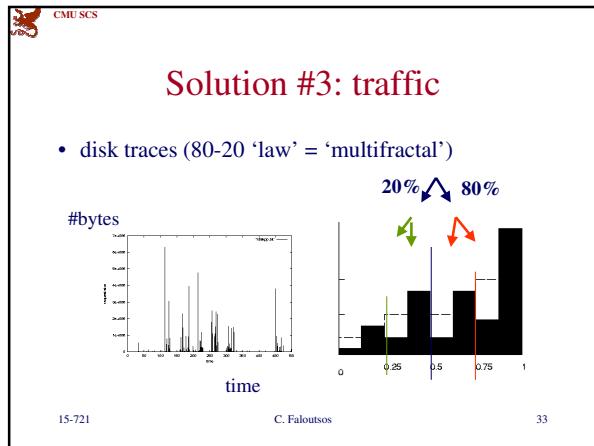
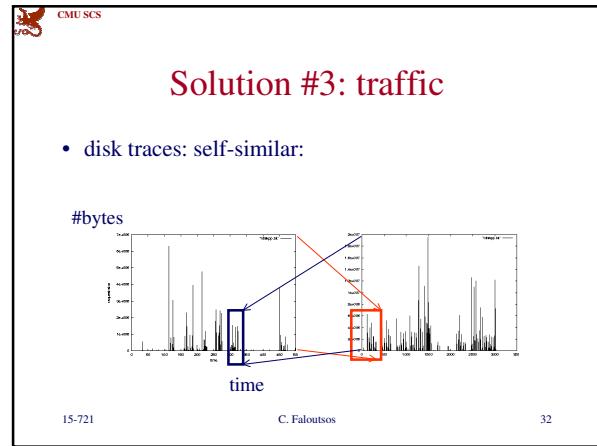
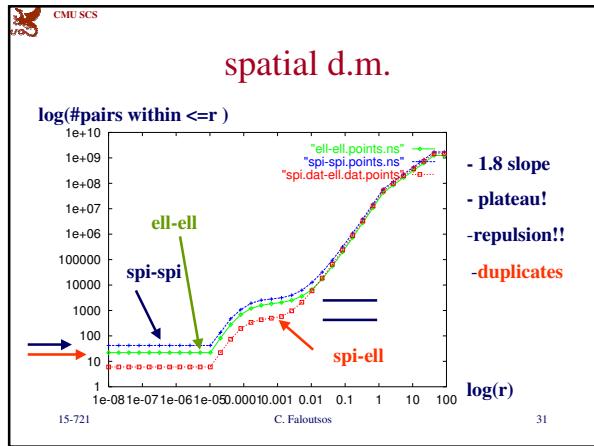
$\log(\#\text{pairs within } \leq r)$



15-721

C. Faloutsos

30



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

15-721

C. Faloutsos

37

More tools

- Zipf's law
- Korcak's law / "fat fractals"

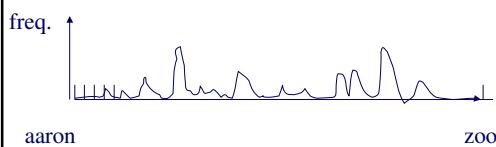
15-721

C. Faloutsos

38

A famous power law: Zipf's law

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



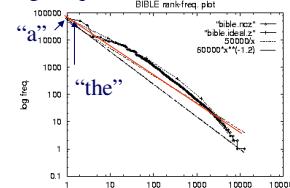
15-721

C. Faloutsos

39

A famous power law: Zipf's law

log(freq)



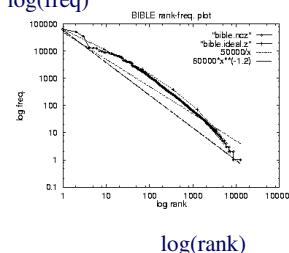
log(rank)

15-721

C. Faloutsos

40

A famous power law: Zipf's law



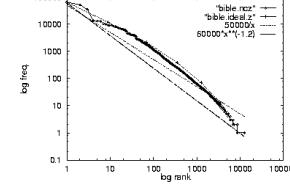
15-721

C. Faloutsos

41

A famous power law: Zipf's law

log(freq)



log(rank)

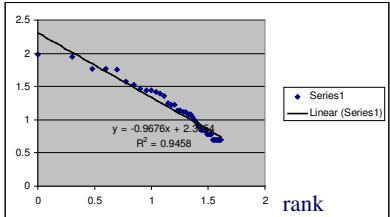
15-721

C. Faloutsos

42

Olympic medals (Sidney):

log(#medals)



15-721

C. Faloutsos

43

More power laws: areas – Korcak's law

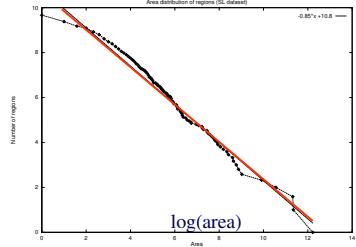
Scandinavian lakes
Any pattern?

15-721

C. Faloutsos

44

More power laws: areas – Korcak's law

log(count(\geq area))

15-721

C. Faloutsos

45

More power laws: Korcak

Japan islands



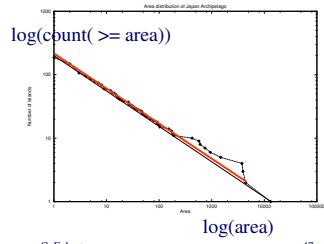
15-721

C. Faloutsos

46

More power laws: Korcak

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



C. Faloutsos

47

(Korcak's law: Aegean islands)



15-721

C. Faloutsos

48

Korcak's law & "fat fractals"



How to generate such regions?

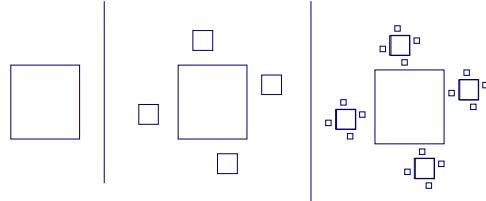
15-721

C. Faloutsos

49

Korcak's law & "fat fractals"

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



15-721

C. Faloutsos

50

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)

15-721

C. Faloutsos

51

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

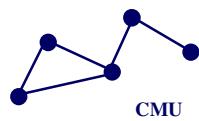
15-721

C. Faloutsos

52

Other applications: Internet

- How does the internet look like?



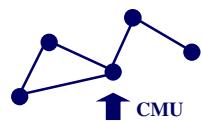
15-721

C. Faloutsos

53

Other applications: Internet

- How does the internet look like?
- Internet routers: how many neighbors within h hops?



15-721

C. Faloutsos

54

(reminder: our tool-box:)

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

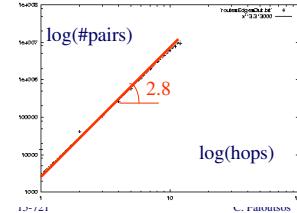
15-721

C. Faloutsos

55

Internet topology

- Internet routers: how many neighbors within h hops?

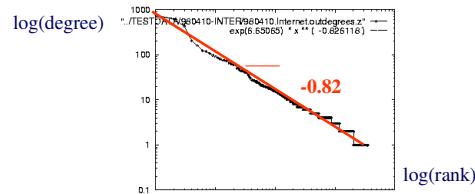


Reachability function:
number of neighbors
within r hops, vs r (log-log).

Mbone routers, 1995

56

More power laws on the Internet



degree vs rank, for Internet domains
(log-log) [sigcomm99]

15-721

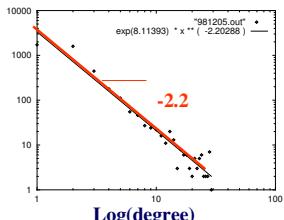
C. Faloutsos

57

More power laws - internet

- pdf of degrees: (slope: 2.2)

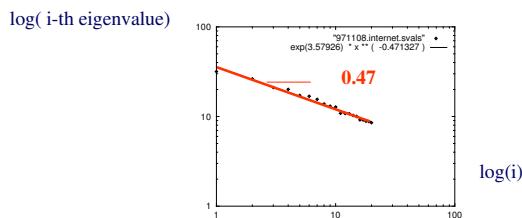
Log(count)



C. Faloutsos

58

Even more power laws on the Internet



Scree plot for Internet domains (log-log) [sigcomm99]

15-721

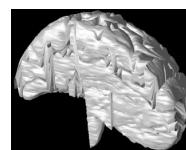
C. Faloutsos

59

More apps: Brain scans

- Oct-trees; brain-scans

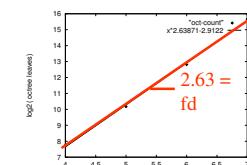
Log(#octants)



15-721

C. Faloutsos

octree levels 60



More apps: Medical images

[Burdett et al, SPIE '93]:

- benign tumors: $fd \sim 2.37$
- malignant: $fd \sim 2.56$

15-721

C. Faloutsos

61

More fractals:

- cardiovascular system: 3 (!) 
- stock prices (LYCOS) - random walks: 1.5



- Coastlines: 1.2-1.58 (Norway!)

15-721

C. Faloutsos

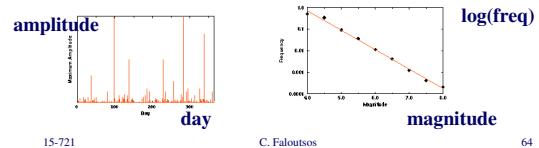
62



15-721

More power laws

- duration of UNIX jobs
- Energy of earthquakes (Gutenberg-Richter law) [simsience.org]



15-721

C. Faloutsos

64

Even more power laws:

- publication counts (Lotka's law)
- Distribution of UNIX file sizes
- Income distribution (Pareto's law)
- web hit counts [Huberman]

15-721

C. Faloutsos

65

Power laws, cont'd

- In- and out-degree distribution of web sites [Barabasi], [IBM-CLEVER]
- length of file transfers [Bestavros+]
- Click-stream data (w/ A. Montgomery (CMU-GSIA) + MediaMetrix)

15-721

C. Faloutsos

66

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

15-721

C. Faloutsos

67

Settings for fractals:

Points; areas (-> fat fractals), eg:

15-721

C. Faloutsos

68

Settings for fractals:

Points; areas, eg:

- cities/stores/hospitals, over earth's surface
- time-stamps of events (customer arrivals, packet losses, criminal actions) over time
- regions (sales areas, islands, patches of habitats) over space

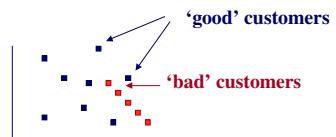
15-721

C. Faloutsos

69

Settings for fractals:

- customer feature vectors (age, income, frequency of visits, amount of sales per visit)



15-721

C. Faloutsos

70

Some uses of fractals:

- Detect non-existence of rules (if points are uniform)
- Detect non-homogeneous regions (eg., legal login time-stamps may have different fd than intruders')
- Estimate number of neighbors / customers / competitors within a radius

15-721

C. Faloutsos

71

Multi-Fractals

Setting: points or objects, w/ some value, eg:

- cities w/ populations
- positions on earth and amount of gold/water/oil underneath
- product ids and sales per product
- people and their salaries
- months and count of accidents

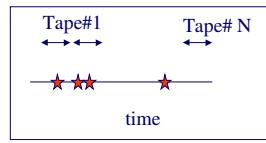
15-721

C. Faloutsos

72

Use of multifractals:

- Estimate tape/disk accesses
 - how many of the 100 tapes contain my 50 phonecall records?*
 - how many days without an accident?*

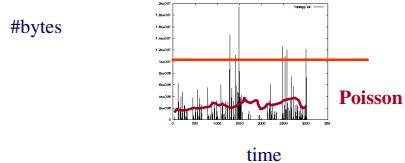


15-721

73

Use of multifractals

- how often do we exceed the threshold?



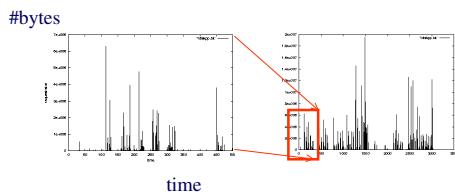
15-721

C. Faloutsos

74

Use of multifractals cont'd

- Extrapolations for/from samples



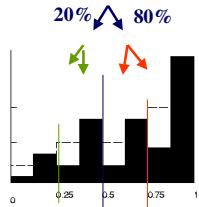
15-721

C. Faloutsos

75

Use of multifractals cont'd

- How many distinct products account for 90% of the sales?*



15-721

C. Faloutsos

76

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

15-721

C. Faloutsos

77

Conclusions

- Real data often **disobey** textbook assumptions (Gaussian, Poisson, uniformity, independence)
 - avoid ‘mean’ - use median, or even better, use: fractals, self-similarity, and power laws, to find patterns - specifically:

15-721

C. Faloutsos

78

Conclusions

- tool#1: (for points) ‘correlation integral’:** (#pairs within $\leq r$) vs (distance r)
- tool#2: (for categorical values) rank-frequency plot (a’la Zipf)**
- tool#3: (for numerical values) CCDF:** Complementary cumulative distr. function (#of elements with value $\geq a$)

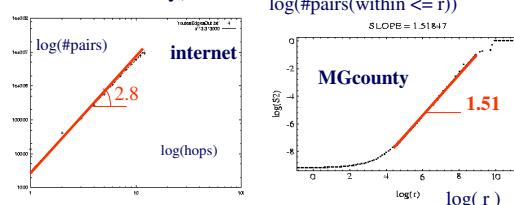
15-721

C. Faloutsos

79

Practitioner’s guide:

- tool#1: #pairs vs distance, for a set of objects,** with a distance function (slope = intrinsic dimensionality)



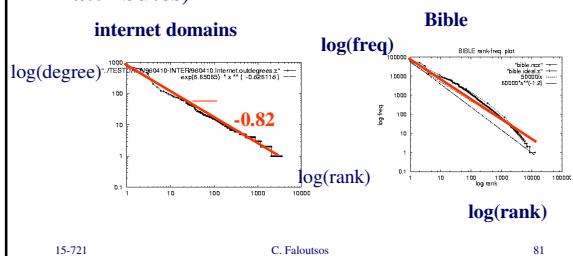
15-721

C. Faloutsos

80

Practitioner’s guide:

- tool#2: rank-frequency plot (for categorical attributes)**



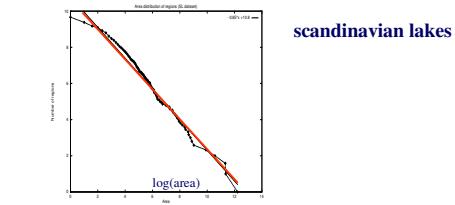
15-721

C. Faloutsos

81

Practitioner’s guide:

- tool#3: CCDF, for (skewed) numerical attributes, eg. areas of islands/lakes, UNIX jobs...)**

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

15-721

C. Faloutsos

82

Resources:

- Software for fractal dimension
 - <http://www.cs.cmu.edu/~christos>
 - christos@cs.cmu.edu

15-721

C. Faloutsos

83

Books

- Strongly recommended intro book:
 - Manfred Schroeder *Fractals, Chaos, Power Laws: Minutes from an Infinite Paradise* W.H. Freeman and Company, 1991
- Classic book on fractals:
 - B. Mandelbrot *Fractal Geometry of Nature*, W.H. Freeman, 1977

15-721

C. Faloutsos

84

References

- [ieeeTN94] W. E. Leland, M.S. Taqqu, W. Willinger, D.V. Wilson, *On the Self-Similar Nature of Ethernet Traffic*, IEEE Transactions on Networking, 2, 1, pp 1-15, Feb. 1994.
- [pods94] Christos Faloutsos and Ibrahim Kamel, *Beyond Uniformity and Independence: Analysis of R-trees Using the Concept of Fractal Dimension*, PODS, Minneapolis, MN, May 24-26, 1994, pp. 4-13

15-721

C. Faloutsos

85

References

- [vlbd95] Alberto Belussi and Christos Faloutsos, *Estimating the Selectivity of Spatial Queries Using the 'Correlation' Fractal Dimension* Proc. of VLDB, p. 299-310, 1995
- [vlbd96] Christos Faloutsos, Yossi Matias and Avi Silberschatz, *Modeling Skewed Distributions Using Multifractals and the '80-20 Law'* Conf. on Very Large Data Bases (VLDB), Bombay, India, Sept. 1996.

15-721

C. Faloutsos

86

References

- [vlbd96] Christos Faloutsos and Volker Gaede *Analysis of the Z-Ordering Method Using the Hausdorff Fractal Dimension* VLD, Bombay, India, Sept. 1996
- [sigcomm99] Michalis Faloutsos, Petros Faloutsos and Christos Faloutsos, *What does the Internet look like? Empirical Laws of the Internet Topology*, SIGCOMM 1999

15-721

C. Faloutsos

87

References

- [icde99] Guido Proietti and Christos Faloutsos, *I/O complexity for range queries on region data stored using an R-tree* International Conference on Data Engineering (ICDE), Sydney, Australia, March 23-26, 1999
- [sigmod2000] Christos Faloutsos, Bernhard Seeger, Agma J. M. Traina and Caetano Traina Jr., *Spatial Join Selectivity Using Power Laws*, SIGMOD 2000

15-721

C. Faloutsos

88

Appendix - Gory details

- Bad news: There are more than one fractal dimensions
 - Minkowski fd; Hausdorff fd; Correlation fd; Information fd
- Great news:
 - they can all be computed fast!
 - they usually have nearby values

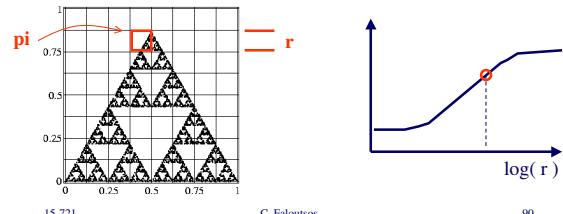
15-721

C. Faloutsos

89

Fast estimation of fd(s):

- How, for the (correlation) fractal dimension?
- A: Box-counting plot: $\log(\sum(\pi_i^2))$



15-721

C. Faloutsos

90

Definitions

- p_i : the percentage (or count) of points in the i -th cell
- r : the side of the grid

15-721

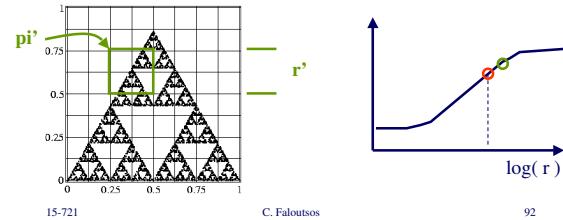
C. Faloutsos

91

Fast estimation of fd(s):

- compute $\text{sum}(p_i^2)$ for another grid side, r'

$$\log(\text{sum}(p_i^2))$$



15-721

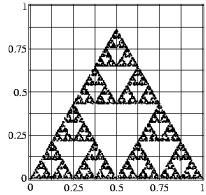
C. Faloutsos

92

Fast estimation of fd(s):

- etc; if the resulting plot has a linear part, its slope is the correlation fractal dimension D_2

$$\log(\text{sum}(p_i^2))$$



15-721

C. Faloutsos

93

Definitions (cont'd)

- Many more fractal dimensions D_q (related to Renyi entropies):

$$D_q = \frac{1}{q-1} \frac{\partial \log(\sum p_i^q)}{\partial \log(r)} \quad q \neq 1$$

$$D_1 = \frac{\partial \sum p_i \log(p_i)}{\partial \log(r)}$$

15-721

C. Faloutsos

94

Hausdorff or box-counting fd:

- Box counting plot: $\log(N(r))$ vs $\log(r)$
- r : grid side
- $N(r)$: count of non-empty cells
- (Hausdorff) fractal dimension D_0 :

$$D_0 = -\frac{\partial \log(N(r))}{\partial \log(r)}$$

15-721

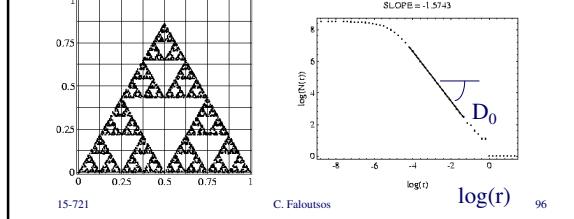
C. Faloutsos

95

Definitions (cont'd)

- Hausdorff fd:

$$r ___ \log(\#\text{non-empty cells})$$



15-721

C. Faloutsos

96

Observations

- $q=0$: Hausdorff fractal dimension
- $q=2$: Correlation fractal dimension
(**identical** to the exponent of the number of neighbors vs radius)
- $q=1$: Information fractal dimension

15-721

C. Faloutsos

97

Observations, cont'd

- in general, the D_q 's take similar, but not identical, values.
- except for perfectly self-similar point-sets, where $D_q=D_{q'}$ for any q, q'

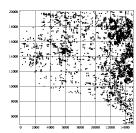
15-721

C. Faloutsos

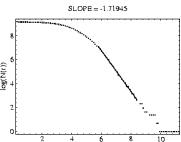
98

Examples:MG county

- Montgomery County of MD (road end-points)



15-721

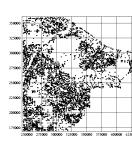


C. Faloutsos

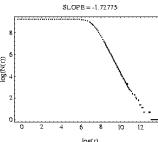
99

Examples:LB county

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



15-721



C. Faloutsos

100

Conclusions

- many fractal dimensions, with nearby values
- can be computed quickly
($O(N)$ or $O(N \log(N))$)
- (code: on the web)

15-721

C. Faloutsos

101