

CarnegieMellon

# 15-826: Multimedia Databases and Data Mining

Lecture #10: Fractals - case studies

*C. Faloutsos*

CarnegieMellon

## Must-read Material - I

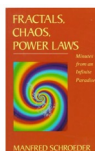
- Christos Faloutsos and Ibrahim Kamel, [\*Beyond Uniformity and Independence: Analysis of R-trees Using the Concept of Fractal Dimension\*](#), Proc. ACM SIGACT-SIGMOD-SIGART PODS, May 1994, pp. 4-13, Minneapolis, MN.

## Must-read Material - II

- Bernd-Uwe Pagel, Flip Korn and Christos Faloutsos, [\*Deflating the Dimensionality Curse using Multiple Fractal Dimensions\*](#), ICDE 2000, San Diego, CA, Feb. 2000.

## Optional Material

Optional, but **very** useful: Manfred Schroeder *Fractals, Chaos, Power Laws: Minutes from an Infinite Paradise* W.H. Freeman and Company, 1991 (on reserve in the WeH library)




## Reminder

- Code at [www.cs.cmu.edu/~christos/SRC/fdnq\\_h.zip](http://www.cs.cmu.edu/~christos/SRC/fdnq_h.zip)

Also, in 'R'  
> library(fdim);

## 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

(c) 2019 C. Faloutsos

7

## Indexing - Detailed outline

- fractals
  - intro
  - applications
    - disk accesses for R-trees (range queries)
    - dim. curse revisited
    - nearest neighbors estimation




15-826

(c) 2019 C. Faloutsos

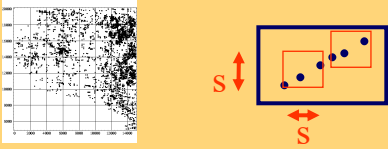
8

CarnegieMellon




## Problem:

- Selectivity of a range query in R-trees?



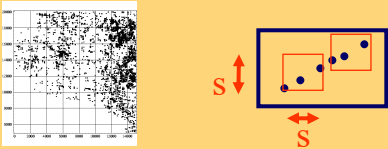
15-826 (c) 2019 C. Faloutsos 9

CarnegieMellon



## Solution:

- Selectivity of a range query in R-trees?
- Depends on *\*fractal\** dimension

$$s = (C/N)^{1/D_0}$$


15-826 (c) 2019 C. Faloutsos 10

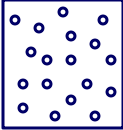
CarnegieMellon

## Case study#1: R-tree performance

Problem

- Given
  - N points in E-dim space
- Estimate # disk accesses for a range query  
( $q_1 \times \dots \times q_E$ )

(assume: ‘good’ R-tree, with tight, cube-like MBRs)



15-826 (c) 2019 C. Faloutsos 11

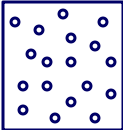
CarnegieMellon

## Case study#1: R-tree performance

Problem

- Given
  - N points in E-dim space
- Estimate # disk accesses for a range query  
( $q_1 \times \dots \times q_E$ )

(assume: ‘good’ R-tree, with tight, cube-like MBRs)  
Typically, in DB Q-opt?



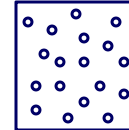
15-826 (c) 2019 C. Faloutsos 12

CarnegieMellon

## Case study#1: R-tree performance

### Problem

- Given
  - N points in E-dim space
- Estimate # disk accesses for a range query  
( $q_1 \times \dots \times q_E$ )



(assume: 'good' R-tree, with tight, cube-like MBRs)  
Typically, in DB Q-opt: uniformity + independence

15-826

(c) 2019 C. Faloutsos

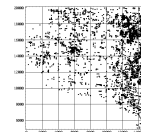
13

CarnegieMellon

## Case study#1: R-tree performance

### Problem

- Given
  - N points in E-dim space
  - – with fractal dimension D
- Estimate # disk accesses for a range query  
( $q_1 \times \dots \times q_E$ )



(assume: 'good' R-tree, with tight, cube-like MBRs)  
Typically, in DB Q-opt: uniformity + independence

15-826

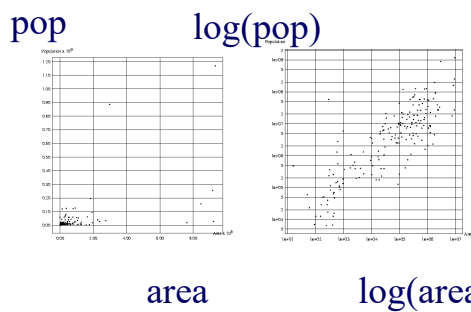
(c) 2019 C. Faloutsos

14

CarnegieMellon

## Examples: World's countries

- BUT: area vs population for  $\sim 200$  countries (1991 CIA fact-book).



15-826

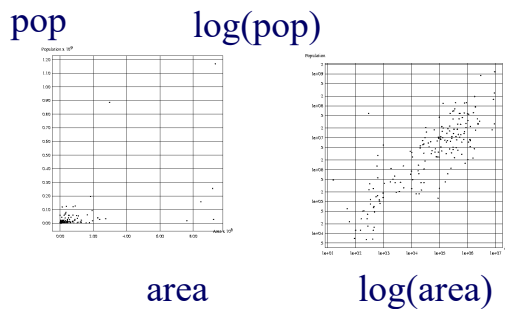
(c) 2019 C. Faloutsos

15

CarnegieMellon

## Examples: World's countries

- neither uniform, nor independent!

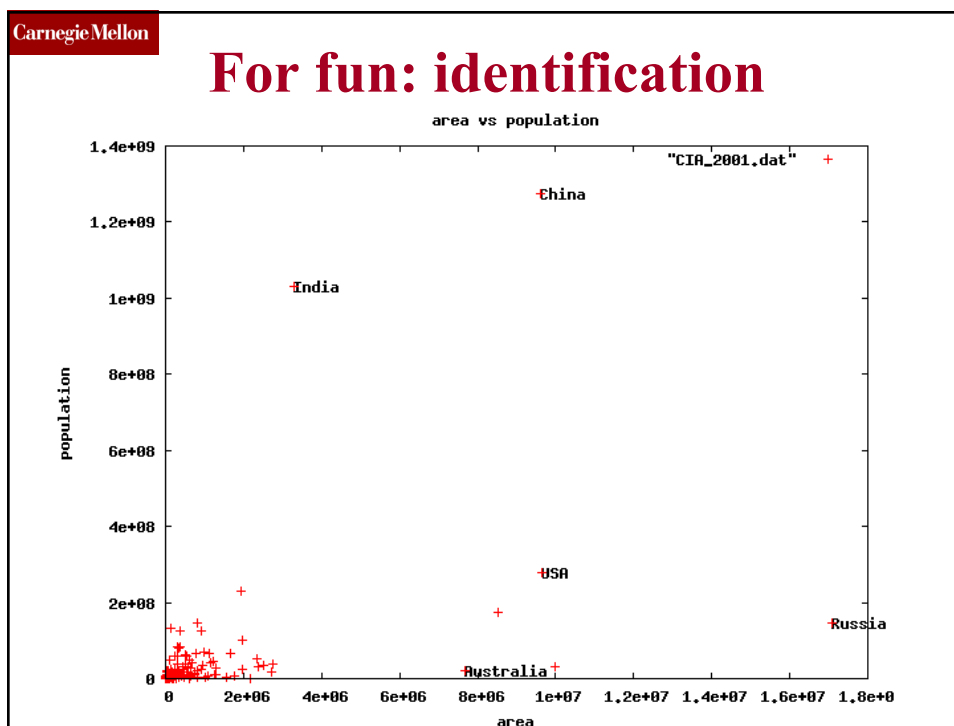
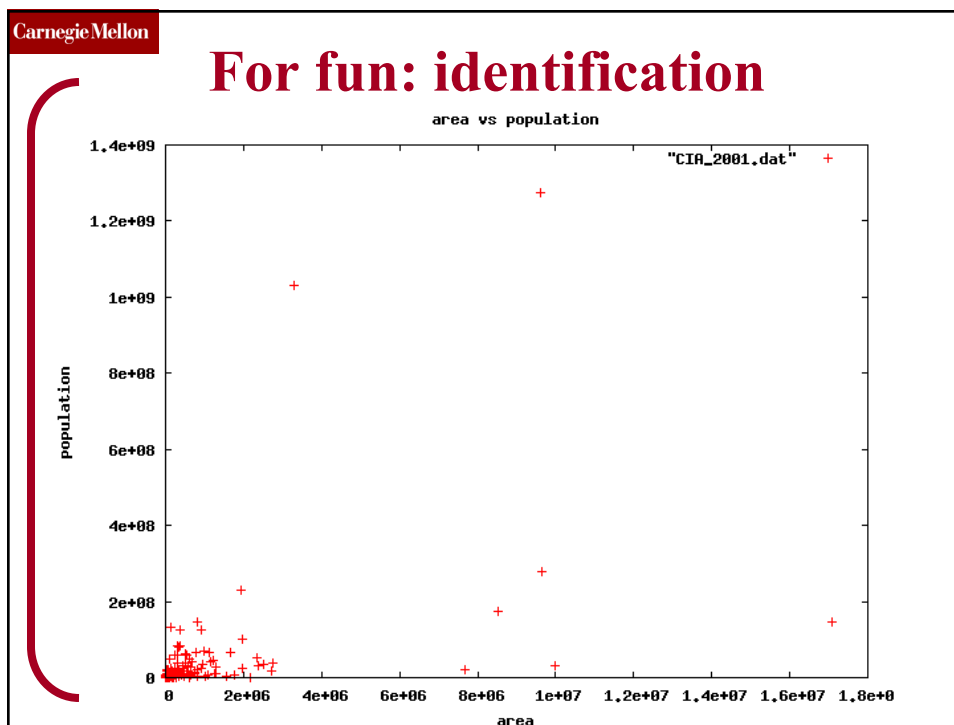


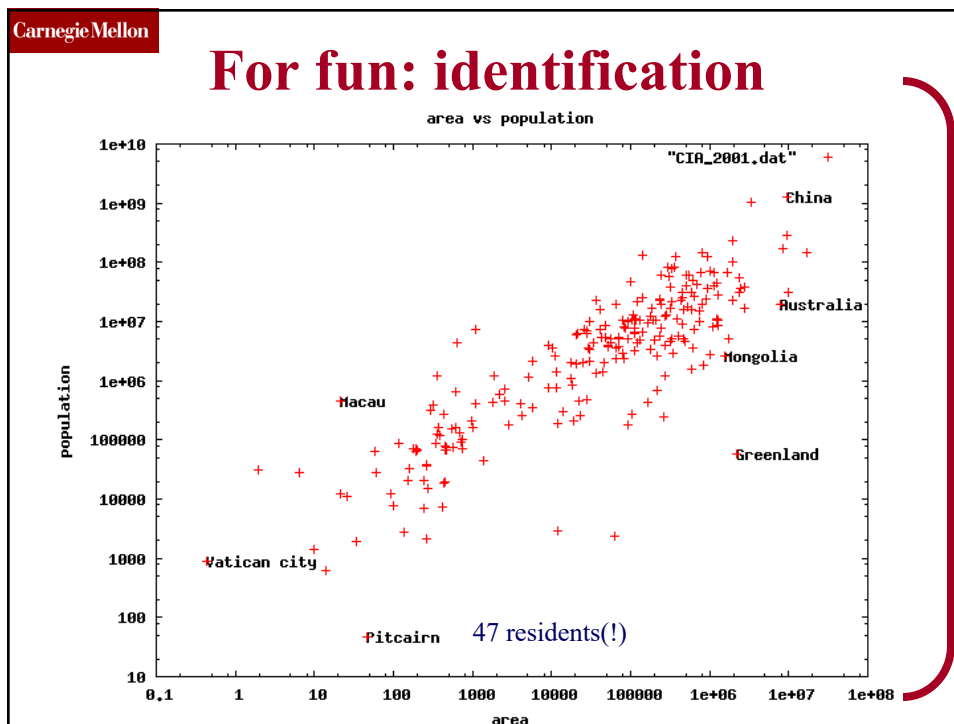
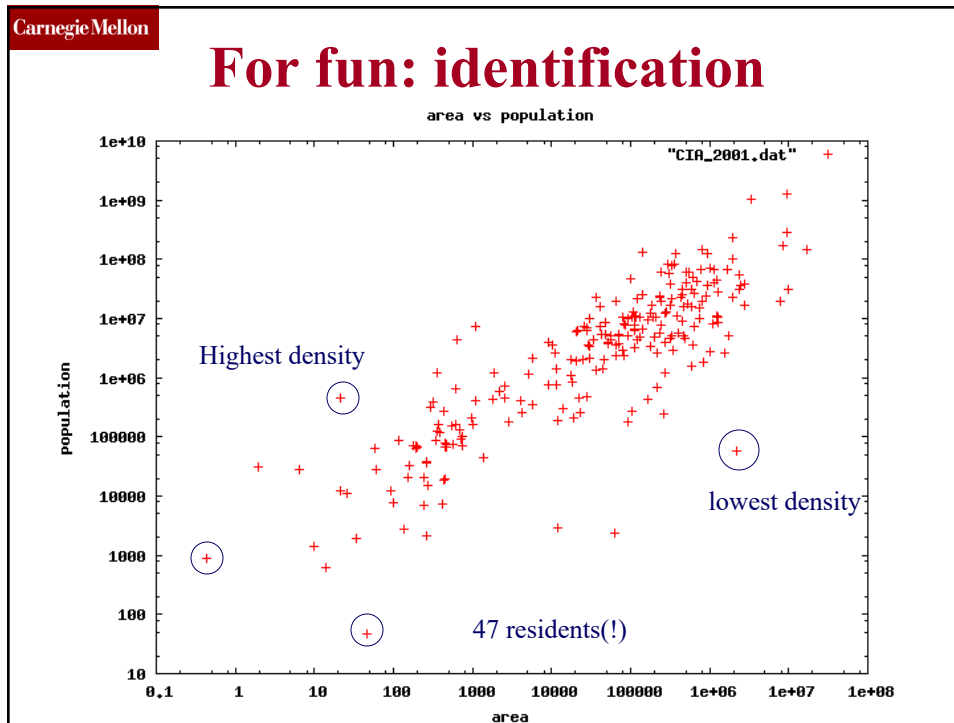
15-826

(c) 2019 C. Faloutsos

16



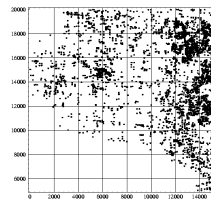




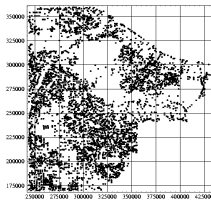
## Examples: TIGER files

- neither uniform, nor independent!

MG county



LB county



15-826

(c) 2019 C. Faloutsos

21

## How to proceed?

- recall the [Pagel+] formula, for range queries of size  $q1 \times q2$

$$\#DiskAccesses(q1, q2) = \sum (x_{i,1} + q1) * (x_{i,2} + q2)$$

15-826

(c) 2019 C. Faloutsos

22

CarnegieMellon Proof

## R-trees - performance analysis

- How many times will P1 be retrieved (unif. queries of size  $q_1 \times q_2$ )?

15-826 (c) 2019 C. Faloutsos #23

CarnegieMellon

## How to proceed?

- recall the [Pagel+] formula, for range queries of size  $q_1 \times q_2$

$$\#DiskAccesses(q_1, q_2) = \sum (x_{i,1} + q_1) * (x_{i,2} + q_2)$$

But:

formula needs to know the  $x_{i,j}$  sizes of MBRs!

15-826 (c) 2019 C. Faloutsos 24

CarnegieMellon

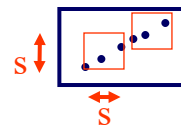
## How to proceed?

But:

formula needs to know the  $x_{i,j}$  sizes of MBRs!

Answer (jumping ahead):

$$s = (C/N)^{1/D_0}$$



15-826

(c) 2019 C. Faloutsos

25

CarnegieMellon

## How to proceed?

But:

formula needs to know the  $x_{i,j}$  sizes of MBRs!

Answer (jumping ahead):

$$s = (C/N)^{1/D_0}$$

← Hausdorff fd  
 ← # of data points  
 ← page capacity  
 ← side of (parent) MBR

15-826

(c) 2019 C. Faloutsos

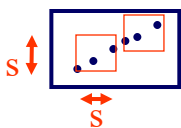
26

CarnegieMellon

### ‘smell’ tests:

- C ↗ s ↘
- N ↗ s ↘
- D0 ↗ s ↘

$$s = (C/N)^{1/D0}$$



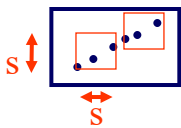
15-826 (c) 2019 C. Faloutsos 27

CarnegieMellon

### ‘smell’ tests:

- C ↗ s ↗
- N ↗ s ↘
- D0 ↗ s ↗

$$s = (C/N)^{1/D0}$$



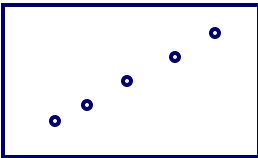
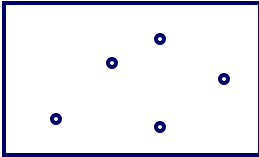
15-826 (c) 2019 C. Faloutsos 28

CarnegieMellon
PROOF

## R-trees - performance analysis

I.e: for range queries - how many disk accesses,  
if we just now that we have  
-  $N$  points in  $E$ -d space?

A: can not tell! need to know distribution

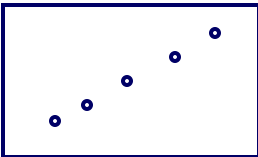
15-826
(c) 2019 C. Faloutsos
29

CarnegieMellon
PROOF

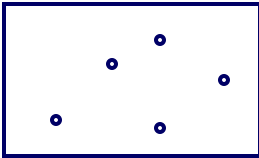
## R-trees - performance analysis

Q: OK - so we are told that the **Hausdorff** fractal  
dim. =  $D_0$  - Next step?  
(also know that there are at most  $C$  points per  
page)

$D_0=1$



$D_0=2$



15-826
(c) 2019 C. Faloutsos
30

CarnegieMellon PROOF

## R-trees - performance analysis

Assumption 1: square-like parents ( $s*s$ )  
 Assumption 2: fully packed (C points each)  
 Assumption 3: non-overlapping

$D_0=1$

$s_1=s_2=s$

$D_0=2$

15-826 (c) 2019 C. Faloutsos 31

CarnegieMellon PROOF

## R-trees - performance analysis

Assumption 1: square-like parents ( $s*s$ )  
 Assumption 2: fully packed (N/C non-empty)  
 Assumption 3: non-overlapping

$D_0=1$

$s_1=s_2=s$

15-826 (c) 2019 C. Faloutsos 32



CarnegieMellon

## R-trees - performance analysis

Hint: dfn of Hausdorff f.d.:



Felix Hausdorff (1868-1942)

15-826

(c) 2019 C. Faloutsos

33

CarnegieMellon

PROOF

### Reminder:

### Hausdorff or box-counting fd:

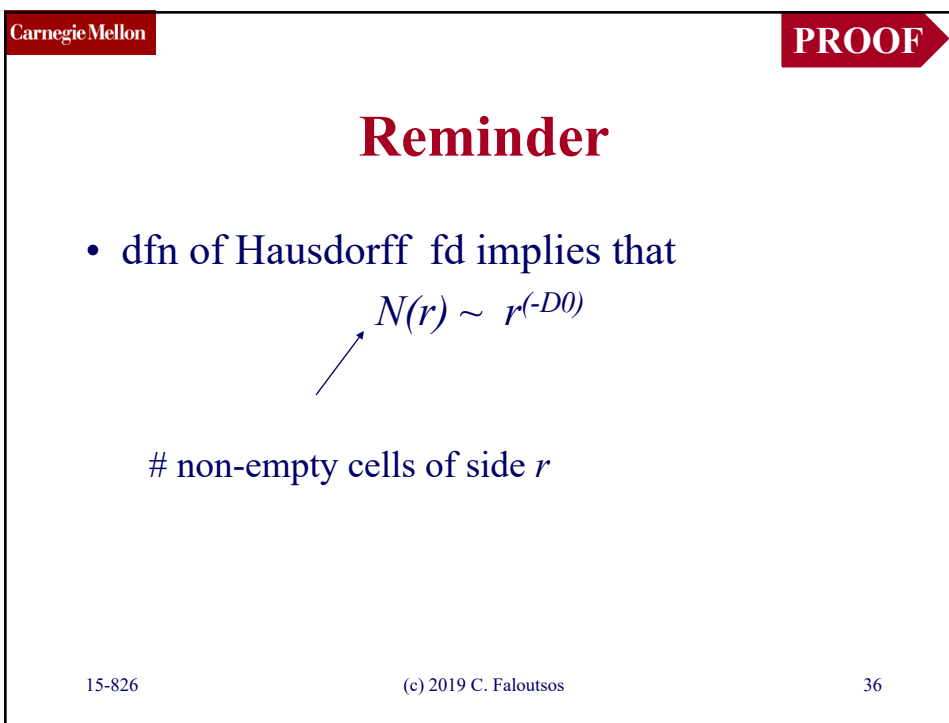
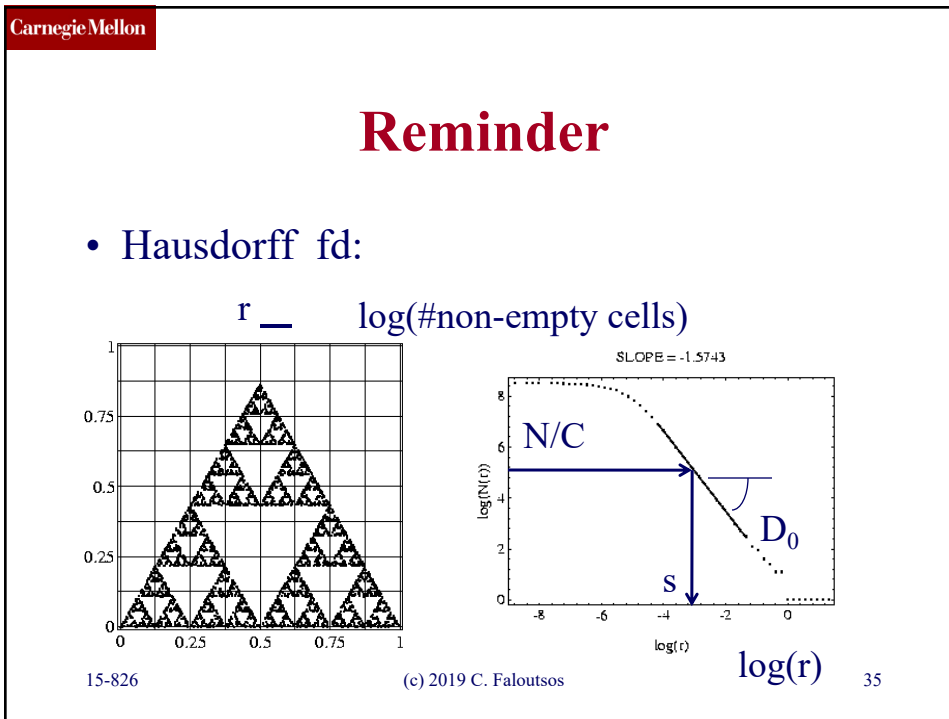
- Box counting plot:  $\text{Log}(N(r))$  vs  $\text{Log}(r)$
- $r$ : grid side
- $N(r)$ : count of non-empty cells
- (Hausdorff) fractal dimension  $D_0$ :

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

15-826

(c) 2019 C. Faloutsos

34

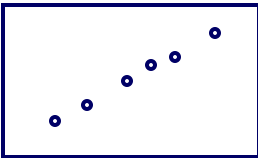


CarnegieMellon PROOF

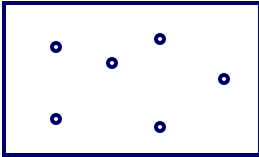
## R-trees - performance analysis

Q (rephrased): what is the side  $s_1, s_2, \dots$  of parent nodes, given  $N$  data points, packed by  $C$ , with f.d. =  $D0$

$D0=1$



$D0=2$



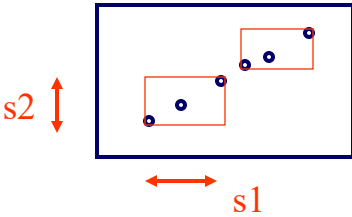
15-826 (c) 2019 C. Faloutsos 37

CarnegieMellon PROOF

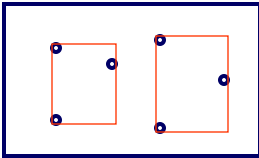
## R-trees - performance analysis

Q (rephrased): what is the side  $s_1, s_2, \dots$  of parent nodes, given  $N$  data points, packed by  $C$ , with f.d. =  $D0$

$D0=1$



$D0=2$



15-826 (c) 2019 C. Faloutsos 38

CarnegieMellon PROOF

## R-trees - performance analysis

Q (rephrased): what is the side  $s_1, s_2, \dots$  of parent nodes, given  $N$  data points, packed by  $C$ , with f.d. =  $D_0$

$D_0=1$

$s_1=s_2=s$

$D_0=2$

(c) 2019 C. Faloutsos 39

CarnegieMellon PROOF

## R-trees - performance analysis

A: (educated guess)

- $s=s_1=s_2$  (= ... ) - square-like MBRs
- $N/C$  non-empty cells =  $K * s^{(-D_0)}$

$D_0=1$

$D_0=2$

$\log(\#cells)$

$\log(s)$

(c) 2019 C. Faloutsos 40

CarnegieMellon

## R-trees - performance analysis

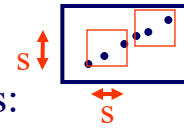
PROOF of derivations: in [PODS 94].

Finally, expected side  $s$  of parent MBRs:

$$s = (C/N)^{1/D0}$$

Q: sanity check: how does  $s$  change with  $D0$ ?

A:



15-826

(c) 2019 C. Faloutsos

41

CarnegieMellon



## R-trees - performance analysis

PROOF of derivations: in [Kamel+, PODS 94].

Finally, expected side  $s$  of parent MBRs:

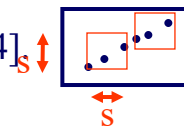
$$s = (C/N)^{1/D0}$$

Q: sanity check: how does  $s$  change with  $D0$ ?

A:  $s$  grows with  $D0$

Q: does it make sense?

Q: does it suffer from (intrinsic) dim. curse?



15-826

(c) 2019 C. Faloutsos

42

CarnegieMellon

PROOF

## R-trees - performance analysis

Q: Final-final formula (# disk accesses for range queries  $q1 \times q2 \times \dots$ ):

A:

15-826

(c) 2019 C. Faloutsos

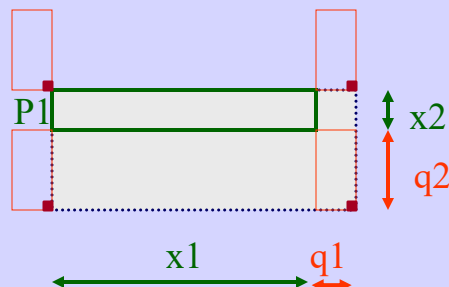
43

CarnegieMellon

Proof

## R-trees - performance analysis

- How many times will P1 be retrieved (unif. queries of size  $q1 \times q2$ )?



15-826

(c) 2019 C. Faloutsos

#44

CarnegieMellon

PROOF

## R-trees - performance analysis

Q: Final-final formula (# disk accesses for range queries  $q1 \times q2 \times \dots$ ):

A: # of parent-node accesses:

$$N/C * (s + q1) * (s + q2) * \dots * (s + qE)$$

A: # of grand-parent node accesses

15-826

(c) 2019 C. Faloutsos

45

CarnegieMellon

PROOF

## R-trees - performance analysis

Q: Final-final formula (# disk accesses for range queries  $q1 \times q2 \times \dots$ ):

A: # of parent-node accesses:

$$N/C * (s + q1) * (s + q2) * \dots * (s + qE)$$

A: # of grand-parent node accesses

$$N/(C^2) * (s' + q1) * (s' + q2) * \dots * (s' + qE)$$

$$s' = ??$$

15-826

(c) 2019 C. Faloutsos

46

## R-trees - performance analysis

Q: Final-final formula (# disk accesses for range queries  $q1 \times q2 \times \dots$ ):

A: # of parent-node accesses:

$$N/C * (s + q1) * (s + q2) * \dots * (s + qE)$$

A: # of grand-parent node accesses

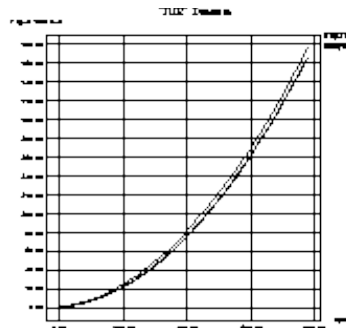
$$N/(C^2) * (s' + q1) * (s' + q2) * \dots * (s' + qE)$$

$$s' = (C^2/N)^{1/D0}$$

## R-trees - performance analysis

Results: IUE (x-y star coordinates)

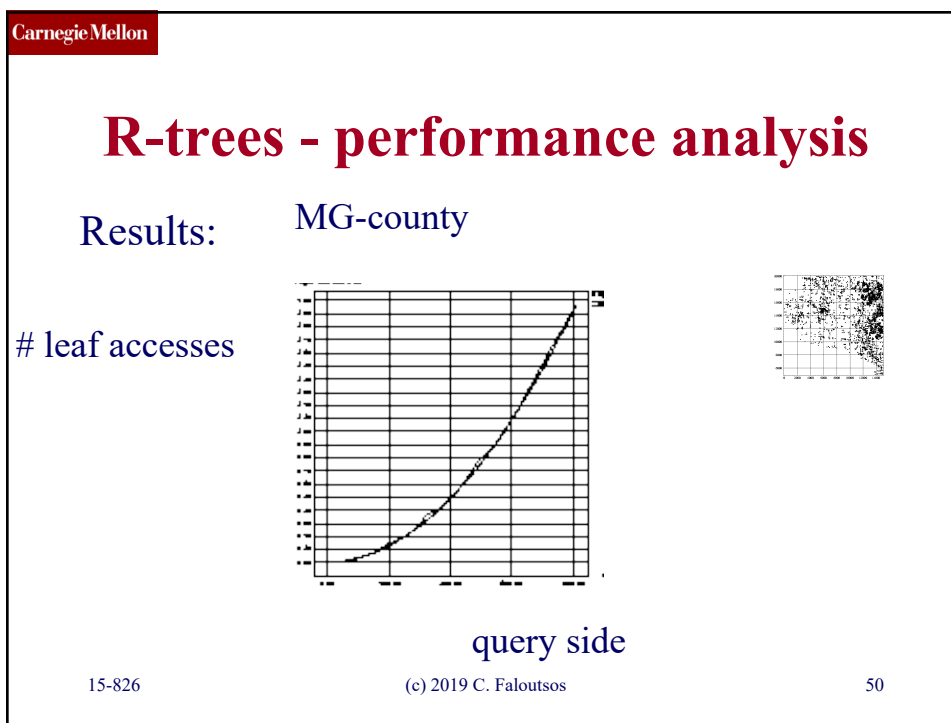
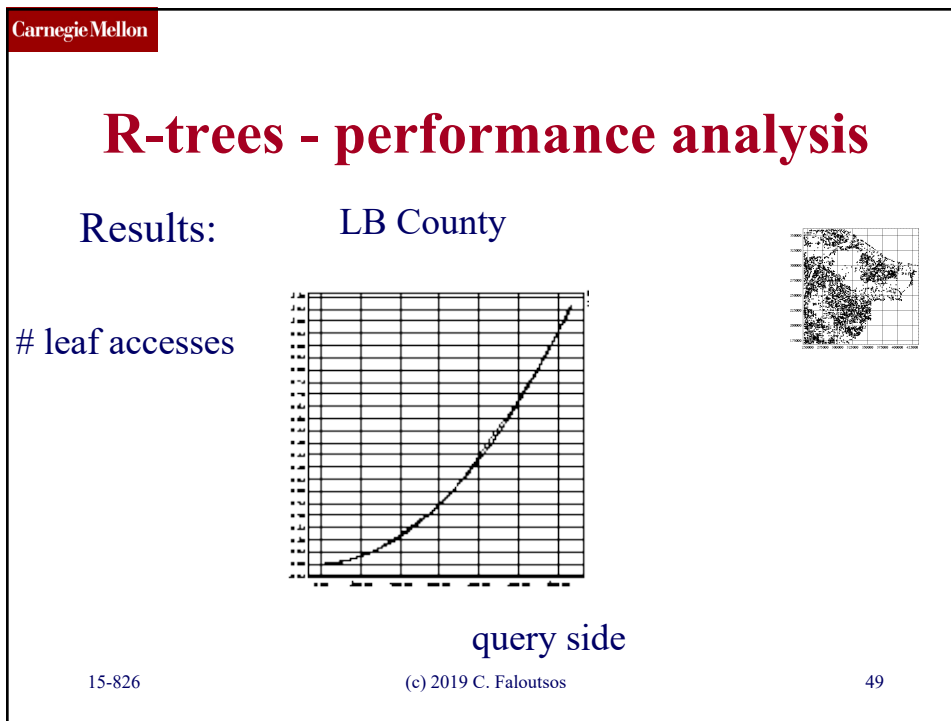
# leaf accesses

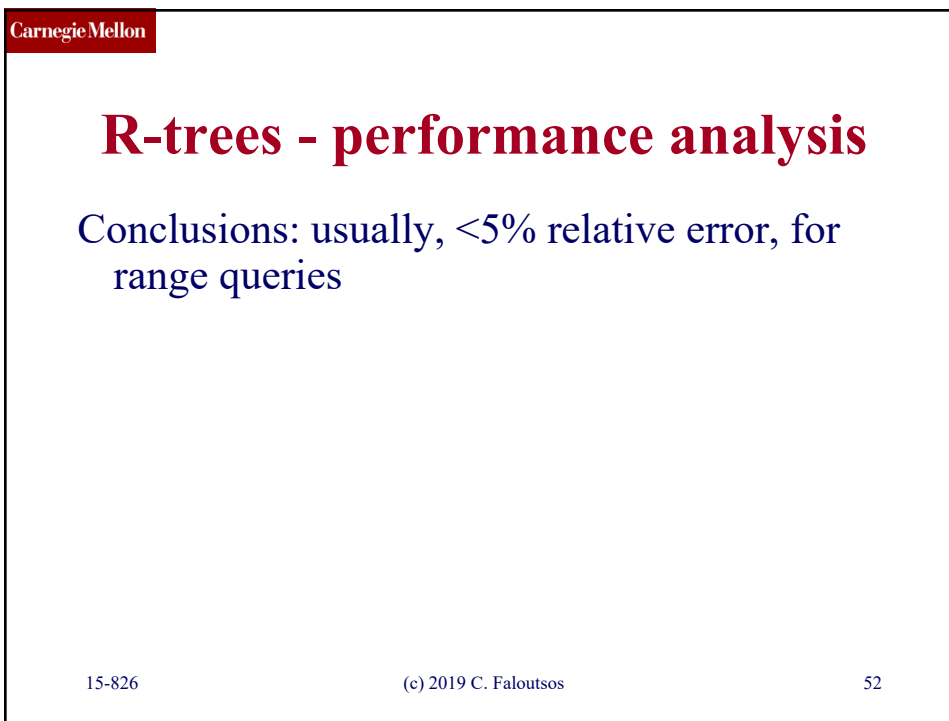
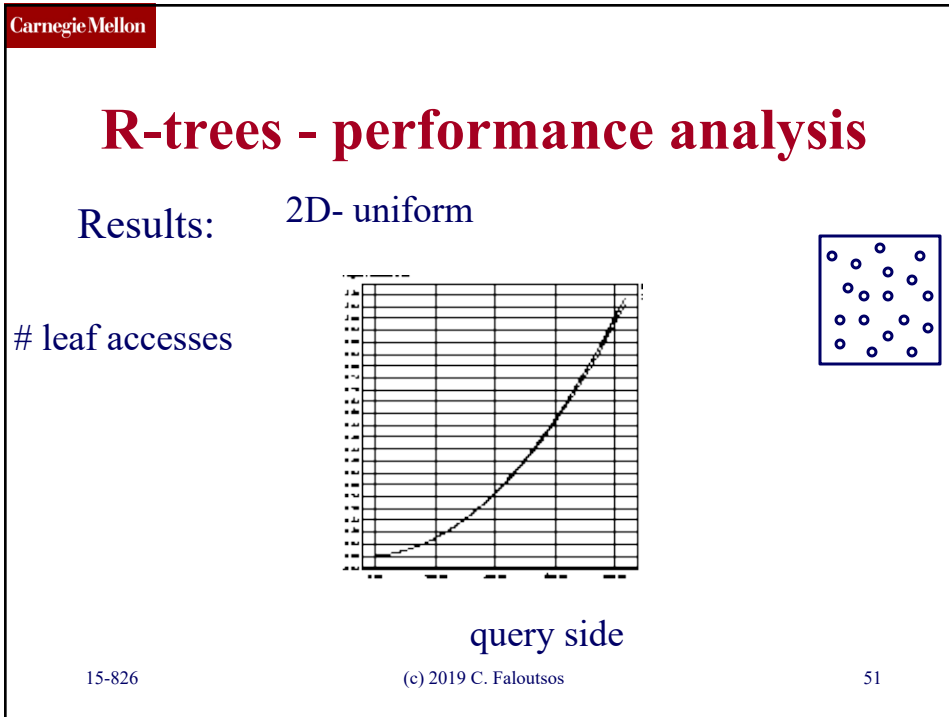


[a.] IUE - Leaf accesses vs. query side


query side







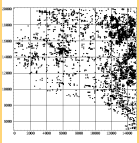
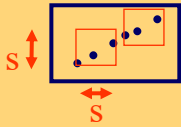
CarnegieMellon



## Solution:

- Selectivity of a range query in R-trees?
- Depends on \*fractal\* dimension

$$s = (C/N)^{1/D_0}$$

15-826 (c) 2019 C. Faloutsos 53

CarnegieMellon

## 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 (c) 2019 C. Faloutsos 54

## Indexing - Detailed outline


- fractals
  - intro
  - applications
    - disk accesses for R-trees (range queries)
    - dim. curse revisited
    - nearest neighbors estimation



## Must-read Material

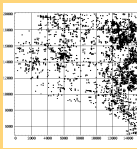
- Bernd-Uwe Pagel, Flip Korn and Christos Faloutsos, [\*Deflating the Dimensionality Curse using Multiple Fractal Dimensions\*](#), ICDE 2000, San Diego, CA, Feb. 2000.

CarnegieMellon




## Problem:

- Q: Do all S.A.M. suffer in high dimensions?
- Q: what to do?



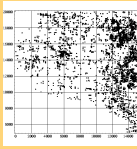
15-826 (c) 2019 C. Faloutsos 57

CarnegieMellon



## Solutions:

- Q: Do all S.A.M. suffer in high dimensions?
- A: Only in high \*fractal\* dimensions
- Q: what to do?
- A: dim-reduction; approximate knn; etc



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

15-826 (c) 2019 C. Faloutsos 58

## Indexing - Detailed outline

- fractals
  - intro
  - applications
    - disk accesses for R-trees (range queries)
    - dim. curse revisited
    - nearest neighbors estimation



## Dimensionality ‘curse’

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

CarnegieMellon

## 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

(c) 2019 C. Faloutsos

#61

CarnegieMellon

- 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

(c) 2019 C. Faloutsos

#62

CarnegieMellon

## Dimensionality ‘curse’

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

15-826

(c) 2019 C. Faloutsos

#63

CarnegieMellon

## Dimensionality ‘curse’

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

15-826

(c) 2019 C. Faloutsos

#64



CarnegieMellon

## Dimensionality ‘curse’

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

15-826

(c) 2019 C. Faloutsos

#65

CarnegieMellon

## Dimensionality ‘curse’

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

15-826

(c) 2019 C. Faloutsos

#66

CarnegieMellon

## Dim. reduction

a.k.a. feature selection/extraction:

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

15-826

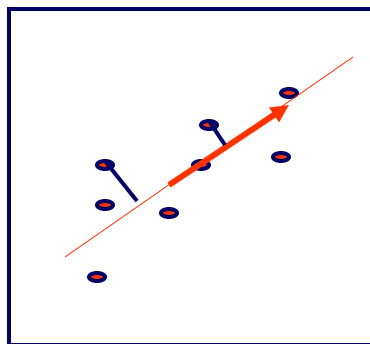
(c) 2019 C. Faloutsos

#67

CarnegieMellon

## Singular Value Decomposition (SVD)

- SVD (~LSI ~ KL ~ PCA ~ spectral analysis...)



LSI: S. Dumais; M. Berry

KL: eg, Duda+Hart

PCA: eg., Jolliffe

MANY more PROOF: soon

15-826

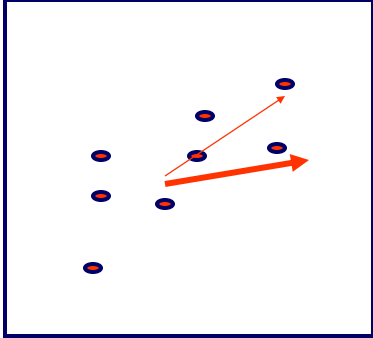
(c) 2019 C. Faloutsos

#68

CarnegieMellon

## Random projections

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



The diagram shows a square frame containing several blue dots representing data points. Two red arrows originate from a central point, pointing towards the right and slightly upwards, representing random directions for projection.

15-826 (c) 2019 C. Faloutsos #69

CarnegieMellon

## 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 (c) 2019 C. Faloutsos #70

CarnegieMellon

## Dim. reduction - w/ fractals

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

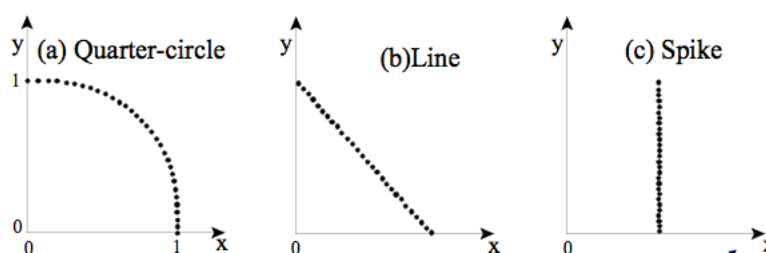
15-826

(c) 2019 C. Faloutsos

#71

CarnegieMellon

## Dim. reduction - w/ fractals

global  $FD=1$ 

15-826

(c) 2019 C. Faloutsos

#72

CarnegieMellon

## Dimensionality ‘curse’

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

15-826

(c) 2019 C. Faloutsos

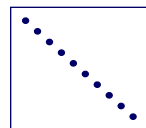
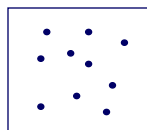
#73

CarnegieMellon

## Intrinsic dimensionality

- before we give up, compute the intrinsic dim.:
- the lower, the better... [Page1+, ICDE 2000]
- more PROOF: in a few foils

intr. d = 2



intr. d = 1

15-826

(c) 2019 C. Faloutsos

#74

CarnegieMellon

## Dimensionality ‘curse’

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

15-826

(c) 2019 C. Faloutsos

#75

CarnegieMellon

## 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

(c) 2019 C. Faloutsos

#76

CarnegieMellon

## Dimensionality ‘curse’

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

15-826

(c) 2019 C. Faloutsos

#77

CarnegieMellon

## Indexing - Detailed outline

- fractals
    - intro
    - applications
      - disk accesses for R-trees (range queries)
      - dim. curse revisited
      - nearest neighbors estimation
- ➔

15-826

(c) 2019 C. Faloutsos

78

## Estimation of knn effort

- (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

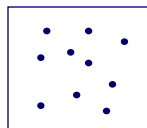


(c) 2019 C. Faloutsos

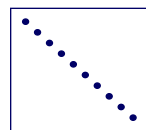
79

## (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



(c) 2019 C. Faloutsos

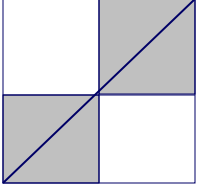
80



CarnegieMellon
PROOF

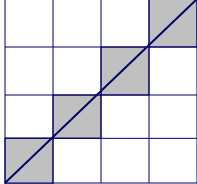
## Reminder: Hausdorff Dimension ( $D_0$ )

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



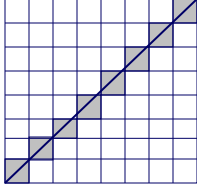
$r = 1/2 \quad B = 2$

$\log r = -1$   
 $\log B = 1$



$r = 1/4 \quad B = 4$

$\log r = -2$   
 $\log B = 2$



$r = 1/8 \quad B = 8$

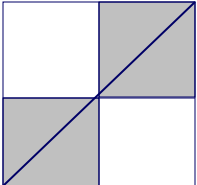
$\log r = -3$   
 $\log B = 3$

15-826
(c) 2019 C. Faloutsos
81

CarnegieMellon
PROOF

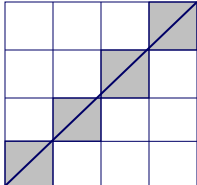
## Reminder: Correlation Dimension ( $D_2$ )

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



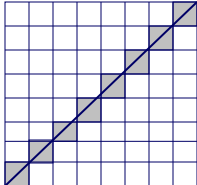
$r = 1/2 \quad S = 1/2$

$\log r = -1$   
 $\log S = -1$



$r = 1/4 \quad S = 1/4$

$\log r = -2$   
 $\log S = -2$



$r = 1/8 \quad S = 1/8$

$\log r = -3$   
 $\log S = -3$

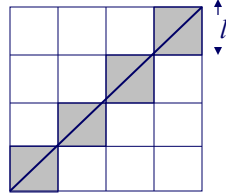
15-826
(c) 2019 C. Faloutsos
82

CarnegieMellon

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

(c) 2019 C. Faloutsos

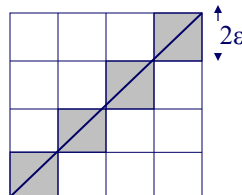
83

CarnegieMellon

PROOF

## Observation #2

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



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

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

15-826

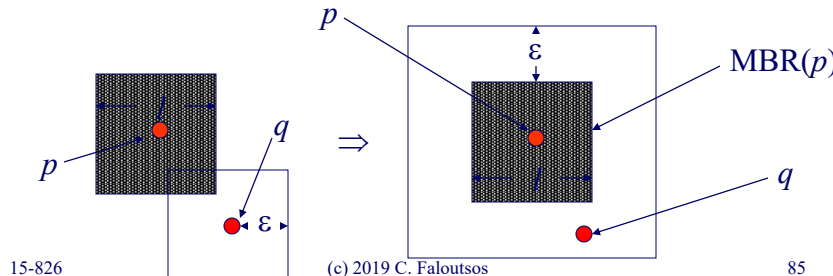
(c) 2019 C. Faloutsos

84

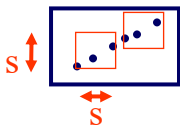
CarnegieMellon
PROOF

## Observation #3

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



15-826
(c) 2019 C. Faloutsos
85

CarnegieMellon


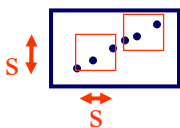
## Asymptotic Formula

- $k$ -NN page accesses as  $N \rightarrow \infty$ 
  - $C$  = page capacity
  - $D$  = fractal dimension ( $=D0 \sim D2$ )
  - $h$  = height of tree

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

15-826
(c) 2019 C. Faloutsos
86

CarnegieMellon



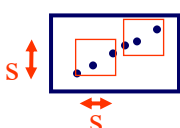
## Asymptotic Formula

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

- Observations?

15-826 (c) 2019 C. Faloutsos 87

CarnegieMellon

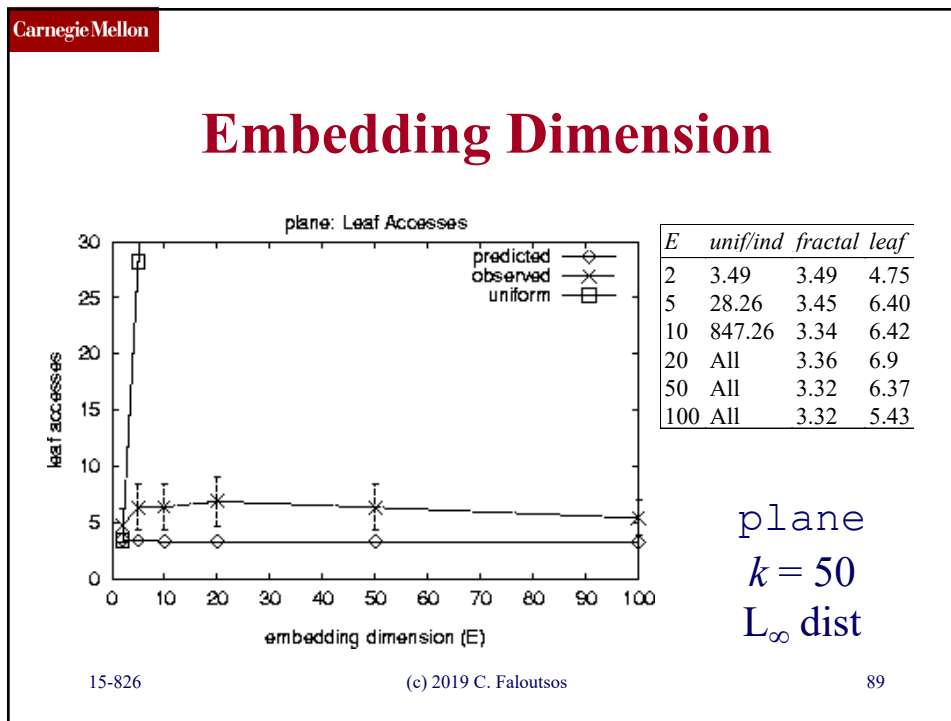


## Asymptotic Formula

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

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

15-826 (c) 2019 C. Faloutsos 88



CarnegieMellon

## A word of caution:

Nearest neighbors: may be meaningless!

Norio Katayama, Shin'ichi Satoh:  
Distinctiveness-Sensitive Nearest Neighbor Search for  
Efficient Similarity Retrieval of Multimedia Information.  
ICDE 2001: 493-502

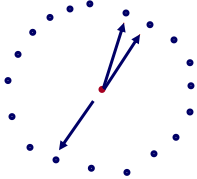
15-826 (c) 2019 C. Faloutsos 90

CarnegieMellon

## A word of caution:

Nearest neighbors: may be meaningless!

Norio Katayama, Shin'ichi Satoh:  
Distinctiveness-Sensitive Nearest Neighbor Search for  
Efficient Similarity Retrieval of Multimedia Information.  
ICDE 2001: 493-502



15-826 (c) 2019 C. Faloutsos 91

CarnegieMellon

## 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 (c) 2019 C. Faloutsos #92

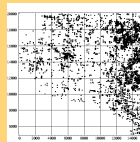
## Conclusions – cont' d

- 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

## Solutions:



- Q: Do all S.A.M. suffer in high dimensions?
- A: Only in high \*fractal\* dimensions
- Q: what to do?
- A: dim-reduction; approximate knn; etc



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

## 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/>

## References

- Berchtold, S., D. A. Keim, et al. (1996). The X-tree : An Index Structure for High-Dimensional Data. VLDB, Mumbai (Bombay), India.



## References cnt' d

- ➔ • 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.

## 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.

CarnegieMellon

## 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

(c) 2019 C. Faloutsos

99