


15-826: Multimedia Databases and Data Mining


Lecture #11: Fractals: M-trees and dim. curse (case studies – Part II)
C. Faloutsos



Must-read Material

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

15-826 Copyright: C. Faloutsos (2010) 2



Optional Material

Optional, but **very** useful: Manfred Schroeder
Fractals, Chaos, Power Laws: Minutes from an Infinite Paradise W.H. Freeman and Company, 1991

15-826 Copyright: C. Faloutsos (2010) 3

CMU SCS

Outline

Goal: 'Find similar / interesting things'

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

15-826 Copyright: C. Faloutsos (2010) 4

CMU SCS

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 (2010) 5

CMU SCS

Indexing - Detailed outline

- fractals
 - intro
 - applications
 - disk accesses for R-trees (range queries)
 - dimensionality reduction
 - ➔ • selectivity in M-trees
 - dim. curse revisited
 - "fat fractals"
 - quad-tree analysis [Gaede+]

15-826 Copyright: C. Faloutsos (2010) 6

CMU SCS

What else can they solve?

- ✓ separability [KDD'02]
 - forecasting [CIKM'02]
- ✓ dimensionality reduction [SBBD'00]
 - non-linear axis scaling [KDD'02]
- ✓ disk trace modeling [Wang+'02]
- ▶ selectivity of spatial/multimedia queries [PODS'94, VLDB'95, ICDE'00]
- ...

15-826 Copyright: C. Faloutsos (2010) 7

CMU SCS

Metric trees - analysis

- Problem: How many disk accesses, for an M-tree?
- Given:
 - N (# of objects)
 - C (fanout of disk pages)
 - r (radius of range query - BIASED model)

15-826 Copyright: C. Faloutsos (2010) 8

CMU SCS

Metric trees - analysis

- Problem: How many disk accesses, for an M-tree?
- Given:
 - N (# of objects)
 - C (fanout of disk pages)
 - r (radius of range query - BIASED model)
- NOT ENOUGH - what else do we need?

15-826 Copyright: C. Faloutsos (2010) 9

CMU SCS

Metric trees - analysis

- A: something about the distribution



15-826 Copyright: C. Faloutsos (2010) 10

CMU SCS

Metric trees - analysis

- A: something about the distribution

[Ciaccia, Patella, Zezula, PODS98]: assumed that the distance distribution is the same, for every object:

Paolo Ciaccia Marco Patella

15-826 Copyright: C. Faloutsos (2010) 11

CMU SCS

Metric trees - analysis

- A: something about the distribution

[Ciaccia+, PODS98]: assumed that the distance distribution is the same, for every object:

$F1(d) = \text{Prob}(\text{an object is within } d \text{ from object \#1})$
 $= F2(d) = \dots = F(d)$

15-826 Copyright: C. Faloutsos (2010) 12

CMU SCS

Metric trees - analysis

- A: something about the distribution
- Given our ‘fractal’ tools, we could try them - which one?

15-826 Copyright: C. Faloutsos (2010) 13

CMU SCS

Metric trees - analysis

- A: something about the distribution
- Given our ‘fractal’ tools, we could try them - which one?
- A: Correlation integral [Traina+, ICDE2000]

15-826 Copyright: C. Faloutsos (2010) 14

CMU SCS

Metric trees - analysis

English dictionary

log(#pairs)

EnglishWords dataset

log(d)

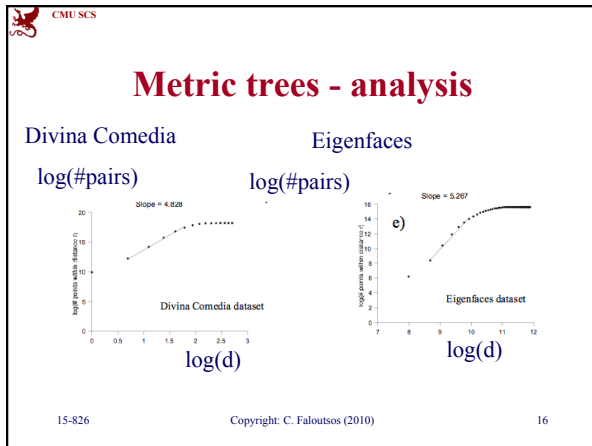
Portuguese dictionary

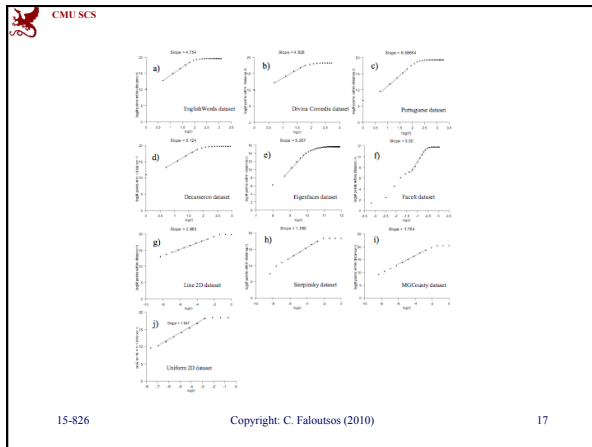
log(#pairs)

Portuguese dataset

log(d)

15-826 Copyright: C. Faloutsos (2010) 15





Metric trees - analysis

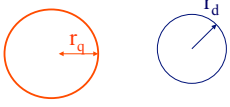
	Data Set	N (# Objects)	Dimension	Distance Function	Distance Exponent, D
Real Metric datasets	English	25,143	NA	L_{Edit}	4.753
	Divina Commedia	12,701	NA	L_{Edit}	4.827
	Decamerone	18,719	NA	L_{Edit}	5.124
	Portuguese	21,473	NA	L_{Edit}	6.686
	Faceit	1,056	NA	Not divulged	6.821
Real vector datasets	MGCounty	15,559	2	L_2	1.752
	Eigenfaces	11,900	16	L_2	5.267
Synthetic datasets	Sierpinski	9,841	2	L_2	1.584
	2D Line	20,000	2	L_2	0.989
	Uniform 2D	10,000	2	L_2	1.947

15-826 Copyright: C. Faloutsos (2010) 18

CMU SCS

Metric trees - analysis

- So, what is the # of disk accesses, for a node of radius r_d , on a query of radius r_q ?

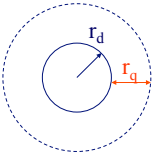


15-826 Copyright: C. Faloutsos (2010) 19

CMU SCS

Metric trees - analysis

- So, what is the # of disk accesses, for a node of radius r_d , on a query of radius r_q ?
- A: $\sim (r_d+r_q) \dots$

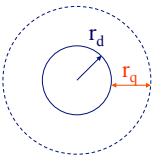


15-826 Copyright: C. Faloutsos (2010) 20

CMU SCS

Metric trees - analysis

- So, what is the # of disk accesses, for a node of radius r_d , on a query of radius r_q ?
- A: $\sim (r_d+r_q)^D$



15-826 Copyright: C. Faloutsos (2010) 21

CMU SCS

Accuracy of selectivity formulas

$\log(\#d.a.)$

Legend: — DA Method (a) — DA Method (b) — DA Method (c) — DA Method (d) — DA Method (e) — DA Method (f)

15-826 Copyright: C. Faloutsos (2010) $\log(rq)$ 22

CMU SCS

Fast estimation of D

- Normally, D takes $O(N^2)$ time
- Anything faster? suppose we have already built an M-tree

15-826 Copyright: C. Faloutsos (2010) 23

CMU SCS

Fast estimation of D

- Hint:

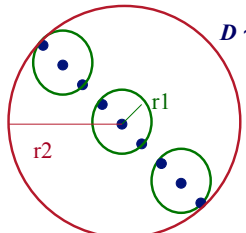
15-826 Copyright: C. Faloutsos (2010) 24

CMU SCS

Fast estimation of D

- Hint:

ratio of radii:
 $r1^D * C = r2^D$
 $D \sim \log(C) / \log(r2/r1)$



15-826 Copyright: C. Faloutsos (2010) 25

CMU SCS

Indexing - Detailed outline

- fractals
 - intro
 - applications
 - disk accesses for R-trees (range queries)
 - dimensionality reduction
 - selectivity in M-trees
 - ➔ dim. curse revisited
 - “fat fractals”
 - quad-tree analysis [Gaede+]


15-826 Copyright: C. Faloutsos (2010) 26

CMU SCS

Dim. curse revisited

- (Q: how serious is the dim. curse, e.g.:)
- Q: what is the search effort for k-nn?
 - given N points, in E dimensions, in an R-tree, with k-nn queries (‘biased’ model)

[Pagel, Korn + ICDE 2000]



15-826 Copyright: C. Faloutsos (2010) 27

CMU SCS

(Overview of proofs)

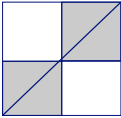
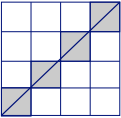
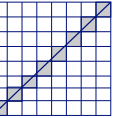
- assume that your points are uniformly distributed in a d -dimensional manifold (= hyper-plane)
- derive the formulas
- substitute d for the fractal dimension

15-826 Copyright: C. Faloutsos (2010) 28

CMU SCS proof

Reminder: Hausdorff Dimension (D_0)

- r = side length (each dimension)
- $B(r)$ = # boxes containing points $\propto r^{D_0}$

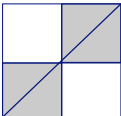
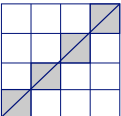
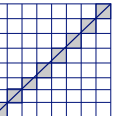
		
$r = 1/2 \quad B = 2$	$r = 1/4 \quad B = 4$	$r = 1/8 \quad B = 8$
$\log r = -1$ $\log B = 1$	$\log r = -2$ $\log B = 2$	$\log r = -3$ $\log B = 3$

15-826 Copyright: C. Faloutsos (2010) 29

CMU SCS proof

Reminder: Correlation Dimension (D_2)

- $S(r) = \sum p_i^2$ (squared % pts in box) $\propto r^{D_2}$
 \propto #pairs(within $\leq r$)

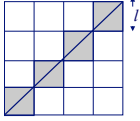
		
$r = 1/2 \quad S = 1/2$	$r = 1/4 \quad S = 1/4$	$r = 1/8 \quad S = 1/8$
$\log r = -1$ $\log S = -1$	$\log r = -2$ $\log S = -2$	$\log r = -3$ $\log S = -3$

15-826 Copyright: C. Faloutsos (2010) 30

CMU SCS proof

Observation #1

- How to determine avg MBR side l ?
 - $N = \#pts, C = \text{MBR capacity}$



Hausdorff dimension: $B(r) \propto r^{D_0}$

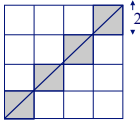
$B(l) = N/C = l^{-D_0} \Rightarrow l = (N/C)^{-1/D_0}$

15-826 Copyright: C. Faloutsos (2010) 31

CMU SCS proof

Observation #2

- k -NN query $\rightarrow \epsilon$ -range query
 - For k pts, what radius ϵ do we expect?



Correlation dimension: $S(r) \propto r^{D_2}$

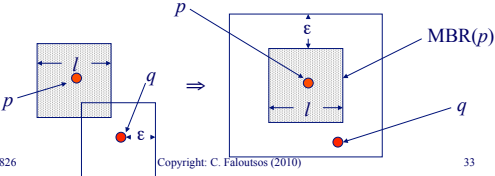
$S(\epsilon) = \frac{k}{N-1} = (2\epsilon)^{D_2}$

15-826 Copyright: C. Faloutsos (2010) 32

CMU SCS proof

Observation #3

- Estimate avg # query-sensitive anchors:
 - How many **expected** q will touch **avg** page?
 - Page touch: q stabs ϵ -dilated MBR(p)



15-826 Copyright: C. Faloutsos (2010) 33

CMU SCS

Asymptotic Formula

- k -NN page accesses as $N \rightarrow \infty$
 - C = page capacity
 - D = fractal dimension ($=D_0 \sim D_2$)

$$P_{all}^{L\infty}(k) \approx \sum_{j=0}^h \left\{ \frac{1}{C^{h-j}} + \left[1 + \left(\frac{k}{C^{h-j}} \right)^{1/D} \right]^D \right\}$$

15-826 Copyright: C. Faloutsos (2010) 34

CMU SCS

Asymptotic Formula

$$P_{all}^{L\infty}(k) \approx \sum_{j=0}^h \left\{ \frac{1}{C^{h-j}} + \left[1 + \left(\frac{k}{C^{h-j}} \right)^{1/D} \right]^D \right\}$$

- NO mention of the embedding dimensionality!!
- Still have dim. curse, but on f.d. D

15-826 Copyright: C. Faloutsos (2010) 35

CMU SCS

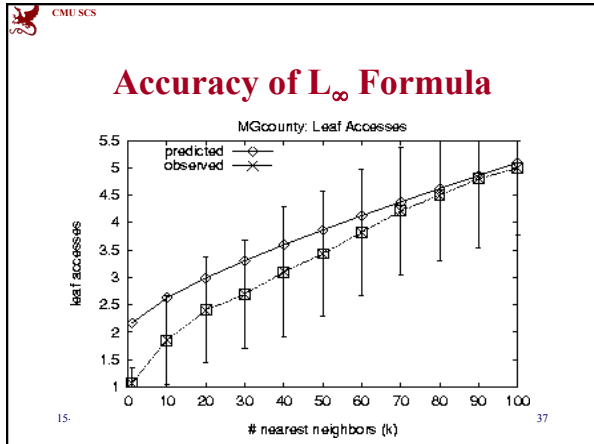
Synthetic Data

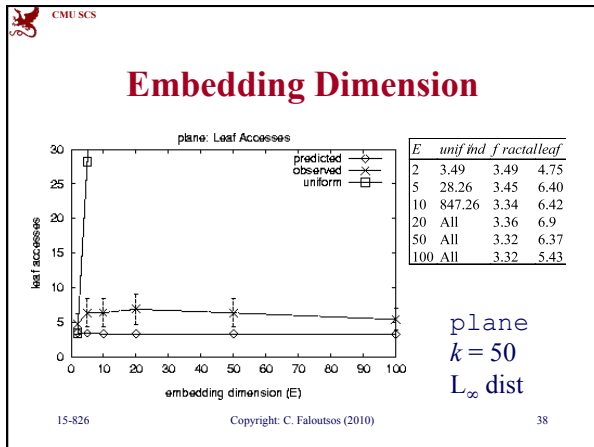
- plane
 - $D_0 = D_2 = 2$
 - embedded in E -space
 - $N = 100K$
- manifold
 - $E = 8$
 - $D_0 = D_2$ varies from 1-6
 - line, plane, etc. (in 8-d)

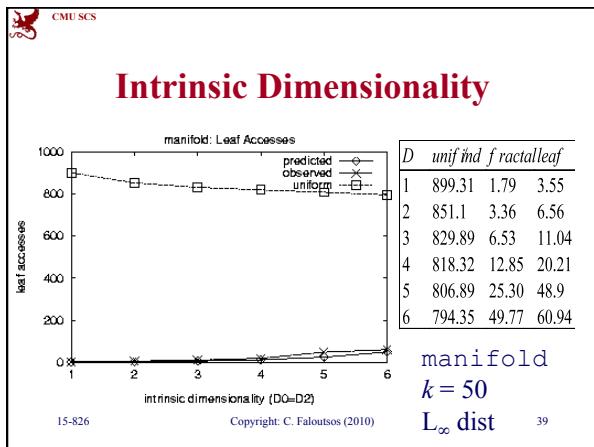
plane in 3-space ($E=3, D_0=D_2=2$)

line in 3-space ($E=3, D_0=D_2=1$)

15-826 Copyright: C. Faloutsos







CMU SCS

Non-Euclidean Data Set

<i>E</i>	<i>unif</i>	<i>ind</i>	<i>f ractal</i>	<i>leaf</i>
2	3.49	2.53	4.72±1.81	
10	847.26	2.53	6.42±2.11	
20	all	2.53	7.76±4.12	
50	all	2.53	6.15±2.82	
100	all	2.53	5.64±2.32	

15-826 sierpinski, $k = 50$, L_∞ dist 40

CMU SCS

Conclusions

- Worst-case theory is **over-pessimistic**
- High dimensional data can exhibit good performance if **correlated, non-uniform**
- Many real data sets are **self-similar**
- Determinant is **intrinsic** dimensionality
 - multiple fractal dimensions (D_0 and D_2)
 - indication of how far one can go

15-826 Copyright: C. Faloutsos (2010) 41

CMU SCS

References

- Ciaccia, P., M. Patella, et al. (1998). *A Cost Model for Similarity Queries in Metric Spaces*. PODS.
- Pagel, B.-U., F. Korn, et al. (2000). *Deflating the Dimensionality Curse Using Multiple Fractal Dimensions*. ICDE, San Diego, CA.
- Traina, C., A. J. M. Traina, et al. (2000). *Distance Exponent: A New Concept for Selectivity Estimation in Metric Trees*. ICDE, San Diego, CA.

15-826 Copyright: C. Faloutsos (2010) 42
