

## Powerful Tools for Data Mining Fractals, Power laws, SVD

*C. Faloutsos*  
Carnegie Mellon University

## PART II

### Fractals

[www.cs.cmu.edu/~christos/](http://www.cs.cmu.edu/~christos/)

### Intro to fractals - outline

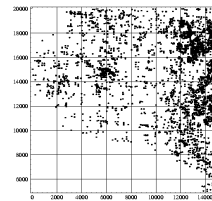
- ➔ • Motivation – 3 problems / case studies
- Definition of fractals and power laws
- Solutions to posed problems
- More examples and tools
- Discussion - putting fractals to work!
- Conclusions – practitioner's guide
- Appendix: gory details - boxcounting plots

MDIC 04

Copyright: C. Faloutsos (2004)

3

### 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??

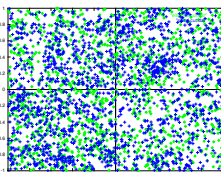
MDIC 04

Copyright: C. Faloutsos (2004)

4

### 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'?

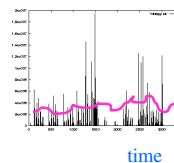
MDIC 04

Copyright: C. Faloutsos (2004)

5

### Problem #3: traffic

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



Poisson

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

MDIC 04

Copyright: C. Faloutsos (2004)

6

## Common answer:

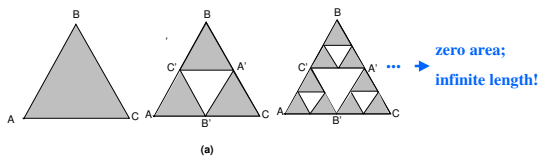
- Fractals / self-similarities / power laws
- Seminal works from Hilbert, Minkowski, Cantor, Mandelbrot, (Hausdorff, Lyapunov, Ken Wilson, ...)

## 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

## What is a fractal?

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

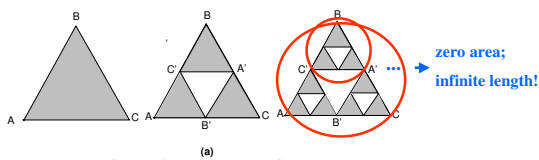


## Definitions (cont'd)

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

## Dfn of fd:

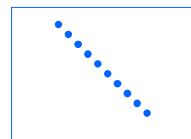
ONLY for a perfectly self-similar point set:



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

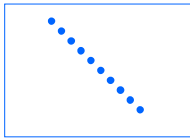
## Intrinsic ('fractal') dimension

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



### Intrinsic ('fractal') dimension

- Q: dfn for a given set of points?



x	y
5	1
4	2
3	3
2	4

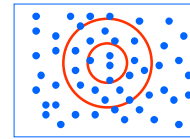
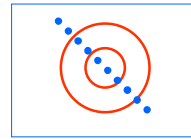
MDIC 04

Copyright: C. Faloutsos (2004)

13

### Intrinsic ('fractal') dimension

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



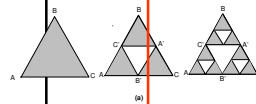
MDIC 04

Copyright: C. Faloutsos (2004)

14

### Intrinsic ('fractal') dimension

- Algorithm, to estimate it?
- Notice
- $avg\ nn(<=r)$  is exactly  $tot\#pairs(<=r) / (N)$



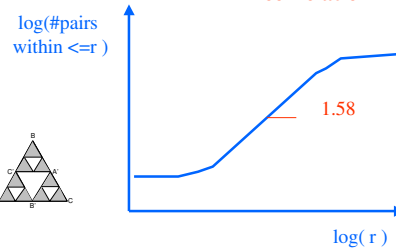
MDIC 04

Copyright: C. Faloutsos (2004)

15

### Sierpinsky triangle

== 'correlation integral'



MDIC 04

Copyright: C. Faloutsos (2004)

16

### Observations:

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



MDIC 04

Copyright: C. Faloutsos (2004)

17

### Important properties

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



MDIC 04

Copyright: C. Faloutsos (2004)

18

## 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

MDIC 04

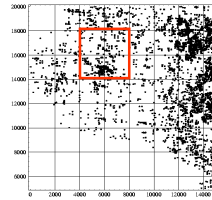
Copyright: C. Faloutsos (2004)

19

## Problem #1: GIS points

Cross-roads of Montgomery county:

- any rules?



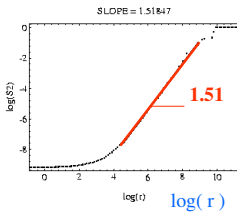
MDIC 04

Copyright: C. Faloutsos (2004)

20

## Solution #1

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



MDIC 04

Copyright: C. Faloutsos (2004)

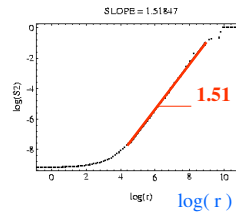
21

A: self-similarity ->

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

## Solution #1

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



MDIC 04

Copyright: C. Faloutsos (2004)

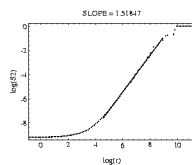
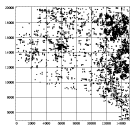
22

A: self-similarity

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

## Examples:MG county

- Montgomery County of MD (road end-points)



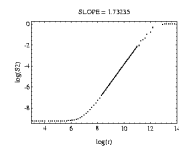
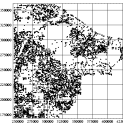
MDIC 04

Copyright: C. Faloutsos (2004)

23

## Examples:LB county

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

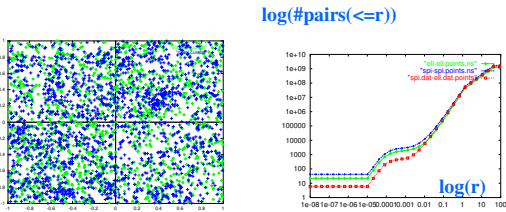


MDIC 04

Copyright: C. Faloutsos (2004)

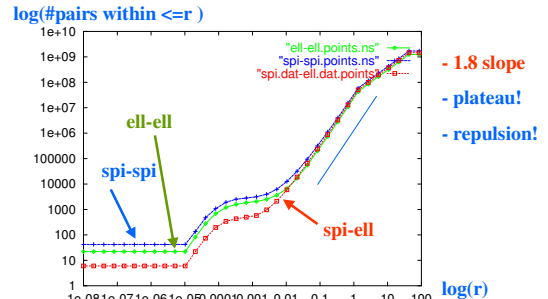
24

### Solution#2: spatial d.m. Galaxies ( 'BOPS' plot - [sigmod2000])



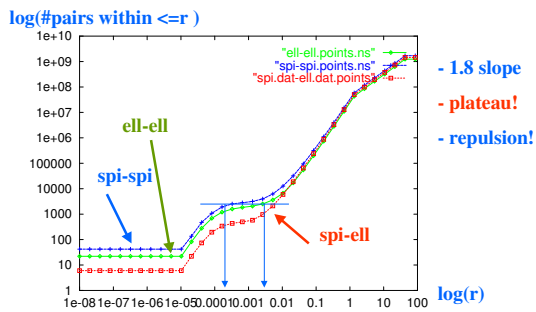
MDIC 04 Copyright: C. Faloutsos (2004) 25

### Solution#2: spatial d.m.



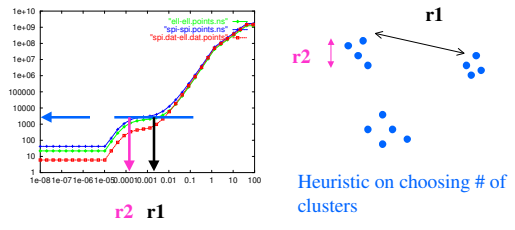
MDIC 04 Copyright: C. Faloutsos (2004) 26

### spatial d.m.



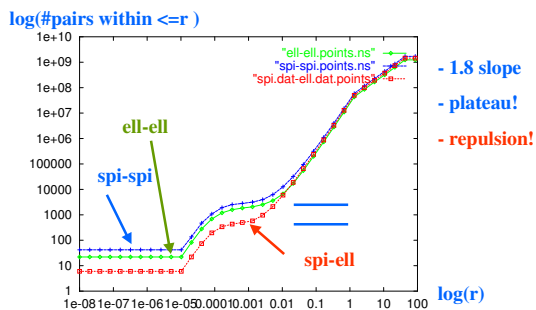
MDIC 04 Copyright: C. Faloutsos (2004) 27

### spatial d.m.



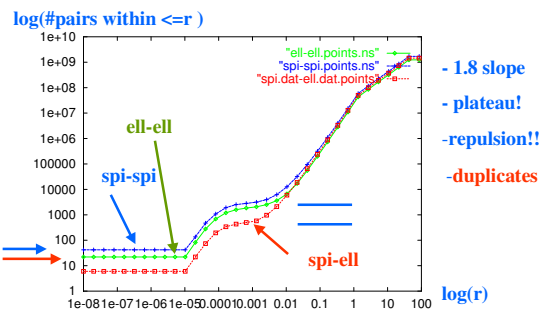
MDIC 04 Copyright: C. Faloutsos (2004) 28

### spatial d.m.



MDIC 04 Copyright: C. Faloutsos (2004) 29

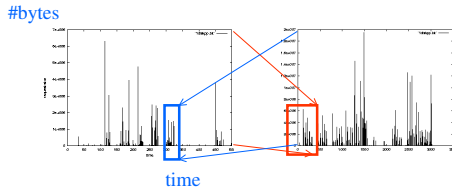
### spatial d.m.



MDIC 04 Copyright: C. Faloutsos (2004) 30

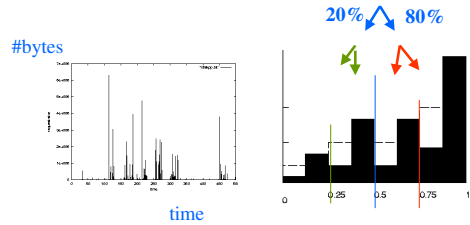
### Solution #3: traffic

- disk traces: self-similar:



### Solution #3: traffic

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



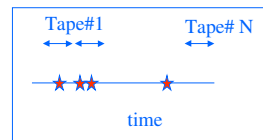
### Solution#3: traffic

Clarification:

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

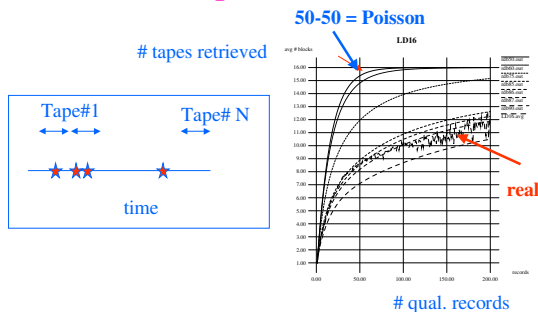
Many other time-sequences are bursty/clustered: (such as?)

### Tape accesses



# tapes needed, to retrieve  $n$  records?  
 (# days down, due to failures / hurricanes / communication noise...)

### Tape accesses



### Road map

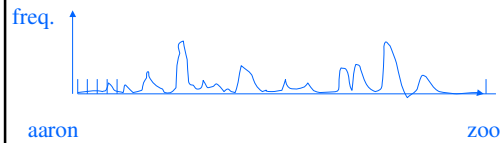
- Motivation – 3 problems / case studies
- Definition of fractals and power laws
- Solutions to posed problems
- ➔ More **tools** and examples
- Discussion - putting fractals to work!
- Conclusions – practitioner's guide
- Appendix: gory details - boxcounting plots

## More tools

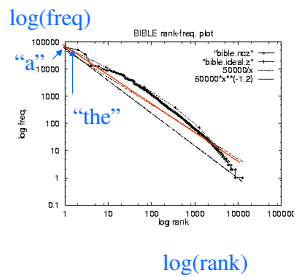
- Zipf's law
- Korcak's law / "fat fractals"

## A famous power law: Zipf's law

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

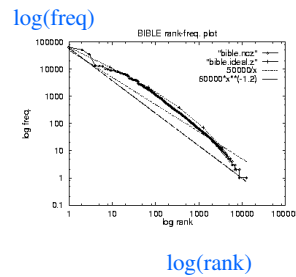


## A famous power law: Zipf's law



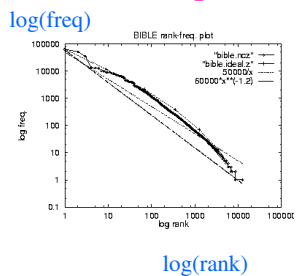
- Bible - rank vs frequency (log-log)

## A famous power law: Zipf's law



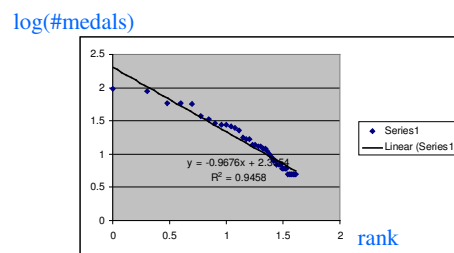
- Bible - rank vs frequency (log-log)
- similarly, in **many other** languages; for customers and sales volume; city populations etc

## A famous power law: Zipf's law




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

## Olympic medals (Sidney):



Carnegie Mellon

## More power laws: areas – Korcak’s law



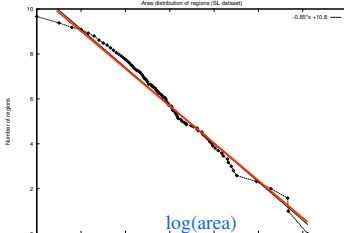
Scandinavian lakes  
Any pattern?

MDIC 04 Copyright: C. Faloutsos (2004) 43

Carnegie Mellon

## More power laws: areas – Korcak’s law

log(count( >= area))




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

MDIC 04 Copyright: C. Faloutsos (2004) 44

Carnegie Mellon

## More power laws: Korcak

Japan islands

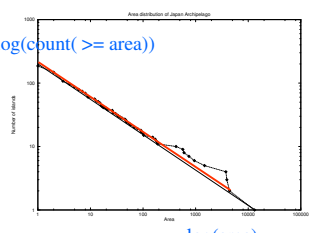


MDIC 04 Copyright: C. Faloutsos (2004) 45

Carnegie Mellon

## More power laws: Korcak

log(count( >= area))





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

MDIC 04 Copyright: C. Faloutsos (2004) 46

Carnegie Mellon

## (Korcak’s law: Aegean islands)





MDIC 04 Copyright: C. Faloutsos (2004) 47

Carnegie Mellon

## Korcak’s law & “fat fractals”

How to generate such regions?



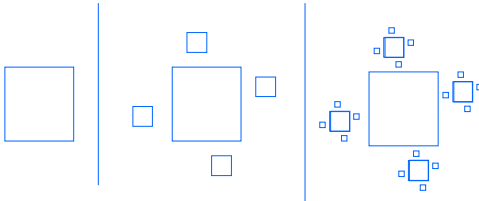
MDIC 04 Copyright: C. Faloutsos (2004) 48



## Korcak's law & "fat fractals"

Q: How to generate such regions?

A: recursively, from a single region



MDIC 04

Copyright: C. Faloutsos (2004)

49

## 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)

MDIC 04

Copyright: C. Faloutsos (2004)

50

## 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

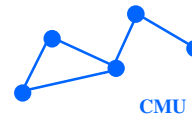
MDIC 04

Copyright: C. Faloutsos (2004)

51

## Other applications: Internet

- How does the internet look like?



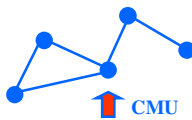
MDIC 04

Copyright: C. Faloutsos (2004)

52

## Other applications: Internet

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



MDIC 04

Copyright: C. Faloutsos (2004)

53

## (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)

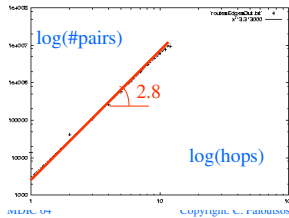
MDIC 04

Copyright: C. Faloutsos (2004)

54

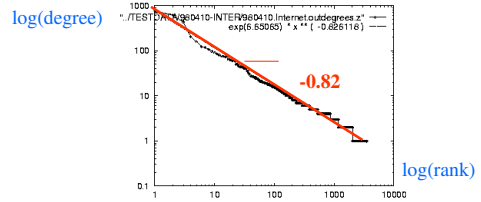
## Internet topology

- Internet routers: how many neighbors within  $h$  hops?



Reachability function:  
number of neighbors within  $r$  hops, vs  $r$  (log-log).  
Mbone routers, 1995

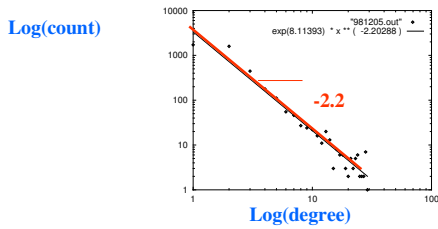
## More power laws on the Internet



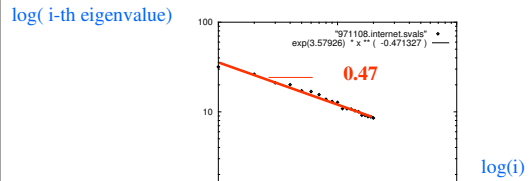
degree vs rank, for Internet domains (log-log) [sigcomm99]

## More power laws - internet

- pdf of degrees: (slope: 2.2)



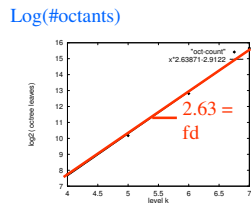
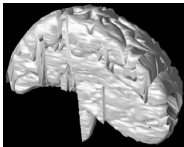
## Even more power laws on the Internet



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

## More apps: Brain scans

- Oct-trees; brain-scans




octree levels

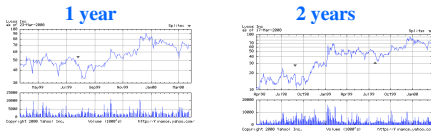
## More apps: Medical images

[Burdett et al, SPIE '93]:

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

## More fractals:

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

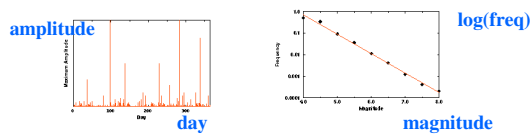


- Coastlines: 1.2-1.58 (Norway!)



## More power laws

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



## Even more power laws:

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

## Power laws, cont'ed

- 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)

## 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

## Settings for fractals:

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

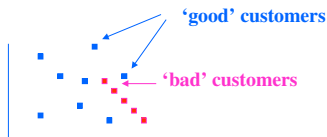
## Settings for fractals:

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

## Settings for fractals:

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



## 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

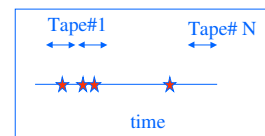
## Multi-Fractals

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

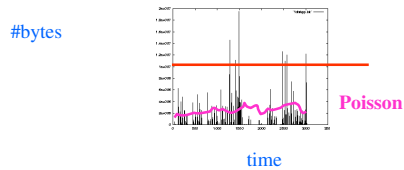
## Use of multifractals:

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



## Use of multifractals

- how often do we exceed the threshold?



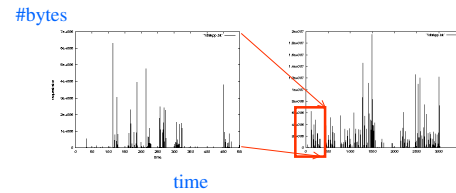
MDIC 04

Copyright: C. Faloutsos (2004)

73

## Use of multifractals cont'd

- Extrapolations for/from samples



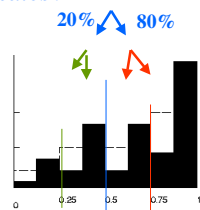
MDIC 04

Copyright: C. Faloutsos (2004)

74

## Use of multifractals cont'd

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



MDIC 04

Copyright: C. Faloutsos (2004)

75

## 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

MDIC 04

Copyright: C. Faloutsos (2004)

76

## 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:

MDIC 04

Copyright: C. Faloutsos (2004)

77

## 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  $\geq a$ )

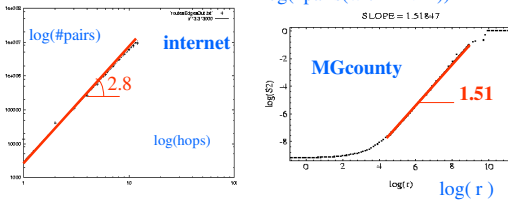
MDIC 04

Copyright: C. Faloutsos (2004)

78

## Practitioner's guide:

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



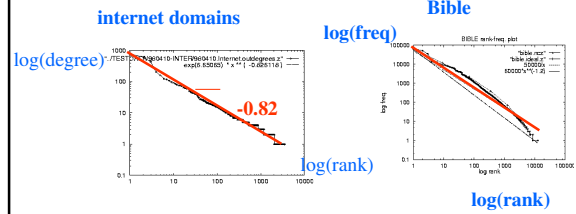
MDIC 04

Copyright: C. Faloutsos (2004)

79

## Practitioner's guide:

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



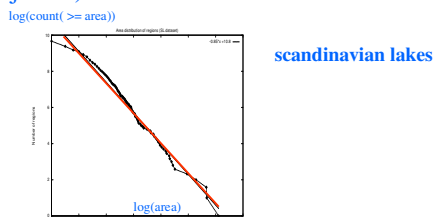
MDIC 04

Copyright: C. Faloutsos (2004)

80

## Practitioner's guide:

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



MDIC 04

Copyright: C. Faloutsos (2004)

81

## Resources:

- Software for fractal dimension
  - <http://www.cs.cmu.edu/~christos>
  - [christos@cs.cmu.edu](mailto:christos@cs.cmu.edu)

MDIC 04

Copyright: C. Faloutsos (2004)

82

## 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

MDIC 04

Copyright: C. Faloutsos (2004)

83

## 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

MDIC 04

Copyright: C. Faloutsos (2004)

84

## References

- [vlb95] Alberto Belussi and Christos Faloutsos, *Estimating the Selectivity of Spatial Queries Using the 'Correlation' Fractal Dimension* Proc. of VLDB, p. 299-310, 1995
- [vlb96] 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.

## References

- [vlb96] 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

## References

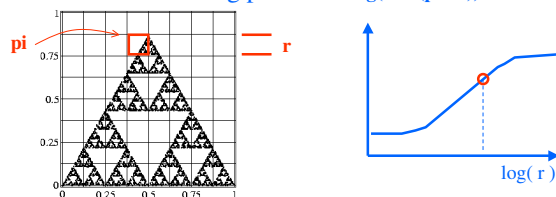
- [icde99] Guido Proietti and Christos Faloutsos, *I/O complexity for range queries on region data stored using an R-tree* International Conference on Data Engineering (ICDE), Sydney, Australia, March 23-26, 1999
- [sigmod2000] Christos Faloutsos, Bernhard Seeger, Agma J. M. Traina and Caetano Traina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

## 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

## Fast estimation of fd(s):

- How, for the (correlation) fractal dimension?
- A: Box-counting plot:

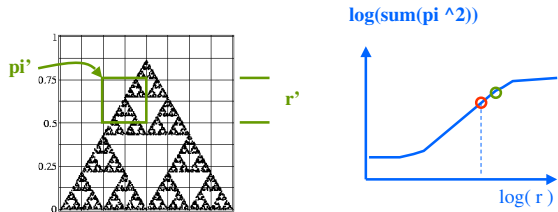


## Definitions

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

### Fast estimation of fd(s):

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



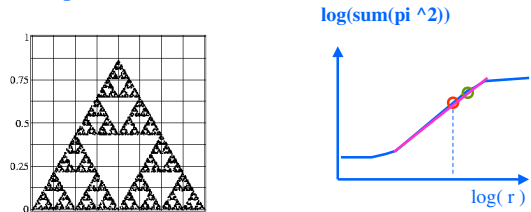
MDIC 04

Copyright: C. Faloutsos (2004)

91

### Fast estimation of fd(s):

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



MDIC 04

Copyright: C. Faloutsos (2004)

92

### Definitions (cont'd)

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

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

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

MDIC 04

Copyright: C. Faloutsos (2004)

93

### 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  $D_0$ :

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

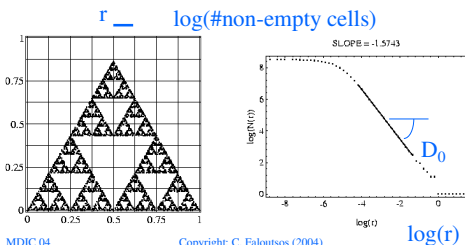
MDIC 04

Copyright: C. Faloutsos (2004)

94

### Definitions (cont'd)

- Hausdorff fd:



MDIC 04

Copyright: C. Faloutsos (2004)

95

### 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

MDIC 04

Copyright: C. Faloutsos (2004)

96



## 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  $q, q'$

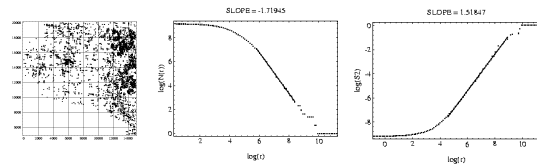
MDIC 04

Copyright: C. Faloutsos (2004)

97

## Examples:MG county

- Montgomery County of MD (road end-points)



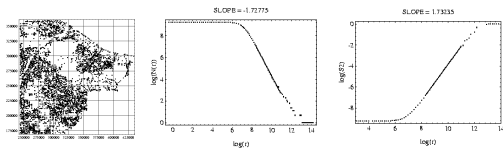
MDIC 04

Copyright: C. Faloutsos (2004)

98

## Examples:LB county

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



MDIC 04

Copyright: C. Faloutsos (2004)

99

## Conclusions

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

MDIC 04

Copyright: C. Faloutsos (2004)

100