


CMU SCS

15-826: Multimedia Databases and Data Mining

Lecture #10: Fractals: M-trees and dim.
curse (case studies – Part II)
C. Faloutsos




CMU SCS

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 (2013) 2




CMU SCS

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 (2013) 3




CMU SCS

Outline

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

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

15-826 Copyright: C. Faloutsos (2013) 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 (2013) 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 (2013) 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 (2013) 7

CMU SCS

Metric trees - analysis

Optional

- 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 (2013) 8

CMU SCS

Metric trees - analysis

Optional

- 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 (2013) 9

CMU SCS

Metric trees - analysis

Optional

- A: something about the distribution

15-826 Copyright: C. Faloutsos (2013) 10



CMU SCS

Metric trees - analysis

Optional

- 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 (2013) 11

CMU SCS

Metric trees - analysis

Optional


- 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 (2013) 12

CMU SCS

Optional


Metric trees - analysis

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

15-826

Copyright: C. Faloutsos (2013)

13

CMU SCS

Optional


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 (2013)

14

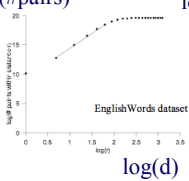
CMU SCS

Optional

Metric trees - analysis

English dictionary

log(#pairs)

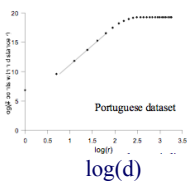


EnglishWords dataset

log(d)

Portuguese dictionary

log(#pairs)




Portuguese dataset

log(d)

15-826

Copyright: C. Faloutsos (2013)

15



CMU SCS

Optional

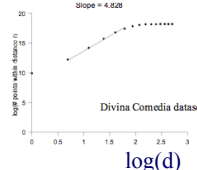
Metric trees - analysis

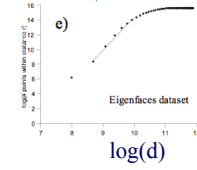
Divina Comedia

Eigenfaces

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

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





Divina Comedia dataset

Eigenfaces dataset


$\log(d)$

$\log(d)$

15-826

Copyright: C. Faloutsos (2013)

16



CMU SCS

a)

b)

c)

English Words dataset

Divina Comedia dataset

Portuguese dataset

d)

e)

f)

Decamerone dataset

Eigenfaces dataset

Facch dataset

g)

h)

i)

Line 2D dataset

Sierpinsky dataset

MdCounty dataset


j)

Uniform 2D dataset

15-826

Copyright: C. Faloutsos (2013)

17



CMU SCS

Optional

Metric trees - analysis

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

15-826

Copyright: C. Faloutsos (2013)

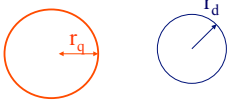
18

CMU SCS

Metric trees - analysis

Optional

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



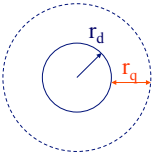
15-826 Copyright: C. Faloutsos (2013) 19

CMU SCS

Metric trees - analysis

Optional

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



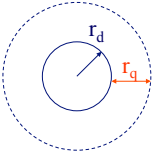
15-826 Copyright: C. Faloutsos (2013) 20

CMU SCS

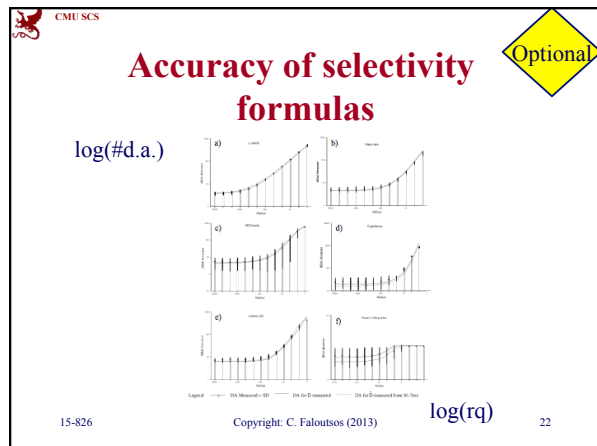
Metric trees - analysis

Optional

- 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 (2013) 21



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 (2013) 23

CMU SCS

Fast estimation of D

- Hint:

15-826 Copyright: C. Faloutsos (2013) 24

CMU SCS

Fast estimation of D

Optional

- Hint:

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

15-826 Copyright: C. Faloutsos (2013) 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 (2013) 26

CMU SCS

Dimensionality ‘curse’

- Q: What is the problem in high-d?

15-826 Copyright: C. Faloutsos (2013) #27




CMU SCS

Dimensionality ‘curse’

- Q: What is the problem in high-d?
- A: indices do not seem to help, for many queries (eg., k-nn)
 - in high-d (& uniform distributions), most points are equidistant -> k-nn retrieves too many near-neighbors
 - [Yao & Yao, '85]: search effort $\sim O(N^{(1-1/d)})$

15-826 Copyright: C. Faloutsos (2013) #28




CMU SCS

Dimensionality ‘curse’

- (counter-intuitive, for db mentality)
- Q: What to do, then?

15-826 Copyright: C. Faloutsos (2013) #29




CMU SCS

Dimensionality ‘curse’

- A1: switch to seq. scanning
- A2: dim. reduction
- A3: consider the ‘intrinsic’/fractal dimensionality
- A4: find approximate nn

15-826 Copyright: C. Faloutsos (2013) #30




CMU SCS

Dimensionality ‘curse’

- A1: switch to seq. scanning
 - X-trees [Kriegel+, VLDB 96]
 - VA-files [Schek+, VLDB 98], ‘test of time’ award

15-826 Copyright: C. Faloutsos (2013) #31




CMU SCS

Dimensionality ‘curse’

- A1: switch to seq. scanning
- ➔ • A2: dim. reduction
- A3: consider the ‘intrinsic’/fractal dimensionality
- A4: find approximate nn

15-826 Copyright: C. Faloutsos (2013) #32



CMU SCS

Dim. reduction

a.k.a. feature selection/extraction:

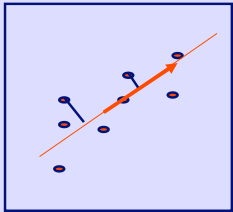
- SVD (optimal, to preserve Euclidean distances)
- random projections
- using the fractal dimension [Traina+ SBBD2000]

15-826 Copyright: C. Faloutsos (2013) #33

CMU SCS

Singular Value Decomposition (SVD)

- SVD (\sim LSI \sim KL \sim PCA \sim spectral analysis...)



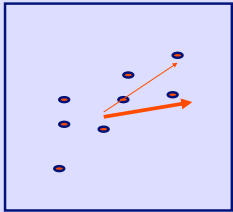
LSI: S. Dumais; M. Berry
KL: eg, Duda+Hart
PCA: eg., Jolliffe
MANY more details: soon

15-826 Copyright: C. Faloutsos (2013) #34

CMU SCS

Random projections

- random projections(Johnson-Lindenstrauss thm [Papadimitriou+ pods98])



15-826 Copyright: C. Faloutsos (2013) #35

CMU SCS

Random projections

- pick 'enough' random directions (will be \sim orthogonal, in high-d!!)
- distances are preserved probabilistically, within epsilon
- (also, use as a pre-processing step for SVD [Papadimitriou+ PODS98])

15-826 Copyright: C. Faloutsos (2013) #36

CMU SCS

Dim. reduction - w/ fractals

- Main idea: drop those attributes that don't affect the intrinsic ('fractal') dimensionality [Traina+, SBBD 2000]

15-826 Copyright: C. Faloutsos (2013) #37

CMU SCS

Dim. reduction - w/ fractals

global FD=1

(a) Quarter-circle (b) Line (c) Spike

15-826 Copyright: C. Faloutsos (2013) #38

CMU SCS

Dimensionality 'curse'

- A1: switch to seq. scanning
- A2: dim. reduction
- ➔ A3: consider the 'intrinsic'/fractal dimensionality
- A4: find **approximate** nn

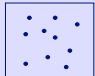
15-826 Copyright: C. Faloutsos (2013) #39


CMU SCS

Intrinsic dimensionality

- before we give up, compute the intrinsic dim.:
- the lower, the better... [Pagel+, ICDE 2000]
- more details: under 'fractals'

intr. $d = 2$





intr. $d = 1$

15-826 Copyright: C. Faloutsos (2013) #40

CMU SCS

Dimensionality 'curse'

- A1: switch to seq. scanning
- A2: dim. reduction
- A3: consider the 'intrinsic'/fractal dimensionality
- ➔ A4: find approximate nn

15-826 Copyright: C. Faloutsos (2013) #41

CMU SCS

Approximate nn

- [Arya + Mount, SODA93], [Patella+ ICDE 2000]
- Idea: find k neighbors, such that the distance of the k -th one is guaranteed to be within epsilon of the actual.

15-826 Copyright: C. Faloutsos (2013) #42

CMU SCS

Dimensionality ‘curse’

- A1: switch to seq. scanning
- A2: dim. reduction
- ➔ • A3: consider the ‘intrinsic’/fractal dimensionality
- A4: find approximate nn



15-826 Copyright: C. Faloutsos (2013) #43

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 (2013) 44

CMU SCS

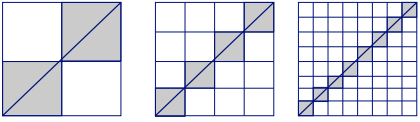
(Overview of proofs)

- 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 (2013) 45


Reminder: Hausdorff Dimension (D_0) 

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

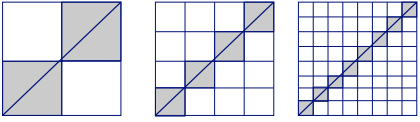


$r = 1/2 \quad B = 2$ $r = 1/4 \quad B = 4$ $r = 1/8 \quad B = 8$
 $\log r = -1$ $\log r = -2$ $\log r = -3$
 $\log B = 1$ $\log B = 2$ $\log B = 3$

15-826 Copyright: C. Faloutsos (2013) 46


Reminder: Correlation Dimension (D_2) 

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

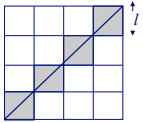


$r = 1/2 \quad S = 1/2$ $r = 1/4 \quad S = 1/4$ $r = 1/8 \quad S = 1/8$
 $\log r = -1$ $\log r = -2$ $\log r = -3$
 $\log S = -1$ $\log S = -2$ $\log S = -3$

15-826 Copyright: C. Faloutsos (2013) 47

Observation #1 

- How to determine avg MBR side l ?
 N = #pts, C = MBR capacity



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

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

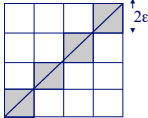
15-826 Copyright: C. Faloutsos (2013) 48

CMU SCS

Observation #2

• k -NN query $\rightarrow \epsilon$ -range query

- For k pts, what radius ϵ do we expect?



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

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

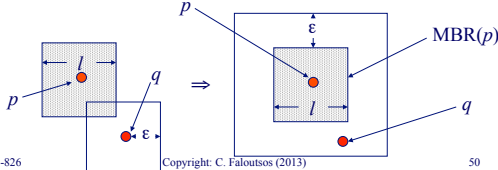
15-826 Copyright: C. Faloutsos (2013) 49

CMU SCS

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 (2013) 50

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 (2013) 51


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

15-826Copyright: C. Faloutsos (2013)52

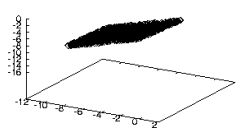
CMU SCS

Synthetic Data

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

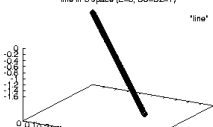
"plane" ◇

- 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)




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

"line" ◇

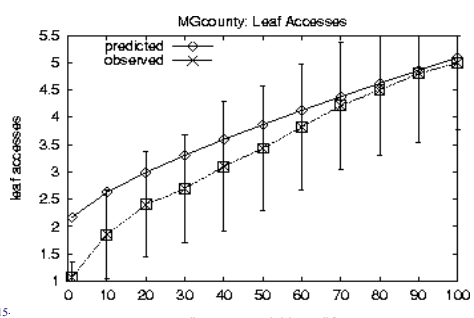


15-826Copyright: C. Faloutsos

CMU SCS

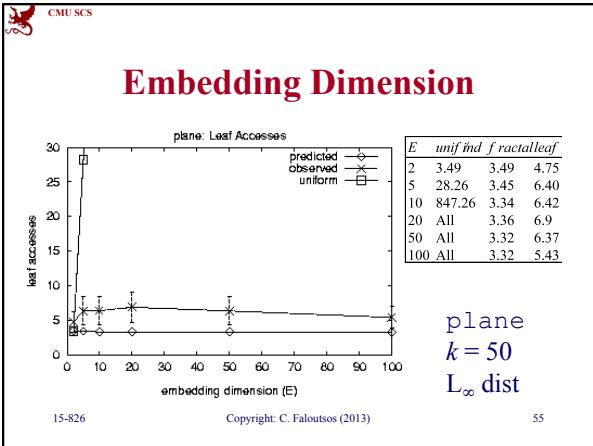
Accuracy of L_{∞} Formula

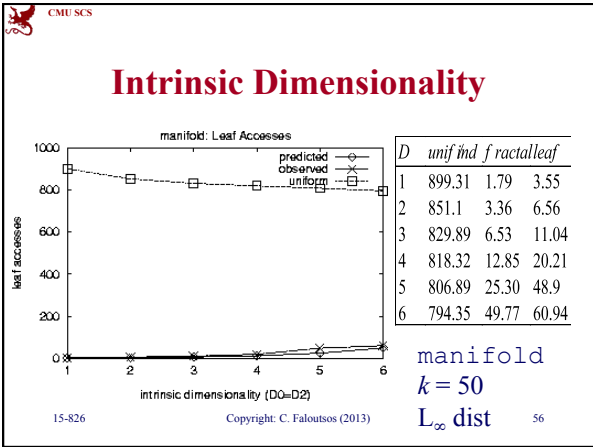
MGcounty: Leaf Accesses

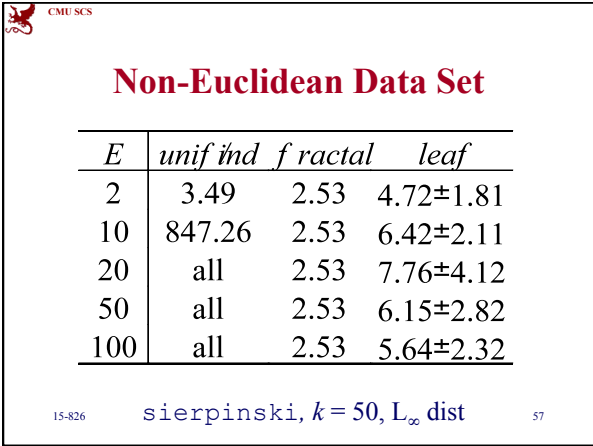



15:54

15:54










CMU SCS

Conclusions

- Dimensionality ‘curse’:
 - for high-d, indices slow down to $\sim O(N)$
- If the **intrinsic** dim. is low, there is hope
- otherwise, do seq. scan, or sacrifice accuracy (approximate nn)

15-826 Copyright: C. Faloutsos (2013) #58




CMU SCS

Conclusions – cont’d

- 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 (2013) 59




CMU SCS

References

- Sunil Arya, David M. Mount: *Approximate Nearest Neighbor Queries in Fixed Dimensions*. SODA 1993: 271-280
ANN library:
<http://www.cs.umd.edu/~mount/ANN/>

15-826 Copyright: C. Faloutsos (2013) #60

CMU SCS


References

- Berchtold, S., D. A. Keim, et al. (1996). The X-tree : An Index Structure for High-Dimensional Data. VLDB, Mumbai (Bombay), India.
- Ciaccia, P., M. Patella, et al. (1998). *A Cost Model for Similarity Queries in Metric Spaces*. PODS.

15-826

Copyright: C. Faloutsos (2013)

#61

CMU SCS


References cnt'd

- Nievergelt, J., H. Hinterberger, et al. (March 1984). "The Grid File: An Adaptable, Symmetric Multikey File Structure." ACM TODS 9(1): 38-71.
- ➡ • Pagel, B.-U., F. Korn, et al. (2000). *Deflating the Dimensionality Curse Using Multiple Fractal Dimensions*. ICDE, San Diego, CA.
- Papadimitriou, C. H., P. Raghavan, et al. (1998). Latent Semantic Indexing: A Probabilistic Analysis. PODS, Seattle, WA.

15-826

Copyright: C. Faloutsos (2013)

#62

CMU SCS


References cnt'd

- ➡ • Traina, C., A. J. M. Traina, et al. (2000). *Distance Exponent: A New Concept for Selectivity Estimation in Metric Trees*. ICDE, San Diego, CA.
- Weber, R., H.-J. Schek, et al. (1998). A Quantitative Analysis and Performance Study for Similarity-Search Methods in high-dimensional spaces. VLDB, New York, NY.

15-826

Copyright: C. Faloutsos (2013)

#63

CMU/SCS

References cnt'd

➡ • Yao, A. C. and F. F. Yao (May 6-8, 1985). A General Approach to d-Dimensional Geometric Queries. Proc. of the 17th Annual ACM Symposium on Theory of Computing (STOC), Providence, RI.

15-826

Copyright: C. Faloutsos (2013)

#64
