**CMU SCS**

# Indexing and Mining Time Sequences

*Christos Faloutsos and Lei Li*
CMU

---

**CMU SCS**

## Outline

➤ • Motivation
• Similarity Search and Indexing
• DSP (Digital Signal Processing)
• Linear Forecasting
• Kalman filters
• fractals and multifractals
• Non-linear forecasting
• Conclusions

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          2

---

**CMU SCS**

## Problem definition

• <u>Given</u>: one or more sequences
  $x_1, x_2, \ldots, x_t, \ldots$
  $(y_1, y_2, \ldots, y_t, \ldots$
  $\ldots )$
• <u>Find</u>
  – similar sequences; forecasts
  – patterns; clusters; outliers

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          3

---

**CMU SCS**

## Motivation - Applications

• Financial, sales, economic series
• Medical
  – reactions to new drugs
  – elderly care

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          4

---

**CMU SCS**

## ECG - physionet.org



KDD 2010          5

---

**CMU SCS**

## EEG - epilepsy



KDD 2010          from wikipedia          6

**CMU SCS**

## Motivation - Applications (cont'd)
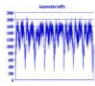
- 'Smart house'
  - sensors monitor temperature, humidity, air quality
- video surveillance

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          7

**CMU SCS**

## Motivation - Applications (cont'd)

- civil/automobile infrastructure
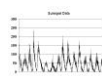  - bridge vibrations [Oppenheim+02]
  - road conditions / traffic monitoring



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          8

**CMU SCS**

## Motivation - Applications (cont'd)

- Weather, environment/anti-pollution
  - volcano monitoring
  - air/water pollutant monitoring



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          9
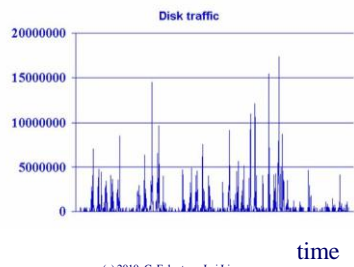
**CMU SCS**

## Motivation - Applications (cont'd)

- Computer systems
  - 'Active Disks' (buffering, prefetching)
  - web servers (ditto)
  - network traffic monitoring
  - ...

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          10

**CMU SCS**

## Stream Data: Disk accesses

#bytes



time

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          11

**CMU SCS**

## Problem #1:

Goal: given a signal (eg., #packets over time)
Find: patterns, periodicities, and/or compress

count



lynx caught per year
(packets per day;
temperature per day)

year

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          12

---

**CMU SCS**

# Problem#2: Forecast

Given $x_t, x_{t-1}, \ldots$, forecast $x_{t+1}$
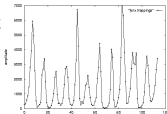


KDD 2010     (c) 2010, C. Faloutsos, Lei Li     13

---

**CMU SCS**

# Problem#2': Similarity search

Eg., Find a 3-tick pattern, similar to the last one



KDD 2010     (c) 2010, C. Faloutsos, Lei Li     14

---

**CMU SCS**

# Problem #3:

- Given: A set of **correlated** time sequences
- Forecast 'Sent(t)'



KDD 2010     (c) 2010, C. Faloutsos, Lei Li     15
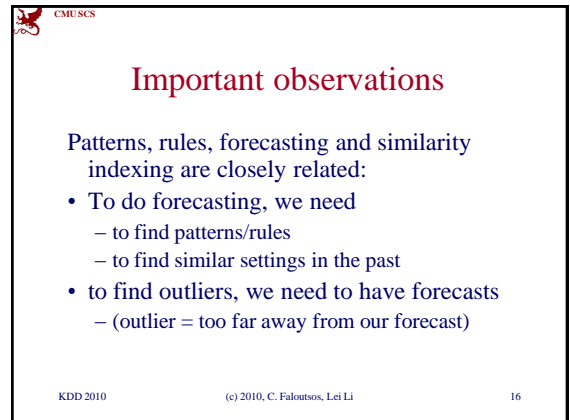
---

**CMU SCS**

# Important observations

Patterns, rules, forecasting and similarity indexing are closely related:

- To do forecasting, we need
  - to find patterns/rules
  - to find similar settings in the past
- to find outliers, we need to have forecasts
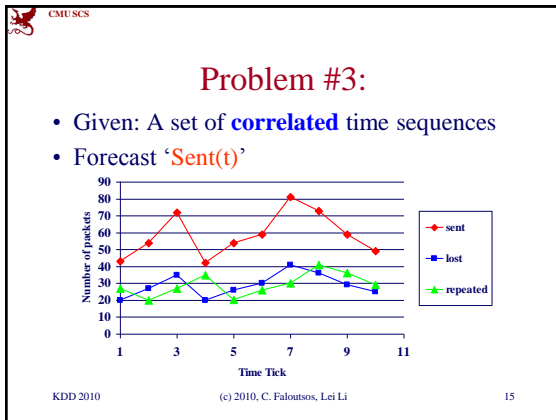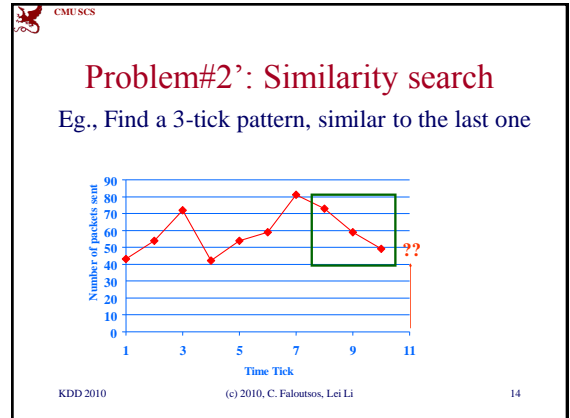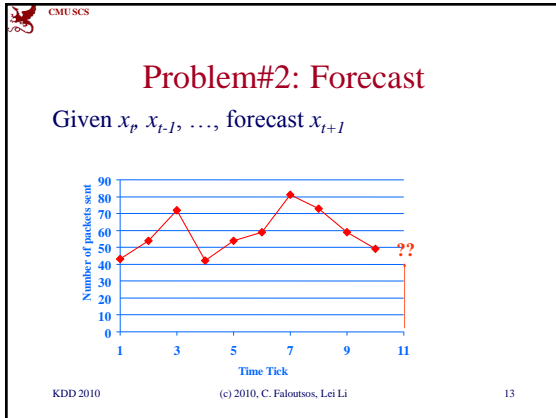  - (outlier = too far away from our forecast)

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     16

---

**CMU SCS**

# Important topics NOT in this tutorial:

- Continuous queries
  - [Babu+Widom ] [Gehrke+] [Madden+]
- Categorical data streams
  - [Hatonen+96]
- Outlier detection (discontinuities)
  - [Breunig+00]

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     17

---

**CMU SCS**

# Outline

- Motivation
- ➡ Similarity Search and Indexing
- DSP
- Linear Forecasting
- Bursty traffic - fractals and multifractals
- Non-linear forecasting
- Conclusions

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     18

**CMU SCS**

## Outline

- Motivation
- Similarity Search and Indexing
  → – distance functions: Euclidean;Time-warping
  – indexing
  – feature extraction
- DSP
- ...

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 19

**CMU SCS**

## Importance of distance functions

Subtle, but **absolutely necessary**:
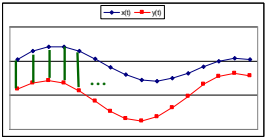- A 'must' for similarity indexing (-> forecasting)
- A 'must' for clustering

Two major families
  – Euclidean and Lp norms
  – Time warping and variations

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 20

**CMU SCS**

## Euclidean and Lp
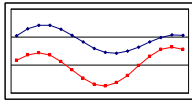
$$D(\vec{x}, \vec{y}) = \sum_{i=1}^{n} (x_i - y_i)^2$$

$$L_p(\vec{x}, \vec{y}) = \sum_{i=1}^{n} |x_i - y_i|^p$$

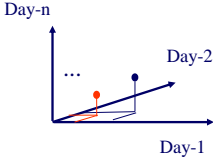- $L_1$: city-block = Manhattan
- $L_2$ = Euclidean
- $L_\infty$

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 21

**CMU SCS**

## Observation #1

- **Time sequence -> n-d vector**

Day-n
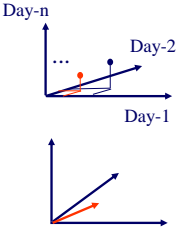
... Day-2

Day-1

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 22

**CMU SCS**

## Observation #2

Euclidean distance is closely related to
  – cosine similarity
  – dot product
  – 'cross-correlation' function

Day-n

... Day-2

Day-1

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 23

**CMU SCS**

## Time Warping

- allow accelerations - decelerations
  – (with or w/o penalty)
- THEN compute the (Euclidean) distance (+ penalty)
- related to the string-editing distance

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 24

## Slide 1

# Time Warping

'stutters':

## Slide 2

**Skip**

# Time Warping

Q: how to compute it?

A: dynamic programming

$D(i, j)$ = cost to match

prefix of length $i$ of first sequence $x$ with prefix of length $j$ of second sequence $y$

## Slide 3

**Skip**

# Time Warping

Thus, with no penalty for stutter, for sequences

$$x_1, x_2, \ldots, x_{i,;} \qquad y_1, y_2, \ldots, y_j$$

$$D(i, j) = \|x[i] - y[j]\| + \min \begin{cases} D(i-1, j-1) & \text{no stutter} \\ D(i, j-1) & \text{x-stutter} \\ D(i-1, j) & \text{y-stutter} \end{cases}$$

## Slide 4

# Time Warping

• Time warping matrix & optimal path:

No stutters

Y

X

## Slide 5

# Time Warping

• Time warping matrix & optimal path:

All stutters
$Y_1$ x N times;
$X_N$ x M times

Y

X

## Slide 6

# Time Warping - variations

• Time warping matrix & optimal path:

At most k stutters:
Sakoe-Chiba band

Y

X

**CMU SCS**

## Time Warping - variations

- Time warping matrix & optimal path:



At most x% stutters:
Itakura parallelogram

Y

X

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          Part2.1 #31
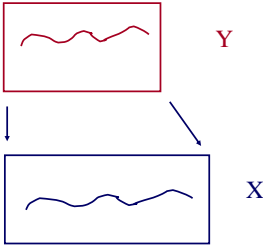
---

**CMU SCS**

## Time warping

- Complexity: O(M*N) - quadratic on the length of the strings
- **<u>Many</u>** variations (penalty for stutters; limit on the number/percentage of stutters; …)
- popular in voice processing [Rabiner+Juang]
- Seems suitable for mo-cap

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          Part2.1 #32

---

**CMU SCS**

## A variation: Uniform axis scaling



Y

X

- Stretch / shrink time axis of Y, up to p%, for free
- THEN compute Euclidean distance
- [Keogh+, VLDB04]

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          Part2.1 #33

---

**CMU SCS**

## Other Distance functions

- piece-wise linear/flat approx.; compare pieces [Keogh+01] [Faloutsos+97]
- 'cepstrum' (for voice [Rabiner+Juang])
  - do DFT; take log of amplitude; do DFT again!
- Allow for small gaps [Agrawal+95]

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          34

---

**CMU SCS**

## More distance functions.

- Chen + Ng [vldb'04]: ERP 'Edit distance with Real Penalty': give a penalty to stutters
- Keogh+ [kdd'04]: VERY NICE, based on information theory: compress each sequence (quantize + Lempel-Ziv), using the **other** sequences' LZ tables

*On The Marriage of Lp-norms and Edit Distance,* Lei Chen, Raymond T. Ng:, VLDB'04

*Towards Parameter-Free Data Mining,* E. Keogh, S. Lonardi, C.A. Ratanamahatana, KDD'04

---

**CMU SCS**

## Conclusions

Prevailing distances:
  - Euclidean and
  - time-warping

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          36

**CMU SCS**

# Outline

- Motivation
- Similarity Search and Indexing
  - distance functions
  ➡ – indexing
  - feature extraction
- DSP
- ...

KDD 2010         (c) 2010, C. Faloutsos, Lei Li         37

---

**CMU SCS**

# Indexing

Problem:
- given a set of time sequences,
- find the ones similar to a desirable query sequence

KDD 2010         (c) 2010, C. Faloutsos, Lei Li         38

---

**CMU SCS**



$price

$price

$price

1    365
day

1    365
day

1    365
day

distance function: by expert

KDD 2010         (c) 2010, C. Faloutsos, Lei Li         39

---

**CMU SCS**

# Idea: 'GEMINI'

Eg., '*find stocks similar to MSFT*'
Seq. scanning: too slow
How to accelerate the search?
[Faloutsos96]

KDD 2010         (c) 2010, C. Faloutsos, Lei Li         40

---

**CMU SCS**

# 'GEMINI' - Pictorially



eg,. std

• F(S1)

S1

1    365
day

•
F(Sn)

Sn

eg, avg

1    365
day

KDD 2010         (c) 2010, C. Faloutsos, Lei Li         41

---

**CMU SCS**

# GEMINI

Solution: Quick-and-dirty' filter:
- extract *n* features (numbers, eg., avg., etc.)
- map into a point in *n*-d feature space
- organize points with off-the-shelf spatial access method ('SAM')
- discard false alarms

KDD 2010         (c) 2010, C. Faloutsos, Lei Li         42

**CMU SCS**

## Examples of GEMINI

- Time sequences: DFT (up to 100 times faster) [SIGMOD94];
- [Kanellakis+], [Mendelzon+]

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          43

---

**CMU SCS**

## Examples of GEMINI

Even on other-than-sequence data:
- Images (QBIC) [JIIS94]
- tumor-like shapes [VLDB96]
- video [Informedia + S-R-trees]
- automobile part shapes [Kriegel+97]

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          44

---

**CMU SCS**

## Indexing - SAMs

Q: How do Spatial Access Methods (SAMs) work?

A: they group nearby points (or regions) together, on nearby disk pages, and answer spatial queries quickly ('range queries', 'nearest neighbor' queries etc)

For example:

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          45

---

**CMU SCS**

Skip

## R-trees

- [Guttman84] eg., w/ fanout 4: group nearby rectangles to parent MBRs; each group -> disk page



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          46

---

**CMU SCS**

Skip

## R-trees

- eg., w/ fanout 4:



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          47

---

**CMU SCS**

Skip

## R-trees

- eg., w/ fanout 4:



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          48

## R-trees - range search?

Skip

P1  P3  I

A C
B
E
P2 D
G F H J
P4

P1 P2 P3 P4
A B C | H I J
D E | F G

## R-trees - range search?

Skip

P1  P3  I

A C
B
E
P2 D
G F H J
P4

P1 P2 P3 P4
A B C | H I J
D E | F G

## Conclusions

- Fast indexing: through GEMINI
  - feature extraction and
  - (off the shelf) Spatial Access Methods [Gaede+98]

## Outline

- Motivation
- Similarity Search and Indexing
  - distance functions
  - indexing
  - feature extraction
- DSP
- ...

## Outline

- Motivation
- Similarity Search and Indexing
  - distance functions
  - indexing
  - feature extraction
    - DFT, DWT, DCT (data independent)
    - SVD, etc (data dependent)
    - MDS, FastMap

## DFT and cousins

- very good for compressing real signals
- more details on DFT/DCT/DWT: later

## DFT and stocks

CMU SCS

• Dow Jones Industrial index, 6/18/2001-12/21/2001

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 55

## DFT and stocks

CMU SCS

• Dow Jones Industrial index, 6/18/2001-12/21/2001
• just 3 DFT coefficients give very good approximation

Log(ampl)

freq

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 56

## Outline

CMU SCS

• Motivation
• Similarity Search and Indexing
  – distance functions
  – indexing
  – feature extraction
    • DFT, DWT, DCT (data independent)
    • SVD etc (data dependent)
    • MDS, FastMap

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 57

## SVD

CMU SCS

• <u>THE</u> optimal method for dimensionality reduction
  – (under the Euclidean metric)

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 58

## Singular Value Decomposition (SVD)

CMU SCS

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

day2

LSI: S. Dumais; M. Berry
KL: eg, Duda+Hart
PCA: eg., Jolliffe
Details: [Press+],
[Faloutsos96]

day1

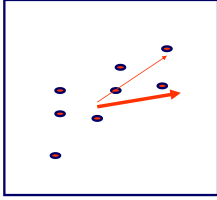KDD 2010 (c) 2010, C. Faloutsos, Lei Li 59

## SVD

CMU SCS

• **<u>Extremely</u>** useful tool
  – (also behind PageRank/google and Kleinberg's algorithm for hubs and authorities)
• But may be slow: O($N * M * M$) if $N>M$
• any approximate, faster method?

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 60

10

## Slide 61

**CMU SCS**

# SVD shorcuts

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

KDD 2010            (c) 2010, C. Faloutsos, Lei Li            61

## Slide 62

**CMU SCS**

# Random projections

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

KDD 2010            (c) 2010, C. Faloutsos, Lei Li            62

## Slide 63

**CMU SCS**

# Outline

- Motivation
- Similarity Search and Indexing
    - distance functions
    - indexing
    - feature extraction
        - DFT, DWT, DCT (data independent)
        - SVD etc (data dependent), **ICA**
        - MDS, FastMap

KDD 2010            (c) 2010, C. Faloutsos, Lei Li            63

## Slide 64

**CMU SCS**

# Citation

- *AutoSplit: Fast and Scalable Discovery of Hidden Variables in Stream and Multimedia Databases,* **Jia-Yu Pan**, Hiroyuki Kitagawa, Christos Faloutsos and Masafumi Hamamoto
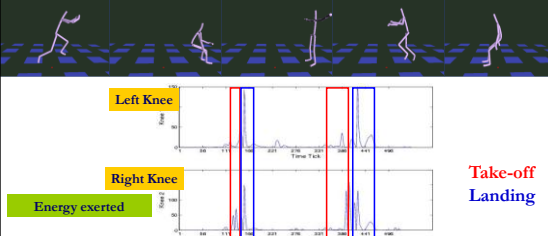PAKDD 2004, Sydney, Australia

KDD 2010            (c) 2010, C. Faloutsos, Lei Li            64

## Slide 65

**CMU SCS**

# Motivation:
# (Q1) Find patterns in data

- Motion capture data (broad jumps)

Left Knee

Right Knee

Energy exerted

**Take-off**
**Landing**
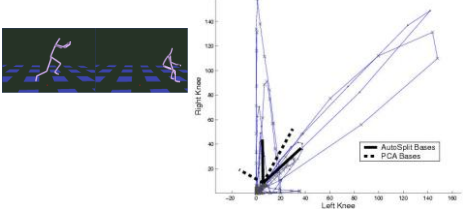
KDD 2010            (c) 2010, C. Faloutsos, Lei Li            65

## Slide 66

**CMU SCS**

# PCA sometimes misses essential features

- Best SVD axis: not always meaningful!

AutoSplit Bases
PCA Bases

KDD 2010            (c) 2010, C. Faloutsos, Lei Li            66

11

**Slide 67**

CMU SCS

# Motivation:
## (Q1) Find patterns in data

- Human would say
  - Pattern 1: along diagonal
  - Pattern 2: along vertical axis
- How to find these automatically?

Right Knee

60:1

1:1

Left Knee

— AutoSplit Bases
- - PCA Bases

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     67

**Slide 68**

CMU SCS

# Motivation:
## (Q2) Find hidden variables

Alcoa

American Express

Boeing

Caterpillar

Citi Group

Find common hidden variables, and weights.

Dow Jones Industrial Average

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     68

**Slide 69**

CMU SCS

# Motivation:
## (Q2) Find hidden variables

CAT

INTC

Caterpillar

Intel

$B_{1,CAT}$   $B_{1,INTC}$   ?   $B_{2,CAT}$   $B_{2,INTC}$

?

?

Hidden variable 1

Hidden variable 2

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     69

**Slide 70**

CMU SCS

# Motivation:
## (Q2) Find hidden variables

CAT

INTC

Caterpillar

Intel

0.94   0.63   0.03   0.64

"Hidden variable 1"

"Hidden variable 2"

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     70

**Slide 71**

CMU SCS

# Motivation:
## (Q2) Find hidden variables

CAT

INTC

Caterpillar

Intel

0.94   0.63   0.03   0.64

"General trend"

"Internet bubble"

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     71

**Slide 72**

CMU SCS

# Motivation:
## Find hidden variables

- ICA: also known as 'Blind Source Separation'
- 'cocktail party problem'
  - in a party, we can hear two concurrent conversations,
  - but separate them (and tune-in on one of them only)
- http://www.cnl.salk.edu/~tewon/Blind/blind_audio.html
- (in stocks: one 'discussion' is the general economy trend; the other 'discussion' is the tech-stock boom

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     72

---

**CMU SCS**

## Problem formulation

- Given n data items, each has m attributes
- Find the m <u>hidden variables</u> and the m <u>bases</u>

$$X_{11}, X_{12}, \ldots, X_{1m} \quad H_{11}, H_{12}, \ldots, H_{1m} \quad B_{11}, B_{12}, \ldots, B_{1m}$$
$$\ldots \quad = \quad \ldots \quad ? \quad \ldots \quad ?$$
$$X_{n1}, X_{n2}, \ldots, X_{nm} \quad H_{n1}, H_{n2}, \ldots, H_{nm} \quad B_{m1}, B_{m2}, \ldots, B_{mm}$$

Samples of the m-th hidden variable

### **X=HB**

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          73

---

**CMU SCS**

## Formulation: (Q1) Find patterns in data

$$X_{11}, X_{12} \quad H_{11}, H_{12} \quad B_{11}, B_{12} \qquad \text{Basis 1}$$
$$\ldots \quad = \quad \ldots \quad ? \quad B_{21}, B_{22}$$
$$X_{n1}, X_{n2} \quad H_{n1}, H_{n2}$$

Left Knee    Right Knee



(c) 2010, C. Faloutsos, Lei Li          74

---
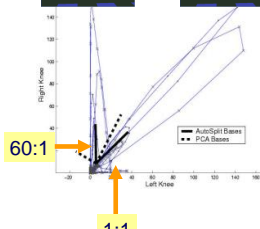
**CMU SCS**

## Q1: Find patterns

Take-off     Landing

- Patterns found

60:1

m=2, n=550

1:1



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          75

---

**CMU SCS**

## Q1: Find patterns

Take-off     Landing

- Patterns found
  - Landing: both knees

60:1

m=2, n=550

1:1



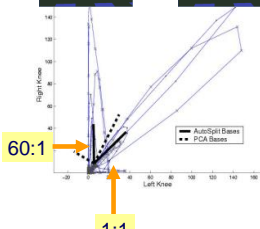KDD 2010          (c) 2010, C. Faloutsos, Lei Li          76

---

**CMU SCS**

## Q1: Find patterns

Take-off     Landing

- Patterns found
  - Landing: both knees
  - Take-off: right knee
  - Right-handed actor

60:1

m=2, n=550

1:1



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          77

---

**CMU SCS**

## Q2: Find hidden variables (DJIA stocks)

- Weekly DJIA closing prices
  - 01/02/1990-08/05/2002, n=660 data points
  - A data point: prices of 29 companies at the time

Alcoa
American Express
Boeing
Caterpillar
Citi Group



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          78

---

---

**CMU SCS**

## (Q2) Characterize hidden variable by the companies it influences

CAT | INTC

Caterpillar | Intel

B$_{1,CAT}$  0.94
B$_{1,INTC}$  0.63    0.03    0.64  B$_{2,INTC}$
B$_{2,CAT}$

"General trend" | "Internet bubble"

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          79

---

**CMU SCS**

## Companies related to hidden variable 1

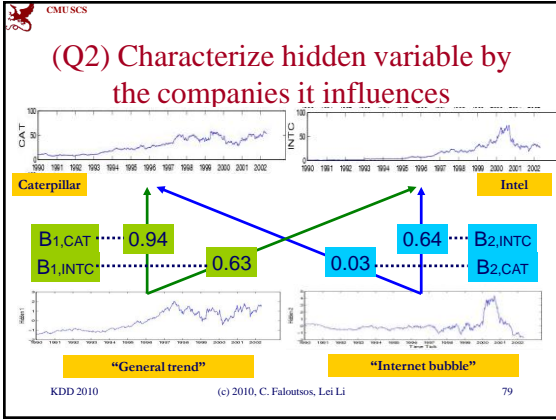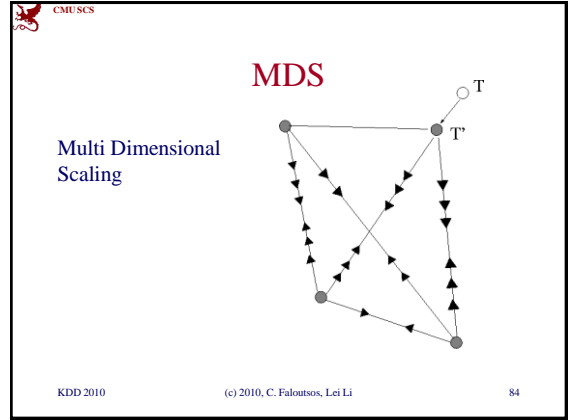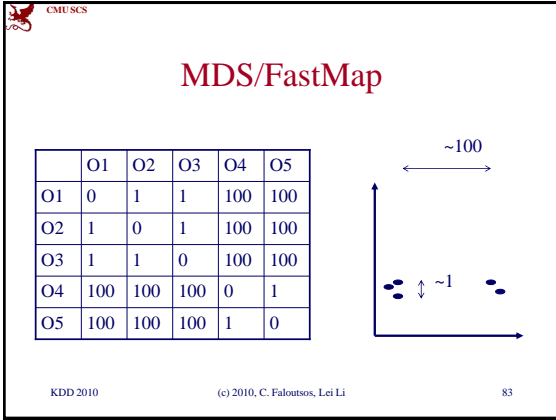| B$_{1,j}$ | | | |
|---|---|---|---|
| Highest | | Lowest | |
| Caterpillar | 0.938512 | AT&T | 0.021885 |
| Boeing | 0.911120 | WalMart | 0.624570 |
| MMM | 0.906542 | Intel | 0.638010 |
| Coca Cola | 0.903858 | Home Depot | 0.647774 |
| Du Pont | 0.900317 | Hewlett-Packard | 0.658768 |

Hidden1

"General trend"

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          80

---

**CMU SCS**

## Outline

- Motivation
- Similarity Search and Indexing
  - distance functions
  - indexing
  - feature extraction
    - DFT, DWT, DCT (data independent)
    - SVD (data dependent)
    - MDS, FastMap

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          81

---

**CMU SCS**

## MDS / FastMap

- but, what if we have NO points to start with?
  (eg. Time-warping distance)
- A: Multi-dimensional Scaling (MDS) ; FastMap

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          82

---

**CMU SCS**

## MDS/FastMap

| | O1 | O2 | O3 | O4 | O5 |
|---|---|---|---|---|---|
| O1 | 0 | 1 | 1 | 100 | 100 |
| O2 | 1 | 0 | 1 | 100 | 100 |
| O3 | 1 | 1 | 0 | 100 | 100 |
| O4 | 100 | 100 | 100 | 0 | 1 |
| O5 | 100 | 100 | 100 | 1 | 0 |

~100

~1

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          83

---

**CMU SCS**

## MDS

Multi Dimensional Scaling

T
T'

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          84
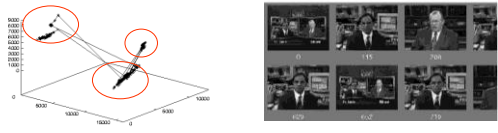
---

14

**CMU SCS**

# FastMap

- Multi-dimensional scaling (MDS) can do that, but in O(N**2) time
- FastMap [Faloutsos+95] takes O(N) time

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    85

---

**CMU SCS**

# FastMap: Application

VideoTrails [Kobla+97]



scene-cut detection (about 10% errors)

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    86

---

**CMU SCS**

# Outline

- Motivation
- Similarity Search and Indexing
  - distance functions
  - indexing
  - feature extraction
    - DFT, DWT, DCT (data independent)
    - SVD (data dependent)
    - MDS, FastMap, **IsoMap** etc

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    87

---

**CMU SCS**

# Variations

- Isomap [Tenenbaum, de Silva, Langford, 2000]
- LLE (Local Linear Embedding) [Roweis, Saul, 2000]
- MVE (Minimum Volume Embedding) [Shaw & Jebara, 2007]

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    #88

---

**CMU SCS**

# Variations

- Isomap [Tenenbaum, de Silva, Langford, 2000]
- LLE (Local Linear Embedding) [Roweis, Saul, 2000]
- MVE (Minimum Volume Embedding) [Shaw & Jebara, 2007]

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    #89

---

**CMU SCS**

# Outline

- Motivation
- Similarity Search and Indexing
  - distance functions
  - indexing
  - feature extraction
    - DFT, DWT, DCT (data independent)
    - SVD (data dependent)
    - MDS, FastMap

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    90

**CMU SCS**

# Conclusions - Practitioner's guide

Similarity search in time sequences

1) establish/choose distance (Euclidean, time-warping,…)

2) extract features (SVD, DWT, MDS), and use an SAM (R-tree/variant) or a Metric Tree (M-tree)

2') for high <u>intrinsic</u> dimensionalities, consider sequential scan (it might win…)

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          91

---

**CMU SCS**

# Books

- William H. Press, Saul A. Teukolsky, William T. Vetterling and Brian P. Flannery: *Numerical Recipes in C*, Cambridge University Press, 1992, 2nd Edition. (Great description, intuition and code for SVD)
- C. Faloutsos: *Searching Multimedia Databases by Content*, Kluwer Academic Press, 1996 (introduction to SVD, and GEMINI)

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          92

---

**CMU SCS**

# References

- Agrawal, R., K.-I. Lin, et al. (Sept. 1995). Fast Similarity Search in the Presence of Noise, Scaling and Translation in Time-Series Databases. Proc. of VLDB, Zurich, Switzerland.
- Babu, S. and J. Widom (2001). "Continuous Queries over Data Streams." SIGMOD Record 30(3): 109-120.
- Breunig, M. M., H.-P. Kriegel, et al. (2000). LOF: Identifying Density-Based Local Outliers. SIGMOD Conference, Dallas, TX.
- Berry, Michael: http://www.cs.utk.edu/~lsi/

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          93

---

**CMU SCS**

# References

- Ciaccia, P., M. Patella, et al. (1997). M-tree: An Efficient Access Method for Similarity Search in Metric Spaces. VLDB.
- Foltz, P. W. and S. T. Dumais (Dec. 1992). "Personalized Information Delivery: An Analysis of Information Filtering Methods." Comm. of ACM (CACM) 35(12): 51-60.
- Guttman, A. (June 1984). R-Trees: A Dynamic Index Structure for Spatial Searching. Proc. ACM SIGMOD, Boston, Mass.

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          94

---

**CMU SCS**

# References

- Gaede, V. and O. Guenther (1998). "Multidimensional Access Methods." Computing Surveys 30(2): 170-231.
- Gehrke, J. E., F. Korn, et al. (May 2001). On Computing Correlated Aggregates Over Continual Data Streams. ACM Sigmod, Santa Barbara, California.

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          95

---

**CMU SCS**

# References

- Gunopulos, D. and G. Das (2001). Time Series Similarity Measures and Time Series Indexing. SIGMOD Conference, Santa Barbara, CA.
- Eamonn J. Keogh, Themis Palpanas, Victor B. Zordan, Dimitrios Gunopulos, Marc Cardle: Indexing Large Human-Motion Databases. VLDB 2004: 780-791

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          Part2.1 #96

16

**CMU SCS**

# References

- Hatonen, K., M. Klemettinen, et al. (1996). Knowledge Discovery from Telecommunication Network Alarm Databases. ICDE, New Orleans, Louisiana.
- Jolliffe, I. T. (1986). Principal Component Analysis, Springer Verlag.

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          97

---

**CMU SCS**

# References

- Keogh, E. J., K. Chakrabarti, et al. (2001). Locally Adaptive Dimensionality Reduction for Indexing Large Time Series Databases. SIGMOD Conference, Santa Barbara, CA.
- Kobla, V., D. S. Doermann, et al. (Nov. 1997). VideoTrails: Representing and Visualizing Structure in Video Sequences. ACM Multimedia 97, Seattle, WA.

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          98

---

**CMU SCS**

# References

- Oppenheim, I. J., A. Jain, et al. (March 2002). A MEMS Ultrasonic Transducer for Resident Monitoring of Steel Structures. SPIE Smart Structures Conference SS05, San Diego.
- Papadimitriou, C. H., P. Raghavan, et al. (1998). Latent Semantic Indexing: A Probabilistic Analysis. PODS, Seattle, WA.
- Rabiner, L. and B.-H. Juang (1993). Fundamentals of Speech Recognition, Prentice Hall.

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          99

---

**CMU SCS**

# References

- Traina, C., A. Traina, et al. (October 2000). Fast feature selection using the fractal dimension,. XV Brazilian Symposium on Databases (SBBD), Paraiba, Brazil.

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          100

---

**CMU SCS**

# References

- Dennis Shasha and Yunyue Zhu *High Performance Discovery in Time Series: Techniques and Case Studies* Springer 2004
- Yunyue Zhu, Dennis Shasha ``*StatStream: Statistical Monitoring of Thousands of Data Streams in Real Time*'' VLDB, August, 2002. pp. 358-369.
- Samuel R. Madden, Michael J. Franklin, Joseph M. Hellerstein, and Wei Hong. *The Design of an Acquisitional Query Processor for Sensor Networks*. SIGMOD, June 2003, San Diego, CA.

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          101

---

**CMU SCS**

# References

- Lawrence Saul & Sam Roweis. *An Introduction to Locally Linear Embedding* (draft)
- Sam Roweis & Lawrence Saul. *Nonlinear dimensionality reduction by locally linear embedding*. Science, v.290 no.5500 , Dec.22, 2000. pp.2323--2326.
- B. Shaw and T. Jebara. "*Minimum Volume Embedding*" . Artificial Intelligence and Statistics, AISTATS, March 2007.

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          #102

**CMU SCS**

## References

- Josh Tenenbaum, Vin de Silva and John Langford. *A Global Geometric Framework for Nonlinear dimensionality Reduction*. Science 290, pp. 2319-2323, 2000.

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     #103

---

**CMU SCS**

# Part 2: DSP (Digital Signal Processing)

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     104

---

**CMU SCS**

## Outline

- Motivation
- Similarity Search and Indexing
- ➡ DSP (DFT, DWT)
- Linear Forecasting
- Kalman filters
- Bursty traffic - fractals and multifractals
- Non-linear forecasting
- Conclusions

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     105

---

**CMU SCS**

## Outline

- ➡ DFT
  - Definition of DFT and properties
  - how to read the DFT spectrum
- DWT
  - Definition of DWT and properties
  - how to read the DWT scalogram
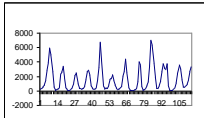
KDD 2010     (c) 2010, C. Faloutsos, Lei Li     106

---

**CMU SCS**

## Introduction - Problem#1

Goal: given a signal (eg., packets over time)

Find: patterns and/or compress

count



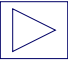lynx caught per year (packets per day; automobiles per hour)

year

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     107

---

**CMU SCS**

## What does DFT do?

A: highlights the periodicities

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     108

**CMU SCS** | Skip

# DFT: definition

- For a sequence $x_0, x_1, \dots x_{n-1}$
- the (**n-point**) Discrete Fourier Transform is
- $X_0, X_1, \dots X_{n-1}$ :

$$X_f = 1/\sqrt{n} \sum_{t=0}^{n-1} x_t * \exp(-j2\pi tf/n) \qquad f=0,\dots,n-1$$

$(j=\sqrt{-1})$

inverse DFT

$$x_t = 1/\sqrt{n} \sum_{t=0}^{n-1} X_f * \exp(+j2\pi tf/n)$$

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 109

**CMU SCS**

# DFT: definition

- **Good** news: Available in **all** symbolic math packages, eg., in 'mathematica'

  x = [1,2,1,2];

  X = Fourier[x];

  Plot[ Abs[X] ];

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 110

**CMU SCS**

# DFT: Amplitude spectrum

Amplitude: $A_f^2 = \text{Re}^2(X_f) + \text{Im}^2(X_f)$

count

Ampl.

freq=0

freq=12

year

Freq.

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 111

**CMU SCS** | Skip

# DFT: examples

flat

Amplitude

time

freq

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 112

**CMU SCS** | Skip

# DFT: examples

Low frequency sinusoid

time

freq

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 113

**CMU SCS** | Skip

# DFT: examples

- Sinusoid - symmetry property: $X_f = X^*_{n-f}$

time

freq

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 114

19

7/23/2010

**CMU SCS**

## DFT: examples

- Higher freq. sinusoid



time

freq

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 115

Skip

**CMU SCS**

## DFT: examples

examples



KDD 2010 (c) 2010, C. Faloutsos, Lei Li 116

Skip

**CMU SCS**

## DFT: examples

examples

Ampl.

Freq.



KDD 2010 (c) 2010, C. Faloutsos, Lei Li 117

Skip

**CMU SCS**

## Outline

- Motivation
- Similarity Search and Indexing
- DSP
- Linear Forecasting
- Bursty traffic - fractals and multifractals
- Non-linear forecasting
- Conclusions

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 118

**CMU SCS**

## Outline

- Motivation
- Similarity Search and Indexing
- DSP
  - DFT
    - Definition of DFT and properties
    - how to read the DFT spectrum
  - DWT

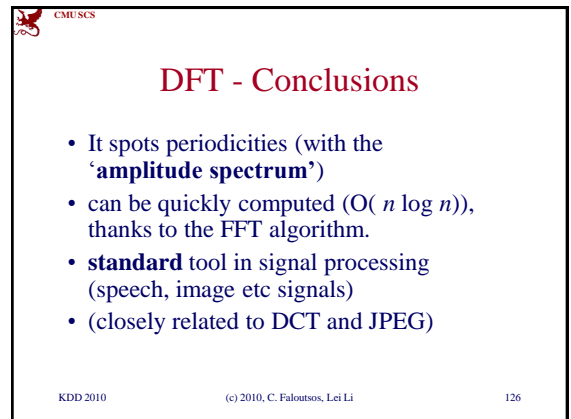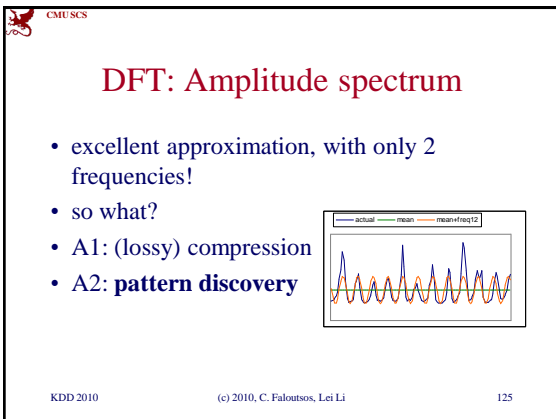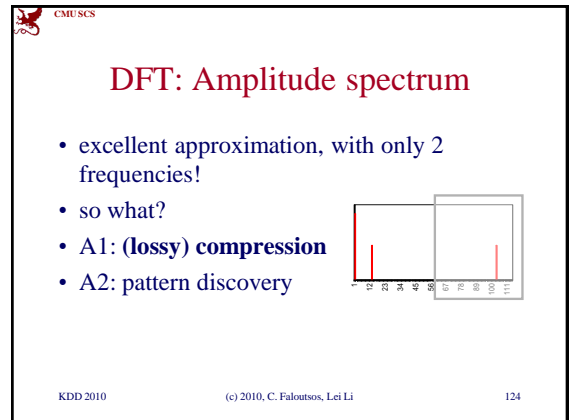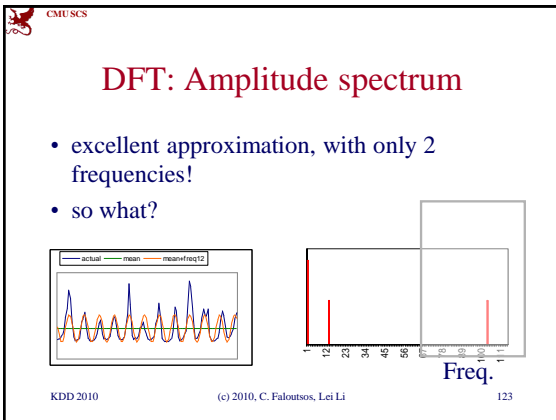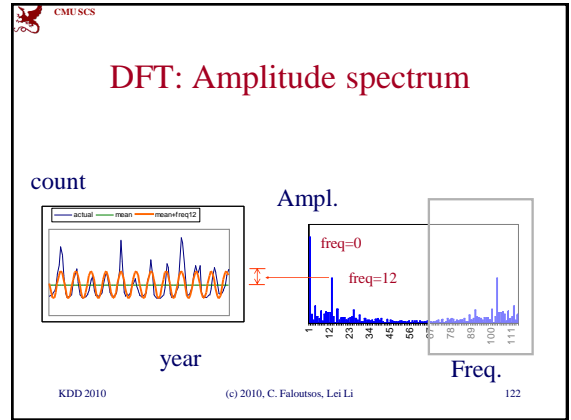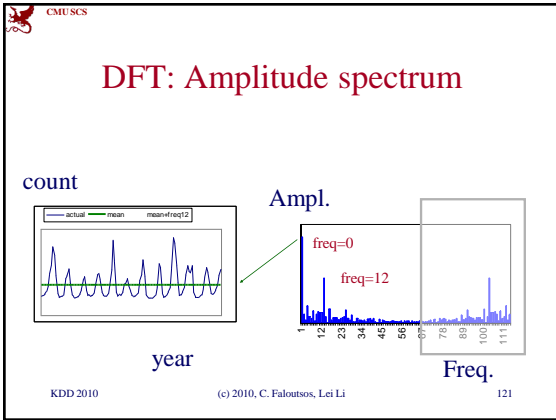KDD 2010 (c) 2010, C. Faloutsos, Lei Li 119

**CMU SCS**

## DFT: Amplitude spectrum

Amplitude: $A_f^2 = \text{Re}^2(X_f) + \text{Im}^2(X_f)$

count

Ampl.

freq=0

freq=12



year

Freq.

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 120

20

**CMU SCS**

# DFT: Amplitude spectrum

count

Ampl.

freq=0

freq=12

year

Freq.

---

**CMU SCS**

# DFT: Amplitude spectrum

count

Ampl.

freq=0

freq=12

year

Freq.

---

**CMU SCS**

# DFT: Amplitude spectrum

- excellent approximation, with only 2 frequencies!
- so what?

Freq.

---

**CMU SCS**

# DFT: Amplitude spectrum

- excellent approximation, with only 2 frequencies!
- so what?
- A1: **(lossy) compression**
- A2: pattern discovery

---

**CMU SCS**

# DFT: Amplitude spectrum

- excellent approximation, with only 2 frequencies!
- so what?
- A1: (lossy) compression
- A2: **pattern discovery**

---

**CMU SCS**

# DFT - Conclusions

- It spots periodicities (with the '**amplitude spectrum'**)
- can be quickly computed (O( $n \log n$)), thanks to the FFT algorithm.
- **standard** tool in signal processing (speech, image etc signals)
- (closely related to DCT and JPEG)

## Slide 127

**CMU SCS**

# Outline

- Motivation
- Similarity Search and Indexing
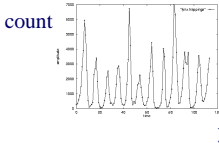- DSP
  - DFT
  - DWT
    - Definition of DWT and properties
    - how to read the DWT scalogram

KDD 2010       (c) 2010, C. Faloutsos, Lei Li       127

## Slide 128

**CMU SCS**

# Problem #1:

Goal: given a signal (eg., #packets over time)
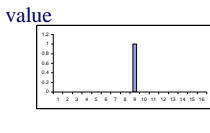Find: patterns, periodicities, and/or **compress**

count



lynx caught per year
(packets per day;
virus infections per month)

year

KDD 2010       (c) 2010, C. Faloutsos, Lei Li       128

## Slide 129

**CMU SCS**

# Wavelets - DWT

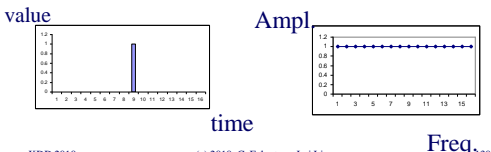- DFT is great - but, how about compressing a spike?

value



time

KDD 2010       (c) 2010, C. Faloutsos, Lei Li       129

## Slide 130

**CMU SCS**

# Wavelets - DWT

- DFT is great - but, how about compressing a spike?
- A: Terrible - all DFT coefficients needed!

value        Ampl



time

KDD 2010       (c) 2010, C. Faloutsos, Lei Li      Freq$_{130}$

## Slide 131

**CMU SCS**

# Wavelets - DWT

- DFT is great - but, how about compressing a spike?
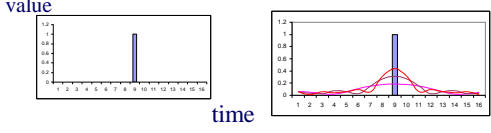- A: Terrible - all DFT coefficients needed!

value



time

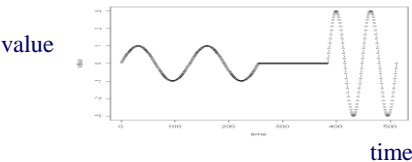KDD 2010       (c) 2010, C. Faloutsos, Lei Li       131

## Slide 132

**CMU SCS**

# Wavelets - DWT

- Similarly, DFT suffers on short-duration waves (eg., baritone, silence, soprano)

value



time

KDD 2010       (c) 2010, C. Faloutsos, Lei Li       132

**CMU SCS**

# Wavelets - DWT

- Solution#1: Short window Fourier transform (SWFT)
- But: how short should be the window?

freq

value

ime

time

**CMU SCS**

# Wavelets - DWT

- Answer: **multiple** window sizes! -> DWT

Time domain    DFT    SWFT    DWT

freq

time

**CMU SCS**

# Haar Wavelets

- subtract sum of left half from right half
- repeat recursively for quarters, eight-ths, ...

**CMU SCS**

Skip

# Wavelets - construction

x0  x1  x2  x3  x4  x5  x6  x7

**CMU SCS**

Skip

# Wavelets - construction

level 1   d1,0   s1,0   d1,1   s1,1   .......

\+

\-

x0  x1  x2  x3  x4  x5  x6  x7

**CMU SCS**

Skip

# Wavelets - construction

level 2   d2,0    s2,0

d1,0   s1,0   d1,1   s1,1   .......

\+

\-

x0  x1  x2  x3  x4  x5  x6  x7

7/23/2010

**Slide 139**

CMU SCS

Skip

# Wavelets - construction

etc ...

d2,0    s2,0

d1,0    s1,0  d1,1  s1,1    .......
        +
     -
        x0  x1  x2  x3  x4  x5  x6  x7

KDD 2010        (c) 2010, C. Faloutsos, Lei Li        139

---

**Slide 140**

CMU SCS

Skip

# Wavelets - construction

Q: map each coefficient

on the time-freq. plane

f

d2,0    s2,0

t

d1,0    s1,0  d1,1  s1,1    .......
        +
     -
        x0  x1  x2  x3  x4  x5  x6  x7

KDD 2010        (c) 2010, C. Faloutsos, Lei Li        140

---

**Slide 141**

CMU SCS

Skip

# Wavelets - construction

Q: map each coefficient

on the time-freq. plane

f

d2,0    s2,0

t

d1,0    s1,0  d1,1  s1,1    .......
        +
     -
        x0  x1  x2  x3  x4  x5  x6  x7

KDD 2010        (c) 2010, C. Faloutsos, Lei Li        141

---

**Slide 142**

CMU SCS

# Haar wavelets - code

```
#!/usr/bin/perl5
# expects a file with numbers
# and prints the dwt transform
# The number of time-ticks should be a power of 2
# USAGE
#   haar.pl <fname>

my @vals=();
my @smooth; # the smooth component of the signal
my @diff;  # the high-freq. component

# collect the values into the array @val
while(<>){
    @vals = ( @vals , split );
}
```

```
my $len = scalar(@vals);
my $half = int($len/2);
while($half >= 1 ){
  for(my $i=0; $i< $half; $i++){
      $diff{$i} = ($vals{2*$i} - $vals{2*$i + 1} )/ sqrt(2);
      print "\t", $diff[$i];
      $smooth {$i} = ($vals{2*$i} + $vals{2*$i + 1} )/ sqrt(2);
  }
  print "\n";
  @vals = @smooth;
  $half = int($half/2);
}
print "\t", $vals[0], "\n" ;    # the final, smooth component
```

KDD 2010        (c) 2010, C. Faloutsos, Lei Li        142

---

**Slide 143**

CMU SCS

# Wavelets - construction

Observation1:

'+' can be some weighted addition

'-' is the corresponding weighted difference
('Quadrature mirror filters')

Observation2: unlike DFT/DCT,

there are *many* wavelet bases: Haar, Daubechies-4, Daubechies-6, Coifman, Morlet, Gabor, ...

KDD 2010        (c) 2010, C. Faloutsos, Lei Li        143

---

**Slide 144**

CMU SCS

# Wavelets - how do they look like?

• E.g., Daubechies-4

KDD 2010        (c) 2010, C. Faloutsos, Lei Li        144

24

**CMU SCS**

# Wavelets - how do they look like?



- E.g., Daubechies-4

?

?

---

**CMU SCS**

# Wavelets - how do they look like?



- E.g., Daubechies-4

---

**CMU SCS**

# Outline

- Motivation
- Similarity Search and Indexing
- DSP
  - DFT
  - DWT
    - Definition of DWT and properties
    - how to read the DWT scalogram

---

**CMU SCS**

# Wavelets - Drill#1:

- Q: baritone/silence/soprano - DWT?

f

t

value

time

---

**CMU SCS**

# Wavelets - Drill#1:

- Q: baritone/silence/soprano - DWT?

f

t

value

time

---

**CMU SCS**

# Wavelets - Drill#2:

- Q: spike - DWT?

f

t

**Wavelets - Drill#2:**
• Q: spike - DWT?

0.00  0.00  **0.71**  0.00
0.00  **0.50**
**-0.35**
**0.35**

KDD 2010 — (c) 2010, C. Faloutsos, Lei Li — 151



**Wavelets - Drill#3:**
• Q: weekly + daily periodicity, + spike - DWT?

KDD 2010 — (c) 2010, C. Faloutsos, Lei Li — 152



**Wavelets - Drill#3:**
• Q: **weekly** + daily periodicity, + spike - DWT?

KDD 2010 — (c) 2010, C. Faloutsos, Lei Li — 153



**Wavelets - Drill#3:**
• Q: weekly + **daily** periodicity, + spike - DWT?

KDD 2010 — (c) 2010, C. Faloutsos, Lei Li — 154



**Wavelets - Drill#3:**
• Q: weekly + daily periodicity, + **spike** - DWT?

KDD 2010 — (c) 2010, C. Faloutsos, Lei Li — 155



**Wavelets - Drill#3:**
• Q: weekly + daily periodicity, + spike - DWT?

KDD 2010 — (c) 2010, C. Faloutsos, Lei Li — 156

26

**CMU SCS**

# Wavelets - Drill#3:

• Q: DFT?

DWT    DFT

f    f

t    t

---

**CMU SCS**

# Advantages of Wavelets

• Better compression (better RMSE with same number of coefficients - used in JPEG-2000)
• fast to compute (usually: O(*n*)!)
• very good for 'spikes'
• mammalian eye and ear: Gabor wavelets

---

**CMU SCS**

# Overall Conclusions

• DFT, DCT spot periodicities
• **DWT** : multi-resolution - matches processing of mammalian ear/eye better
• All three: powerful tools for **compression**, **pattern detection** in real signals
• All three: included in math packages
  – (matlab, 'R', mathematica, … - often in spreadsheets!)

---

**CMU SCS**

# Overall Conclusions

• DWT : very suitable for self-similar traffic
• DWT: used for summarization of streams [Gilbert+01], db histograms etc

---

**CMU SCS**

# Resources - software and urls

• http://www.dsptutor.freeuk.com/jsanalyser/ FFTSpectrumAnalyser.html : Nice java applets for FFT
• http://www.relisoft.com/freeware/freq.html voice frequency analyzer (needs microphone)

---

**CMU SCS**

# Resources: software and urls

• *xwpl:* open source wavelet package from Yale, with excellent GUI
• http://monet.me.ic.ac.uk/people/gavin/java /waveletDemos.html : wavelets and scalograms

## Books

- William H. Press, Saul A. Teukolsky, William T. Vetterling and Brian P. Flannery: *Numerical Recipes in C*, Cambridge University Press, 1992, 2nd Edition. (Great description, intuition and code for DFT, DWT)
- C. Faloutsos: *Searching Multimedia Databases by Content*, Kluwer Academic Press, 1996 (introduction to DFT, DWT)

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      163

## Additional Reading

- [Gilbert+01] Anna C. Gilbert, Yannis Kotidis and S. Muthukrishnan and Martin Strauss, *Surfing Wavelets on Streams: One-Pass Summaries for Approximate Aggregate Queries*, VLDB 2001

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      164

# Part 3: Linear Forecasting

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      165

## Outline

- Motivation
- Similarity Search and Indexing
- DSP
➡ - Linear Forecasting
- Kalman filters
- Bursty traffic - fractals and multifractals
- Non-linear forecasting
- Conclusions

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      166

## Forecasting

"Prediction is very difficult, especially about the future." - Nils Bohr

`http://www.hfac.uh.edu/MediaFutures/thoughts.html`

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      167

## Outline

- Motivation
- ...
- Linear Forecasting
  ➡ – Auto-regression: Least Squares; RLS
  – Co-evolving time sequences
  – Examples
  – Conclusions

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      168

## Slide 169

**CMU SCS**

# Problem#2: Forecast

- Example: give $x_{t-1}$, $x_{t-2}$, …, forecast $x_t$



KDD 2010      (c) 2010, C. Faloutsos, Lei Li      169

## Slide 170

**CMU SCS**

# Forecasting: Preprocessing

MANUALLY:

remove trends      spot periodicities

7 days



time      time

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      170

## Slide 171

**CMU SCS**

# Problem#2: Forecast

- Solution: try to express

  $x_t$

  as a linear function of the past: $x_{t-2}$, $x_{t-2}$, …,

  (up to a window of $w$)

Formally:

$$x_t \approx a_1 x_{t-1} + \ldots + a_w x_{t-w} + noise$$

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      171

## Slide 172

**CMU SCS**

# (Problem: Back-cast; interpolate)

- Solution - interpolate: try to express

  $x_t$

  as a linear function of the past AND the future:

  $x_{t+1}$, $x_{t+2}$, … $x_{t+wfuture;}$ $x_{t-1}$, … $x_{t-wpast}$

  (up to windows of $w_{past}$, $w_{future}$)

- EXACTLY the same algo's

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      172

## Slide 173

**CMU SCS**

# Linear Regression: idea

| patient | weight | height |
|---------|--------|--------|
| 1 | 27 | 43 |
| 2 | 43 | 54 |
| 3 | 54 | 72 |
| … | … | … |
| N | 25 | ?? |



- express what we don't know (= 'dependent variable')
- as a linear function of what we know (= 'indep. variable(s)')

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      173

## Slide 174

**CMU SCS**

# Linear <u>Auto</u> Regression:

| Time | Packets Sent(t) |
|------|-----------------|
| 1 | 43 |
| 2 | 54 |
| 3 | 72 |
| … | … |
| N | ?? |

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      174

**Slide 175**

CMU SCS

## Linear <u>Auto</u> Regression:

| Time | Packets Sent (t-1) | Packets Sent(t) |
|---|---|---|
| 1 | - | 43 |
| 2 | 43 | 54 |
| 3 | 54 | 72 |
| … | … | … |
| N | 25 | ?? |

'lag-plot'

Number of packets sent (t) vs Number of packets sent (t-1)

- lag $w$=1
- <u>Dependent</u> variable = # of packets sent (S [t])
- <u>Independent</u> variable = # of packets sent (S[t-1])

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          175

**Slide 176**

CMU SCS

## Outline

- Motivation
- ...
- Linear Forecasting
  → – Auto-regression: **Least Squares; RLS**
  – Co-evolving time sequences
  – Examples
  – Conclusions

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          176

**Slide 177**

CMU SCS

## More details:

- Q1: Can it work with window $w$>1?
- A1: YES!

$x_t$, $x_{t-1}$, $x_{t-2}$

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          177

**Slide 178**

CMU SCS

## More details:

- Q1: Can it work with window $w$>1?
- A1: YES! (we'll fit a hyper-plane, then!)

$x_t$, $x_{t-1}$, $x_{t-2}$

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          178

**Slide 179**

CMU SCS

## More details:

- Q1: Can it work with window $w$>1?
- A1: YES! (we'll fit a hyper-plane, then!)

$x_t$, $x_{t-1}$, $x_{t-2}$

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          179

**Slide 180**

CMU SCS

Skip

## More details:

- Q1: Can it work with window $w$>1?
- A1: YES! The problem becomes:

$$\mathbf{X}_{[N \times w]} \times \mathbf{a}_{[w \times 1]} = \mathbf{y}_{[N \times 1]}$$

- OVER-CONSTRAINED
  – **a** is the vector of the regression coefficients
  – **X** has the $N$ values of the $w$ indep. variables
  – **y** has the N values of the dependent variable

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          180

30

**CMU SCS**

More details:

$$\mathbf{X}_{[N \times w]} \times \mathbf{a}_{[w \times 1]} = \mathbf{y}_{[N \times 1]}$$

Ind-var1      Ind-var-w

time $\begin{bmatrix} X_{11}, X_{12}, \cdots, X_{1w} \\ X_{21}, X_{22}, \ldots, X_{2w} \\ \vdots \\ \vdots \\ \vdots \\ X_{N1}, X_{N2}, \ldots, X_{Nw} \end{bmatrix} \times \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_w \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \\ \vdots \\ y_N \end{bmatrix}$

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          181

---

**CMU SCS**

More details:

$$\mathbf{X}_{[N \times w]} \times \mathbf{a}_{[w \times 1]} = \mathbf{y}_{[N \times 1]}$$

Ind-var1      Ind-var-w

time $\begin{bmatrix} X_{11}, X_{12}, \cdots, X_{1w} \\ X_{21}, X_{22}, \ldots, X_{2w} \\ \vdots \\ \vdots \\ \vdots \\ X_{N1}, X_{N2}, \ldots, X_{Nw} \end{bmatrix} \times \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_w \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \\ \vdots \\ y_N \end{bmatrix}$

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          182

---

**CMU SCS**

More details

- Q2: How to estimate $a_1, a_2, \ldots a_w = \mathbf{a}$?
- A2: with Least Squares fit

$$\mathbf{a} = ( \mathbf{X}^T \times \mathbf{X} )^{-1} \times (\mathbf{X}^T \times \mathbf{y})$$

- (Moore-Penrose pseudo-inverse)
- $\mathbf{a}$ is the vector that minimizes the RMSE from $\mathbf{y}$

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          183

---

**CMU SCS**

Even more details

- Q3: Can we estimate $\mathbf{a}$ incrementally?
- A3: Yes, with the brilliant, classic method of 'Recursive Least Squares' (RLS) (see, e.g., [Yi+00], for details) - pictorially:

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          184

---

**CMU SCS**

Even more details

- Given:



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          185

---

**CMU SCS**

Even more details



← new point

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          186

**CMU SCS**

## Even more details

### RLS: quickly compute new best fit



← new point

Dependent Variable

Independent Variable

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     187

---

**CMU SCS**

## Even more details

- Straightforward Least Squares
  - Needs huge matrix (**growing** in size) $O(N \times w)$
  - Costly matrix operation $O(N \times w^2)$

- Recursive LS
  - Need much smaller, fixed size matrix $O(w \times w)$
  - Fast, incremental computation $O(1 \times w^2)$

$N = 10^6, \quad w = 1\text{-}100$

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     188

---

**CMU SCS**

## Even more details

- Q4: can we 'forget' the older samples?
- A4: Yes - RLS can easily handle that [Yi+00]:

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     189

---

**CMU SCS**

## Adaptability - 'forgetting'



Dependent Variable eg., #bytes sent

Independent Variable eg., #packets sent

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     190

---

**CMU SCS**

## Adaptability - 'forgetting'



Trend change

(R)LS with no forgetting

Dependent Variable eg., #bytes sent

Independent Variable eg. #packets sent

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     191

---

**CMU SCS**

## Adaptability - 'forgetting'



Trend change

(R)LS with no forgetting

(R)LS **with** forgetting

Dependent Variable

Independent Variable

- RLS: can *trivially* handle 'forgetting'

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     192

## How to choose 'w'?

- goal: capture arbitrary periodicities
- with NO human intervention
- on a semi-infinite stream

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          193

## Answer:

- 'AWSOM' (Arbitrary Window Stream fOrecasting Method) [Papadimitriou+, vldb2003]
- idea: do AR on each wavelet level
- in detail:

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          194

## AWSOM



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          195

## AWSOM



KDD 2010          (c) 2010, C. Faloutsos, Lei Li          196

## AWSOM - idea



$$W_{l,t} = \beta_{l,1}W_{l,t-1} + \beta_{l,2}W_{l,t-2} + \dots$$

$$W_{l',t'} = \beta_{l',1}W_{l',t'-1} + \beta_{l',2}W_{l',t'-2} + \dots$$

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          197

## More details…

- Update of wavelet coefficients   (incremental)
- Update of linear models   (incremental; RLS)
- Feature selection          (single-pass)
  - Not all correlations are significant
  - Throw away the insignificant ones ("noise")

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          198

33

## Results - Synthetic data

AWSOM     AR     Seasonal AR



- Triangle pulse
- Mix (sine + square)
- AR captures wrong trend (or none)
- Seasonal AR estimation fails

CMU SCS

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     199

---

## Results - Real data



- Automobile traffic
  - Daily periodicity
  - Bursty "noise" at smaller scales
- AR fails to capture any trend
- Seasonal AR estimation fails

CMU SCS

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     200

---

## Results - real data



- Sunspot intensity
  - Slightly time-varying "period"
- AR captures wrong trend
- Seasonal ARIMA
  - wrong downward trend, despite help by human!

CMU SCS

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     201

---

## Complexity

Skip

- Model update
  
  Space: $O(lgN + mk^2) \approx O(lgN)$
  
  Time: $O(k^2) \approx O(1)$
- Where
  - $N$: number of points (so far)
  - $k$: number of regression coefficients; fixed
  - $m$: number of linear models; $O(lgN)$

CMU SCS

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     202

---

## Outline

- Motivation
- ...
- Linear Forecasting
  - Auto-regression: Least Squares; RLS
  - Co-evolving time sequences
  - Examples
  - Conclusions

CMU SCS

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     203

---

## Co-Evolving Time Sequences

- Given: A set of **correlated** time sequences
- Forecast '**Repeated(t)**'



CMU SCS

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     204

**CMU SCS**

# Solution:

Q: what should we do?

---

**CMU SCS**

# Solution:

Least Squares, with
- Dep. Variable: Repeated(t)
- Indep. Variables: Sent(t-1) … Sent(t-w); Lost(t-1) …Lost(t-w); Repeated(t-1), ...
- (named: 'MUSCLES' [Yi+00])

---

**CMU SCS**

# Time Series Analysis - Outline

- Auto-regression
- Least Squares; recursive least squares
- Co-evolving time sequences
- Examples
- ➡ Conclusions

---

**CMU SCS**

# Conclusions - Practitioner's guide

- AR(IMA) methodology: prevailing method for linear forecasting
- Brilliant method of Recursive Least Squares for fast, incremental estimation.
- See [Box-Jenkins]
- very recently: AWSOM (no human intervention)

---

**CMU SCS**

# Resources: software and urls

- MUSCLES: Prof. Byoung-Kee Yi:
  `http://www.postech.ac.kr/~bkyi/`
  or `christos@cs.cmu.edu`
- free-ware: 'R' for stat. analysis
  (clone of Splus)
  `http://cran.r-project.org/`

---

**CMU SCS**

# Books

- George E.P. Box and Gwilym M. Jenkins and Gregory C. Reinsel, *Time Series Analysis: Forecasting and Control*, Prentice Hall, 1994 (the classic book on ARIMA, 3rd ed.)
- Brockwell, P. J. and R. A. Davis (1987). Time Series: Theory and Methods. New York, Springer Verlag.

---

**CMU SCS**

## Additional Reading

- [Papadimitriou+ vldb2003] Spiros Papadimitriou, Anthony Brockwell and Christos Faloutsos *Adaptive, Hands-Off Stream Mining* VLDB 2003, Berlin, Germany, Sept. 2003
- [Yi+00] Byoung-Kee Yi et al.: *Online Data Mining for Co-Evolving Time Sequences*, ICDE 2000. (Describes MUSCLES and Recursive Least Squares)

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     211

---

**CMU SCS**

# BREAK!

## Next: Kalman filters

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     212

---

**CMU SCS**

# Part 5: Bursty traffic and multifractals

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     213

---

**CMU SCS**

## Outline

- Motivation
- Similarity Search and Indexing
- DSP
- Linear Forecasting
- Kalman filters
- → Bursty traffic - fractals and multifractals
- Non-linear forecasting
- Conclusions

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     214

---

**CMU SCS**

## Outline

- Motivation
- ...
- Linear Forecasting
- Bursty traffic - fractals and multifractals
  - → Problem
  - Main idea (80/20, Hurst exponent)
  - Results

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     215

---

**CMU SCS**

## Recall: Problem #1:

Goal: given a signal (eg., #bytes over time)
Find: patterns, periodicities, and/or compress

#bytes  Bytes per 30'
(packets per day; earthquakes per year)

time

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     216

---

36

**Problem #1**

- model bursty traffic
- generate realistic traces
- (Poisson does not work)

# bytes

Poisson

time

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 217

**Motivation**

- predict queue length distributions (e.g., to give probabilistic guarantees)
- "learn" traffic, for buffering, prefetching, 'active disks', web servers

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 218

**Q: any 'pattern'?**

- Not Poisson
- spike; silence; more spikes; more silence…
- any rules?

# bytes

time

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 219

**solution: self-similarity**

# bytes       # bytes

time       time

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 220

**solution: self-similarity**

# bytes       # bytes

time       time

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 221

**solution: self-similarity**

# bytes       # bytes

time       time

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 222

37

**CMU SCS**

## solution: self-similarity

**# bytes**

**# bytes**



**time**

**time**

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          223

---

**CMU SCS**

## But:

- Q1: How to generate realistic traces; extrapolate; give guarantees?
- Q2: How to estimate the model parameters?

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          224

---

**CMU SCS**

## Outline

- Motivation
- ...
- Linear Forecasting
- Bursty traffic - fractals and multifractals
  – Problem
  ➡ – Main idea (80/20, Hurst exponent)
  – Results

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          225

---

**CMU SCS**

## Approach

- Q1: How to generate a sequence, that is
  – bursty
  – self-similar
  – and has similar queue length distributions

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          226

---

**CMU SCS**

## Approach

- A: 'binomial multifractal' [Wang+02]
- ~ 80-20 'law':
  – 80% of bytes/queries etc on first half
  – repeat recursively
- *b*: bias factor (eg., 80%)

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          227

---

**CMU SCS**

## binary multifractals

**20** ∧ **80**



KD|          228

---

38

## Slide 229
**CMU SCS**

# binary multifractals

**20** ∧ **80**

∧ ∧



KD|     229

## Slide 230
**CMU SCS**

# Parameter estimation

- Q2: How to estimate the bias factor *b*?

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    230

## Slide 231
**CMU SCS**

# Parameter estimation

- Q2: How to estimate the bias factor *b*?
- A: MANY ways [Crovella+96]
  - Hurst exponent
  - variance plot
  - even DFT amplitude spectrum! ('periodogram')
  - More robust: 'entropy plot' [Wang+02]

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    231

## Slide 232
**CMU SCS**

# Entropy plot

- Rationale:
  - burstiness: inverse of uniformity
  - entropy measures uniformity of a distribution
  - find entropy at several granularities, to see whether/how our distribution is close to uniform.

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    232

## Slide 233
**CMU SCS**

# Entropy plot

p1    p2
% of bytes here

- Entropy *E(n)* after *n* levels of splits
- n=1: $E(1)= - p1 \log_2(p1) - p2 \log_2(p2)$

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    233

## Slide 234
**CMU SCS**

# Entropy plot

$p_{2,1}$   $p_{2,2}$   $p_{2,3}$   $p_{2,4}$

- Entropy *E(n)* after *n* levels of splits
- n=1: $E(1)= - p1 \log(p1) - p2 \log(p2)$
- n=2: $E(2) = - \Sigma_t p_{2,i} * \log_2 (p_{2,i})$

KDD 2010    (c) 2010, C. Faloutsos, Lei Li    234

39

## Slide 235

**CMU SCS**

# Real traffic

Entropy
*E(n)*

0.73

• Has linear entropy plot
  (-> self-similar)

# of levels (*n*)

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          235

## Slide 236

**CMU SCS**      **Skip**

# Observation - intuition:

Entropy
*E(n)*

intuition: slope =
  intrinsic dimensionality =
  info-bits per coordinate-bit
  – unif. Dataset: slope =1
  – multi-point: slope = 0

# of levels (*n*)

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          236

## Slide 237

**CMU SCS**      **Skip**

# Entropy plot - Intuition

• Slope ~ intrinsic dimensionality (in fact, 'Information fractal dimension')
• = info bit per coordinate bit - eg

Dim = 1

Pick a point;
reveal its coordinate bit-by-bit -
how much info is each bit worth to me?

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          237

## Slide 238

**CMU SCS**      **Skip**

# Entropy plot

• Slope ~ intrinsic dimensionality (in fact, 'Information fractal dimension')
• = info bit per coordinate bit - eg

Dim = 1

Is MSB 0?

'info' value = E(1): 1 bit

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          238

## Slide 239

**CMU SCS**      **Skip**

# Entropy plot

• Slope ~ intrinsic dimensionality (in fact, 'Information fractal dimension')
• = info bit per coordinate bit - eg

Dim = 1

Is MSB 0?

Is next MSB =0?

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          239

## Slide 240

**CMU SCS**      **Skip**

# Entropy plot

• Slope ~ intrinsic dimensionality (in fact, 'Information fractal dimension')
• = info bit per coordinate bit - eg

Dim = 1

Info value =1 bit
= E(2) - E(1) =
slope!

Is MSB 0?

Is next MSB =0?

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          240

**CMU SCS**

## Entropy plot

- Repeat, for all points at same position:

Dim=0

---

**CMU SCS**

## Entropy plot

- Repeat, for all points at same position:
- we need 0 bits of info, to determine position
- -> slope = 0 = intrinsic dimensionality

Dim=0

---

**CMU SCS**

## Entropy plot

- Real (and 80-20) datasets can be in-between: bursts, gaps, smaller bursts, smaller gaps, at every scale

Dim = 1

Dim=0

0<Dim<1

---

**CMU SCS**

## (Fractals)

- What set of points could have behavior between point and line?

---

**CMU SCS**

## Cantor dust

- Eliminate the middle third
- Recursively!

---

**CMU SCS**

## Cantor dust

**CMU SCS**

# Cantor dust

**CMU SCS**

# Cantor dust

**CMU SCS**

# Cantor dust

**CMU SCS**

# Cantor dust

Dimensionality?
(no length; infinite # points!)
Answer: log2 / log3 = 0.6

**CMU SCS**

# Some more entropy plots:

- **Poisson** vs real



Poisson: slope = ~1 -> uniformly distributed

**CMU SCS**

# B-model

$E(n)$



$n$

- b-model traffic gives perfectly linear plot
- Lemma: its slope is

  $slope = -b \, log_2 b - (1-b) \, log_2 (1-b)$

- Fitting: do entropy plot; get slope; solve for $b$

---

**CMU SCS**

# Outline

- Motivation
- ...
- Linear Forecasting
- Bursty traffic - fractals and multifractals
  - Problem
  - Main idea (80/20, Hurst exponent)
  ➡ - Experiments - Results

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     253

---

**CMU SCS**

# Experimental setup

- Disk traces (from HP [Wilkes 93])
- web traces from LBL
  `http://repository.cs.vt.edu/`
  `lbl-conn-7.tar.Z`

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     254

---

**CMU SCS**

# Model validation

- Linear entropy plots



Bias factors $b$: 0.6-0.8
smallest $b$ / smoothest: nntp traffic

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     255

---

**CMU SCS**

# Web traffic - results

- LBL, NCDF of queue lengths (log-log scales)

Prob( >$l$)



How to give guarantees?    (queue length $l$)

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     256

---

**CMU SCS**

# Web traffic - results

- LBL, NCDF of queue lengths (log-log scales)

Prob( >$l$)



20% of the requests will see
queue lengths <100

(queue length $l$)

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     257

---

**CMU SCS**

# Conclusions

- Multifractals (80/20, 'b-model', Multiplicative Wavelet Model (MWM)) for analysis and synthesis of bursty traffic
- can give (probabilistic) guarantees

KDD 2010     (c) 2010, C. Faloutsos, Lei Li     258

## Books

- Fractals: Manfred Schroeder: *Fractals, Chaos, Power Laws: Minutes from an Infinite Paradise* W.H. Freeman and Company, 1991 (Probably the BEST book on fractals!)

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      259

## Further reading:

- Crovella, M. and A. Bestavros (1996). Self-Similarity in World Wide Web Traffic, Evidence and Possible Causes. Sigmetrics.
- [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.

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      260

## Further reading

- [Riedi+99] R. H. Riedi, M. S. Crouse, V. J. Ribeiro, and R. G. Baraniuk, *A Multifractal Wavelet Model with Application to Network Traffic*, IEEE Special Issue on Information Theory, 45. (April 1999), 992-1018.
- [Wang+02] Mengzhi Wang, Tara Madhyastha, Ngai Hang Chang, Spiros Papadimitriou and Christos Faloutsos, *Data Mining Meets Performance Evaluation: Fast Algorithms for Modeling Bursty Traffic*, ICDE 2002, San Jose, CA, 2/26/2002 - 3/1/2002.

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      261

# Part 6: chaos and non-linear forecasting

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      262

## Outline

- Motivation
- Similarity Search and Indexing
- DSP
- Linear Forecasting
- Kalman filters
- Bursty traffic - fractals and multifractals
- Non-linear forecasting
- Conclusions

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      263

## Detailed Outline

- Non-linear forecasting
  - Problem
  - Idea
  - How-to
  - Experiments
  - Conclusions

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      264

**CMU SCS**

## Recall: Problem #1

Value



Time

Given a time series $\{x_t\}$, predict its future course, that is, $x_{t+1}$, $x_{t+2}$, ...

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 265

---

**CMU SCS**

## How to forecast?

- ARIMA - but: linearity assumption

- ANSWER: 'Delayed Coordinate Embedding' = Lag Plots [Sauer92]

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 266

---

**CMU SCS**

## General Intuition (Lag Plot)

**Lag = 1, k = 4 NN**

$x_t$

Interpolate these…

To get the final prediction
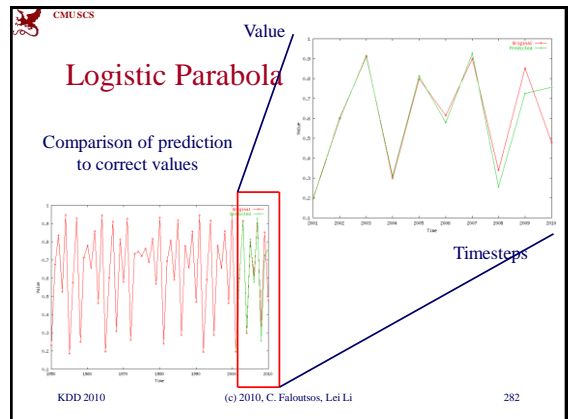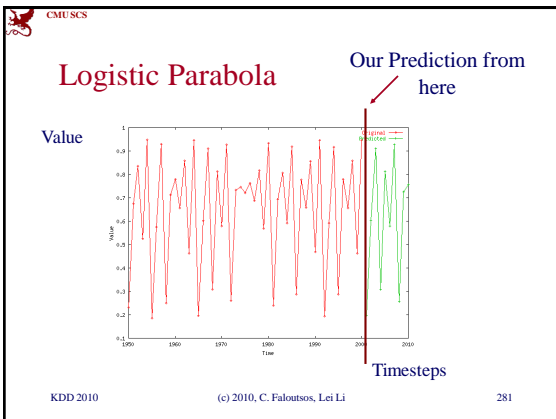
4-NN

New Point

$x_{t-1}$

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 267

---

**CMU SCS**

## Questions:

- Q1: How to choose lag $L$?
- Q2: How to choose $k$ (the # of NN)?
- Q3: How to interpolate?
- Q4: why should this work at all?

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 268

---

**CMU SCS**

## Q1: Choosing lag $L$

- Manually (16, in award winning system by [Sauer94])

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 269

---

**CMU SCS**

## Q2: Choosing number of neighbors $k$

- Manually (typically ~ 1-10)

KDD 2010 (c) 2010, C. Faloutsos, Lei Li 270

45

**CMU SCS**

## Q3: How to interpolate?

How do we interpolate between the *k* nearest neighbors?

A3.1: Average

A3.2: Weighted average (weights drop with distance - how?)

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          271

---

**CMU SCS**

## Q3: How to interpolate?

A3.3: Using SVD - seems to perform best ([Sauer94] - first place in the Santa Fe forecasting competition)



$x_t$

$X_{t-1}$

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          272

---

**CMU SCS**

## Q4: Any theory behind it?

A4: YES!

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          273

---

**CMU SCS**

## Theoretical foundation

- Based on the "Takens' Theorem" [Takens81]
- which says that <u>long enough</u> delay vectors can do prediction, even if there are unobserved variables in the dynamical system (= diff. equations)

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          274

---

**CMU SCS**

Skip

## Theoretical foundation

Example: Lotka-Volterra equations

$$dH/dt = r H - a H*P$$
$$dP/dt = b H*P - m P$$



P

H

H is count of prey (e.g., hare)
P is count of predators (e.g., lynx)

Suppose only P(t) is observed (t=1, 2, …).

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          275

---

**CMU SCS**

Skip

## Theoretical foundation

- But the delay vector space is a faithful reconstruction of the internal system state
- So prediction in **delay vector space** is as good as prediction in **state space**

P

P(t)



KDD 2010          H          (c) 2010, C. Faloutsos, Lei Li          P(t-1)          276

---

46

Solution to Volterra-Lotka eq.

# predators

prey

predators
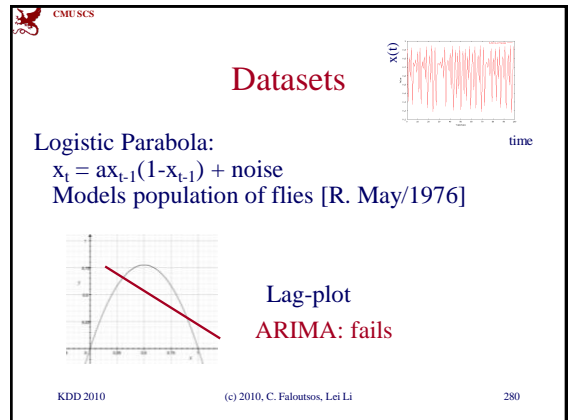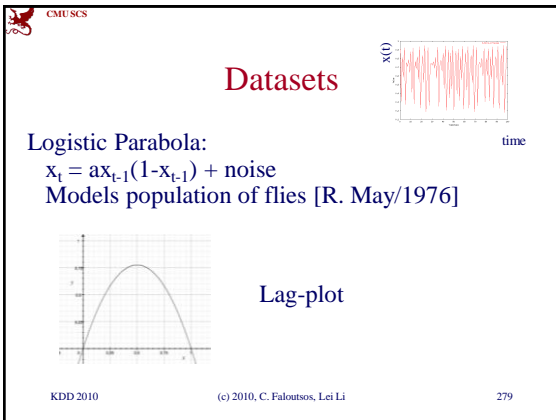
time

# prey

from wikipedia

KDD 2010 · 277



Detailed Outline

- Non-linear forecasting
  - Problem
  - Idea
  - How-to
  - Experiments
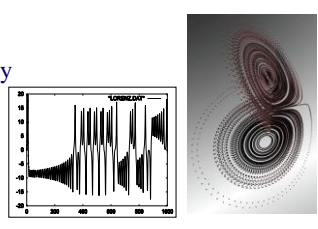  - Conclusions

KDD 2010 · (c) 2010, C. Faloutsos, Lei Li · 278



Datasets

$x(t)$

time

Logistic Parabola:
$x_t = ax_{t-1}(1-x_{t-1}) + noise$
Models population of flies [R. May/1976]

Lag-plot

KDD 2010 · (c) 2010, C. Faloutsos, Lei Li · 279



Datasets

$x(t)$

time

Logistic Parabola:
$x_t = ax_{t-1}(1-x_{t-1}) + noise$
Models population of flies [R. May/1976]

Lag-plot

ARIMA: fails

KDD 2010 · (c) 2010, C. Faloutsos, Lei Li · 280



Logistic Parabola

Our Prediction from here

Value

Timesteps

KDD 2010 · (c) 2010, C. Faloutsos, Lei Li · 281



Logistic Parabola

Value

Comparison of prediction to correct values

Timesteps

KDD 2010 · (c) 2010, C. Faloutsos, Lei Li · 282

47

## Datasets

**Skip**

LORENZ: Models convection
currents in the air
$$dx / dt = a (y - x)$$
$$dy / dt = x (b - z) - y$$
$$dz / dt = xy - c z$$

---

## LORENZ

**Skip**

Value

Comparison of prediction
to correct values

Timesteps

---

Value

## Datasets

**Skip**

- LASER: fluctuations in a Laser over time (used in Santa Fe competition)

Time

---

## Laser

**Skip**

Value

Comparison of prediction
to correct values

Timesteps

---

## Conclusions

- Lag plots for non-linear forecasting (Takens' theorem)
- suitable for 'chaotic' signals

---

## References

- Deepay Chakrabarti and Christos Faloutsos *F4: Large-Scale Automated Forecasting using Fractals* CIKM 2002, Washington DC, Nov. 2002.
- Sauer, T. (1994). *Time series prediction using delay coordinate embedding*. (in book by Weigend and Gershenfeld, below) Addison-Wesley.
- Takens, F. (1981). *Detecting strange attractors in fluid turbulence*. Dynamical Systems and Turbulence. Berlin: Springer-Verlag.

**CMU SCS**

# References

- Weigend, A. S. and N. A. Gerschenfeld (1994). *Time Series Prediction: Forecasting the Future and Understanding the Past*, Addison Wesley. (Excellent collection of papers on chaotic/non-linear forecasting, describing the algorithms behind the winners of the Santa Fe competition.)

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      289

---

**CMU SCS**

# Overall conclusions

- Similarity search: **Euclidean**/time-warping; **feature extraction** and **SAMs**

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      290

---

**CMU SCS**

# Overall conclusions

- Similarity search: **Euclidean**/time-warping; **feature extraction** and **SAMs**
- Signal processing: **DWT** is a powerful tool

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      291

---

**CMU SCS**

# Overall conclusions

- Similarity search: **Euclidean**/time-warping; **feature extraction** and **SAMs**
- Signal processing: **DWT** is a powerful tool
- Linear Forecasting: **AR** (Box-Jenkins)

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      292

---

**CMU SCS**

# Overall conclusions

- Similarity search: **Euclidean**/time-warping; **feature extraction** and **SAMs**
- Signal processing: **DWT** is a powerful tool
- Linear Forecasting: **AR** (Box-Jenkins)
- **Kalman** filters & extensions: forecasting, pattern discovery, segmentation

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      293

---

**CMU SCS**

# Overall conclusions

- Similarity search: **Euclidean**/time-warping; **feature extraction** and **SAMs**
- Signal processing: **DWT** is a powerful tool
- Linear Forecasting: **AR** (Box-Jenkins)
- **Kalman** filters & extensions: forecasting, pattern discovery, segmentation
- Bursty traffic: **multifractals** (80-20 'law')

KDD 2010      (c) 2010, C. Faloutsos, Lei Li      294

**CMU SCS**

## Overall conclusions

- Similarity search: **Euclidean**/time-warping; **feature extraction** and **SAMs**
- Signal processing: **DWT** is a powerful tool
- Linear Forecasting: **AR** (Box-Jenkins)
- **Kalman** filters & extensions: forecasting, pattern discovery, segmentation
- Bursty traffic: **multifractals** (80-20 'law')
- Non-linear forecasting: **lag-plots** (Takens)

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          295

**CMU SCS**

# THANK YOU!

christos AT cs.cmu.edu
www.cs.cmu.edu/~christos

leili AT cs.cmu.edu
www.cs.cmu.edu/~leili
www.cs.cmu.edu/~leili/pubs/dynammo.2.1.2.zip

KDD 2010          (c) 2010, C. Faloutsos, Lei Li          296