**CMU SCS**

# 15-826: Multimedia Databases and Data Mining

Lecture #10: Fractals - case studies - I

*C. Faloutsos*

---
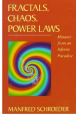
**CMU SCS**

# Must-read Material

- 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.

15-826 Copyright: C. Faloutsos (2017) 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 (on reserve in the WeH library)

15-826 Copyright: C. Faloutsos (2017) 3

---

**CMU SCS**

# Reminder

- Code at

www.cs.cmu.edu/~christos/SRC/fdnq_h.zip

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

15-826 Copyright: C. Faloutsos (2017) 4

**CMU SCS**

# Outline

Goal: 'Find similar / interesting things'

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

15-826 Copyright: C. Faloutsos (2017) 5

**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 (2017) 6
...

**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 (2017) 7

**CMU SCS**

# (Fractals mentioned before:)

- for performance analysis of R-trees
- fractals for dim. reduction
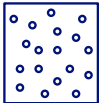
15-826 Copyright: C. Faloutsos (2017) 8

**CMU SCS**

# Case study#1: R-tree performance

Problem
- Given
  - N points in E-dim space

- Estimate # disk accesses for a range query
  (q1 x ... x $q_E$ )

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

**CMU SCS**

# Case study#1: R-tree performance

Problem
- Given
  - N points in E-dim space

- Estimate # disk accesses for a range query
  (q1 x ... x $q_E$ )

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

**CMU SCS**
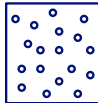
# Case study#1: R-tree performance

Problem
- Given
  - N points in E-dim space

- Estimate # disk accesses for a range query
  (q1 x ... x $q_E$ )

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

**CMU SCS**

# 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
  (q1 x ... x $q_E$ )

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

3

**CMU SCS**

# Examples:World's countries

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

pop        log(pop)



area            log(area)

**CMU SCS**

# Examples:World's countries

- neither uniform, nor independent!

pop        log(pop)



area            log(area)

**CMU SCS**

# For fun: identification



**CMU SCS**

# For fun: identification

**CMU SCS**

## For fun: identification

area vs population



- Highest density
- lowest density
- 47 residents(!)

**CMU SCS**

## For fun: identification

area vs population



- China
- Australia
- Mongolia
- Macau
- Greenland
- Vatican city
- Pitcairn
- 47 residents(!)

**CMU SCS**

## Examples: TIGER files

- neither uniform, nor independent!

MG county        LB county



15-826                    Copyright: C. Faloutsos (2017)                    19

**CMU SCS**

## How to proceed?

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

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

But:

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

15-826                    Copyright: C. Faloutsos (2017)                    20

**CMU SCS**

# How to proceed?

But:

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

Answer (jumping ahead):

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

---

**CMU SCS**

# How to proceed?

But:

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

Answer (jumping ahead):

$$s = (C/N)^{1/D0}$$ ← **Hausdorff fd**

side of (parent) MBR

# of data points

page capacity

---

**CMU SCS**                                                          **DETAILS**

# Let's see the rationale

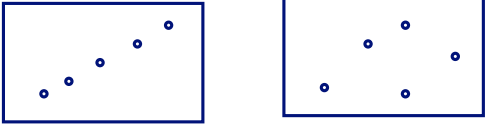$$s = (C/N)^{1/D0}$$

---

**CMU SCS**                                                          **DETAILS**

# 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

6

# R-trees - performance analysis

Q: OK - so we are told that the **Hausdorff** fractal dim. = D0 - Next step?

(also know that there are at most $C$ points per page)

D0=1

D0=2

---

# R-trees - performance analysis

Assumption1: square-like parents (s*s)
Assumption2: fully packed (C points each)
Assumption3: non-overlapping

D0=1

D0=2

s1=s2=s

---

# R-trees - performance analysis

Assumption1: square-like parents (s*s)
Assumption2: fully packed (N/C non-empty)
Assumption3: non-overlapping

D0=1

s1=s2=s

---

# R-trees - performance analysis

Hint: dfn of Hausdorff f.d.:

Felix Hausdorff (1868-1942)

**CMU SCS**

# Reminder:
# 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 D0:

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

15-826          Copyright: C. Faloutsos (2017)          29

---

**CMU SCS**

# Reminder

- Hausdorff  fd:

r ___     log(#non-empty cells)



log(r)

15-826          Copyright: C. Faloutsos (2017)          30

---

**CMU SCS**

# Reminder

- dfn of Hausdorff  fd implies that

$$N(r) \sim r^{(-D0)}$$

# non-empty cells of side *r*

15-826          Copyright: C. Faloutsos (2017)          31

---

**CMU SCS**

# R-trees - performance analysis

Q (rephrased): what is the side s1, s2, ... of parent nodes, given *N* data points, packed by *C,* with f.d. = *D0*

D0=1                    D0=2



15-826          Copyright: C. Faloutsos (2017)          32

# R-trees - performance analysis

Q (rephrased): what is the side s1, s2, ... of parent nodes, given *N* data points, packed by *C,* with f.d. = *D0*

D0=1

D0=2

s2

s1

# R-trees - performance analysis

Q (rephrased): what is the side s1, s2, ... of parent nodes, given *N* data points, packed by *C,* with f.d. = *D0*

D0=1

D0=2

s1=s2=s

# R-trees - performance analysis

A: (educated guess)

• s=s1=s2 (= ... ) - square-like MBRs
• N/C non-empty cells = K * $s^{(-D0)}$

*log(#cells)*

D0=1

D0=2

s2

s1

*log(s)*

# R-trees - performance analysis

Details of derivations: in [PODS 94].

s

Finally, expected side *s* of parent MBRs:

s

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

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

A:

9

**CMU SCS**

# R-trees - performance analysis

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

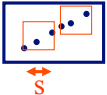Finally, expected side $s$ of parent MBRs:

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

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 Copyright: C. Faloutsos (2017) 37

---

**CMU SCS** **DETAILS**

# R-trees - performance analysis

Q: Final-final formula (# disk accesses for range queries $q1$ x $q2$ x ... ):

A:

15-826 Copyright: C. Faloutsos (2017) 38

---

**CMU SCS** **DETAILS**

# R-trees - performance analysis

Q: Final-final formula (# disk accesses for range queries $q1$ x $q2$ x ... ):

A: # of parent-node accesses:

$$N/C * (s + q1) * (s + q2) * ... (s + q_E)$$

A: # of grand-parent node accesses

15-826 Copyright: C. Faloutsos (2017) 39

---

**CMU SCS** **DETAILS**

# R-trees - performance analysis

Q: Final-final formula (# disk accesses for range queries $q1$ x $q2$ x ... ):

A: # of parent-node accesses:

$$N/C * (s + q1) * (s + q2) * ... (s + q_E)$$

A: # of grand-parent node accesses

$$N/(C^2) * (s' + q1) * (s' + q2) * ... (s' + q_E)$$
$$s' = ??$$

15-826 Copyright: C. Faloutsos (2017) 40

## R-trees - performance analysis

**DETAILS**

Q: Final-final formula (# disk accesses for range
    queries *q1* x *q2* x ... ):

A: # of parent-node accesses:

$$N/C * (s + q1) * (s + q2) * ... (s + q_E)$$

A: # of grand-parent node accesses

$$N/(C^2) * (s' + q1) * (s' + q2) * ... (s' + q_E)$$
$$s' = (C^2/N)^{1/D0}$$

## R-trees - performance analysis

Results:       IUE (x-y star coordinates)

# leaf accesses



(a) IUE - Leaf accesses vs. query s

query side

## R-trees - performance analysis

Results:       LB County

# leaf accesses



query side

## R-trees - performance analysis

Results:       MG-county

# leaf accesses



query side

**CMU SCS**

# R-trees - performance analysis

Results:          2D- uniform

# leaf accesses

query side

15-826          Copyright: C. Faloutsos (2017)          45

---

**CMU SCS**

# R-trees - performance analysis

Conclusions: usually, <5% relative error, for range queries

15-826          Copyright: C. Faloutsos (2017)          46

---

**CMU SCS**

Optional

# Indexing - Detailed outline

- fractals
  - intro
  - applications
    - ✓ disk accesses for R-trees (range queries)
    - dimensionality reduction
    - dim. curse revisited
    - quad-tree analysis [Gaede+]
    - ....

15-826          Copyright: C. Faloutsos (2017)          47

---

**CMU SCS**

Optional

# Case study #2: Dim. reduction

Problem definition: 'Feature selection'
- given $N$ points, with $E$ dimensions
- keep the $k$ most 'informative' dimensions [Traina+,SBBD'00]

Caetano    Agma    Leejay
Traina     Traina   Wu

15-826          Copyright: C. Faloutsos (2017)          48

**CMU SCS**

Optional

## Dim. reduction - w/ fractals

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

not informative

15-826    Copyright: C. Faloutsos (2017)    49

---

**CMU SCS**

Optional

## Dim. reduction

Problem definition: 'Feature selection'
- given $N$ points, with $E$ dimensions
- keep the $k$ most 'informative' dimensions

Re-phrased: spot and drop attributes with strong (non-)linear correlations

Q: how do we do that?

15-826    Copyright: C. Faloutsos (2017)

---

**CMU SCS**

Optional

## Dim. reduction

A: Hint: correlated attributes do not affect the intrinsic/fractal dimension, e.g., if

$$y = f(x,z,w)$$

we can drop $y$

(hence: 'partial fd' (PFD) of a set of attributes = the fd of the dataset, when projected on those attributes)

15-826    Copyright: C. Faloutsos (2017)

---

**CMU SCS**

Optional

## Dim. reduction - w/ fractals

**global FD=1**
PFD=1

y (c) Spike

PFD~0

15-826    Copyright: C. Faloutsos (2017)    52

13

**CMU SCS**

Optional

# Dim. reduction - w/ fractals

**global FD=1**
PFD=1


(b)Line

PFD=1

**CMU SCS**

Optional

# Dim. reduction - w/ fractals

**global FD=1**
PFD~1


(a) Quarter-circle

PFD~1

**CMU SCS**

Optional

# Dim. reduction - w/ fractals

- (problem: given N points in E-d, choose k best dimensions)
- Q: Algorithm?

**CMU SCS**

Optional

# Dim. reduction - w/ fractals

- Q: Algorithm?
- A: e.g., greedy - forward selection:
  - keep the attribute with highest partial fd
  - add the one that causes the highest increase in pfd
  - etc., until we are within *epsilon* from the full f.d.

**CMU SCS**

Optional

# Dim. reduction - w/ fractals

- (backward elimination: ~ reverse)
  – drop the attribute with least impact on the p.f.d.
  – repeat
  – until we are *epsilon* below the full f.d.

**CMU SCS**

Optional

# Dim. reduction - w/ fractals

- Q: what is the smallest # of attributes we should keep?

**CMU SCS**

Optional

# Dim. reduction - w/ fractals

- Q: what is the smallest # of attributes we should keep?
- A: we should keep at least as many as the f.d. (and probably, a few more)
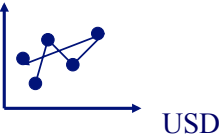
**CMU SCS**
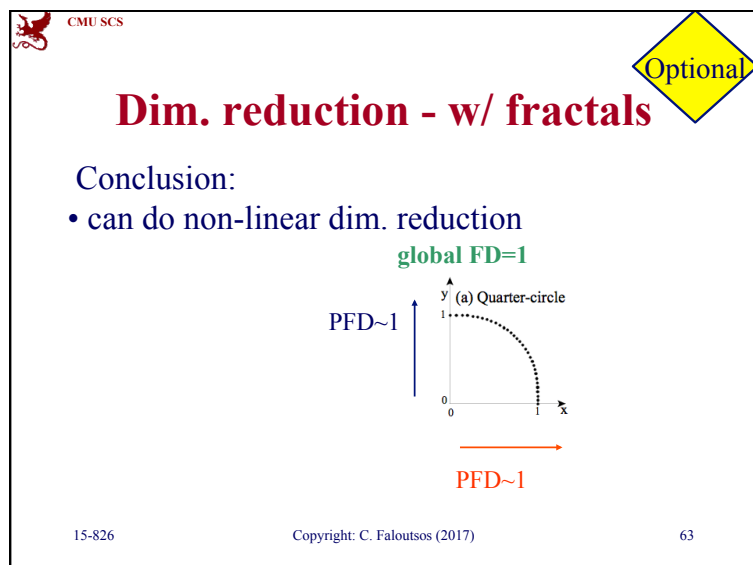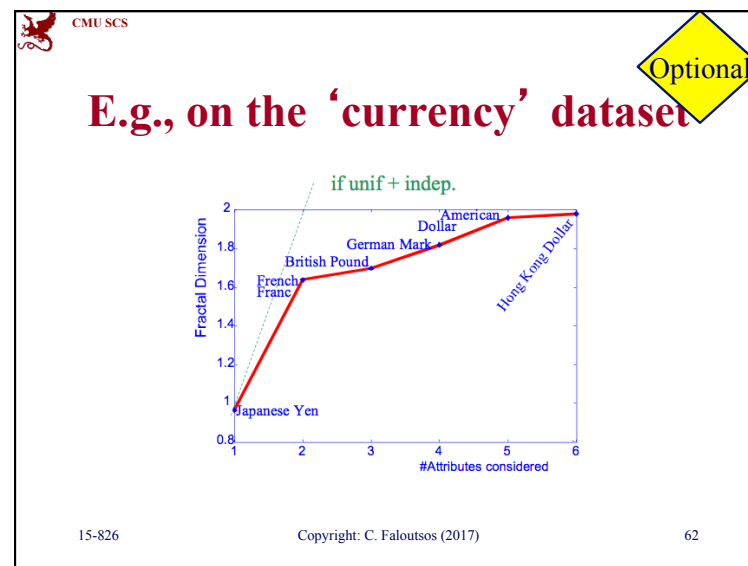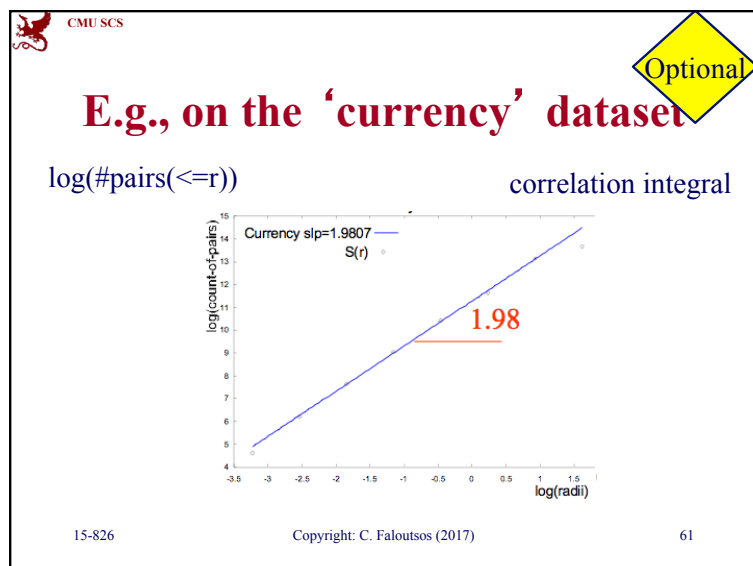
Optional

# Dim. reduction - w/ fractals

- Results: E.g., on the 'currency' dataset
- (daily exchange rates for USD, HKD, BP, FRF, DEM, JPY - i.e., 6-d vectors, one per day - base currency: CAD)

e.g.:  FRF

USD

**E.g., on the 'currency' dataset**

Optional

log(#pairs(<=r))

correlation integral

1.98

**E.g., on the 'currency' dataset**

Optional

if unif + indep.

American Dollar
German Mark
British Pound
French Franc
Hong Kong Dollar
Japanese Yen

Fractal Dimension

#Attributes considered

**Dim. reduction - w/ fractals**

Optional

Conclusion:
• can do non-linear dim. reduction

global FD=1

PFD~1

(a) Quarter-circle

PFD~1

**References**

• [PODS94] Faloutsos, C. and I. Kamel (May 24-26, 1994). *Beyond Uniformity and Independence: Analysis of R-trees Using the Concept of Fractal Dimension*. Proc. ACM SIGACT-SIGMOD-SIGART PODS, Minneapolis, MN.
• [Traina+, SBBD'00] Traina, C., A. Traina, et al. (2000). *Fast feature selection using the fractal dimension*. XV Brazilian Symposium on Databases (SBBD), Paraiba, Brazil.