**CMU SCS**

# Sensor Data Mining: Similarity Search and Pattern Analysis

*Christos Faloutsos*
CMU

---

**CMU SCS**

# Thanks

Deepay Chakrabarti (CMU)

Prof. Dimitris Gunopulos (UCR)

Spiros Papadimitriou (CMU)

Mengzhi Wang (CMU)

Prof. Byoung-Kee Yi (Pohang U.)

CIKM 04     (c) C. Faloutsos, 2004     2

---

**CMU SCS**

# Outline

➡ • Motivation
• Similarity Search and Indexing
• DSP (Digital Signal Processing)
• Linear Forecasting
• Bursty traffic - fractals and multifractals
• Non-linear forecasting
• Conclusions

CIKM 04     (c) C. Faloutsos, 2004     3

---

**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

CIKM 04     (c) C. Faloutsos, 2004     4

---

**CMU SCS**

# Motivation - Applications

• Financial, sales, economic series
• Medical
  – ECGs +; blood pressure etc monitoring
  – reactions to new drugs
  – elderly care

CIKM 04     (c) C. Faloutsos, 2004     5

---

**CMU SCS**

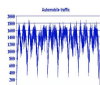# Motivation - Applications (cont'd)

• 'Smart house'
  – sensors monitor temperature, humidity, air quality
• video surveillance

CIKM 04     (c) C. Faloutsos, 2004     6

**CMU SCS**

## Motivation - Applications (cont'd)

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



CIKM 04  (c) C. Faloutsos, 2004  7

**CMU SCS**

## Motivation - Applications (cont'd)

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



CIKM 04  (c) C. Faloutsos, 2004  8

**CMU SCS**

## Motivation - Applications (cont'd)

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

CIKM 04  (c) C. Faloutsos, 2004  9

**CMU SCS**

## Stream Data: Disk accesses

#bytes



time

CIKM 04  (c) C. Faloutsos, 2004  10

**CMU SCS**

## Settings & Applications

- One or more sensors, collecting time-series data

CIKM 04  (c) C. Faloutsos, 2004  11

**CMU SCS**

## Settings & Applications



Each sensor collects data $(x_1, x_2, ..., x_t, ...)$

CIKM 04  (c) C. Faloutsos, 2004  12

**CMU SCS**

## Settings & Applications

Some sensors 'report' to others or to the central site

CIKM 04                    (c) C. Faloutsos, 2004                    13

---

**CMU SCS**

## Settings & Applications

Goal #1:
Finding patterns
in a single time sequence

CIKM 04                    (c) C. Faloutsos, 2004                    14

---

**CMU SCS**

## Settings & Applications

Goal #2:
Finding patterns
in many time
sequences
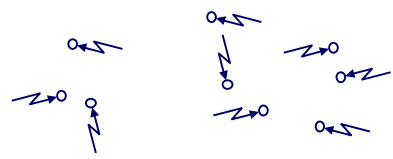
CIKM 04                    (c) C. Faloutsos, 2004                    15

---

**CMU SCS**

## Problem #1:

Goal: given a signal (e.g.., #packets over time)

Find: patterns, periodicities, and/or compress

count

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

year

CIKM 04                    (c) C. Faloutsos, 2004                    16

---

**CMU SCS**

## Problem#2: Forecast

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

??

CIKM 04                    (c) C. Faloutsos, 2004                    17

---

**CMU SCS**

## Problem#2': Similarity search

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

??

CIKM 04                    (c) C. Faloutsos, 2004                    18

---

**CMU SCS**

## Problem #3:

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

**CMU SCS**

## Differences from DSP/Stat

- Semi-infinite streams
  - we need on-line, 'any-time' algorithms
- Can not afford human intervention
  - need automatic methods
- sensors have limited memory / processing / transmitting power
  - need for (lossy) compression

**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)

**CMU SCS**

## Important topics NOT in this tutorial:

- Continuous queries
  - [Babu+Widom ] [Gehrke+] [Madden+]
- Categorical data streams
  - [Hatonen+96]
- Outlier detection (discontinuities)
  - [Breunig+00]
- Related (see D. Shasha's tutorial)

**CMU SCS**

## Outline

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

**CMU SCS**

## Outline

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

**CMU SCS**

## Importance of distance functions

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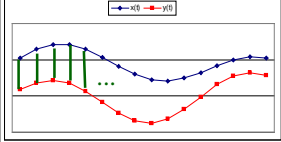- A 'must' for similarity indexing (-> forecasting)
- A 'must' for clustering

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

CIKM 04       (c) C. Faloutsos, 2004       25

---

**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
- $L_\infty$

CIKM 04       (c) C. Faloutsos, 2004       26

---

**CMU SCS**

## Observation #1

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

Day-n

...

Day-2

Day-1

CIKM 04       (c) C. Faloutsos, 2004       27

---

**CMU SCS**

## Observation #2

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

Day-n

...

Day-2

Day-1

CIKM 04       (c) C. Faloutsos, 2004       28

---

**CMU SCS**

## Time Warping

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

CIKM 04       (c) C. Faloutsos, 2004       29

---

**CMU SCS**

## Time Warping

'stutters':

CIKM 04       (c) C. Faloutsos, 2004       30

CMU SCS

**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 $y$

CIKM 04      (c) C. Faloutsos, 2004      31

---

CMU SCS

**Skip**

Time warping    # Time warping

Thus, with no penalty for stutter, for sequences

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

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

CIKM 04      (c) C. Faloutsos, 2004      32

---

CMU SCS

**Skip**

# Time warping

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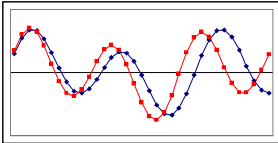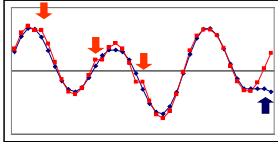- **Many** variations (penalty for stutters; limit on the number/percentage of stutters; …)
- popular in voice processing [Rabiner+Juang]

CIKM 04      (c) C. Faloutsos, 2004      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]

See tutorial by [Gunopulos Das, SIGMOD01]

CIKM 04      (c) C. Faloutsos, 2004      34

---

CMU SCS

# Other Distance functions

- recently: parameter-free, MDL based [Keogh, KDD'04]

CIKM 04      (c) C. Faloutsos, 2004      35

---

CMU SCS

# Conclusions

Prevailing distances:
  – Euclidean and
  – time-warping

CIKM 04      (c) C. Faloutsos, 2004      36

---

**CMU SCS**

# Outline

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

CIKM 04       (c) C. Faloutsos, 2004       37

---

**CMU SCS**

# Indexing

Problem:

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

CIKM 04       (c) C. Faloutsos, 2004       38

---

**CMU SCS**



distance function: by expert

CIKM 04       (c) C. Faloutsos, 2004       39

---

**CMU SCS**

# Idea: 'GEMINI'

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

CIKM 04       (c) C. Faloutsos, 2004       40

---

**CMU SCS**

# 'GEMINI' - Pictorially



CIKM 04       (c) C. Faloutsos, 2004       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

CIKM 04       (c) C. Faloutsos, 2004       42

---

## Examples of GEMINI

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

## 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]

## 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:

## R-trees

**Skip**

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

## R-trees

**Skip**

- eg., w/ fanout 4:

## R-trees

**Skip**

- eg., w/ fanout 4:

**CMU SCS**

# R-trees - range search?

**Skip**

P1    P3    I

A C
B
E
P2 D    F G H
P4 J

P1 P2 P3 P4

A B C
D E
H I J
F G

CIKM 04 (c) C. Faloutsos, 2004 49

---

**CMU SCS**

# R-trees - range search?

**Skip**

P1    P3    I

A C
B
E
P2 D    F G H
P4 J

P1 P2 P3 P4

A B C
D E
H I J
F G

CIKM 04 (c) C. Faloutsos, 2004 50

---

**CMU SCS**

# Conclusions

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

CIKM 04 (c) C. Faloutsos, 2004 51

---

**CMU SCS**

# Outline

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

CIKM 04 (c) C. Faloutsos, 2004 52

---

**CMU SCS**

# Outline

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

CIKM 04 (c) C. Faloutsos, 2004 53

---

**CMU SCS**

# DFT and cousins

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

CIKM 04 (c) C. Faloutsos, 2004 54

**CMU SCS**

## DFT and stocks

Fourier appx — actual

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

CIKM 04      (c) C. Faloutsos, 2004      55

---

**CMU SCS**

## DFT and stocks

Fourier appx — actual

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

Log(ampl)    — Series1

freq

CIKM 04      (c) C. Faloutsos, 2004      56

---

**CMU SCS**

## Outline

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

CIKM 04      (c) C. Faloutsos, 2004      57
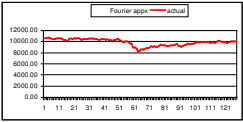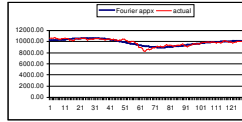
---

**CMU SCS**

## SVD

- THE optimal method for dimensionality reduction
  - (under the Euclidean metric)

CIKM 04      (c) C. Faloutsos, 2004      58

---

**CMU SCS**

## Singular Value Decomposition (SVD)

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

day2

LSI: S. Dumais; M. Berry

KL: eg, Duda+Hart

PCA: eg., Jolliffe

Details: [Press+],

[Faloutsos96]

day1

CIKM 04      (c) C. Faloutsos, 2004      59

---

**CMU SCS**

## SVD

- **Extremely** 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?

CIKM 04      (c) C. Faloutsos, 2004      60

## SVD shorcuts

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

## 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])

## Feature extraction - w/ fractals

Skip

- Main idea: drop those attributes that don't affect the intrinsic ('fractal') dimensionality [Traina+, SBBD 2000]
- i.e.., drop attributes that depend on others (linearly or non-linearly!)

## Fractals

Skip

Fractal dimension
= intrinsic dimension
~ degrees of freedom

Real data: often self-similar, with NON-INTEGER intrinsic dimension (!)

## Feature extraction - w/ fractals

Skip

global FD=1

PFD~1

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

PFD~1

## Outline

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

**CMU SCS**

## MDS / FastMap

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

CIKM 04 (c) C. Faloutsos, 2004 67

**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

CIKM 04 (c) C. Faloutsos, 2004 68

**CMU SCS**

## MDS

Multi Dimensional Scaling



CIKM 04 (c) C. Faloutsos, 2004 69

**CMU SCS**

## FastMap

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

CIKM 04 (c) C. Faloutsos, 2004 70

**CMU SCS**

## FastMap: Application

VideoTrails [Kobla+97]



scene-cut detection (about 10% errors)

CIKM 04 (c) C. Faloutsos, 2004 71

**CMU SCS**

## Outline

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

CIKM 04 (c) C. Faloutsos, 2004 72

**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 intrinsic dimensionalities, consider sequential scan (it might win…)

CIKM 04     (c) C. Faloutsos, 2004     73

**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)

CIKM 04     (c) C. Faloutsos, 2004     74

**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/

CIKM 04     (c) C. Faloutsos, 2004     75

**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.

CIKM 04     (c) C. Faloutsos, 2004     76

**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.

CIKM 04     (c) C. Faloutsos, 2004     77

**CMU SCS**

## References

- Gunopulos, D. and G. Das (2001). Time Series Similarity Measures and Time Series Indexing. SIGMOD Conference, Santa Barbara, CA.
- 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.

CIKM 04     (c) C. Faloutsos, 2004     78

**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.
- Eamonn J. Keogh, Stefano Lonardi, Chotirat (Ann) Ratanamahatana: Towards parameter-free data mining. KDD 2004: 206-215
- Kobla, V., D. S. Doermann, et al. (Nov. 1997). VideoTrails: Representing and Visualizing Structure in Video Sequences. ACM Multimedia 97, Seattle, WA.

CIKM 04          (c) C. Faloutsos, 2004          79

**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.

CIKM 04          (c) C. Faloutsos, 2004          80

**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.

CIKM 04          (c) C. Faloutsos, 2004          81

**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.

CIKM 04          (c) C. Faloutsos, 2004          82

**CMU SCS**

# Part 2: DSP (Digital Signal Processing)

CIKM 04          (c) C. Faloutsos, 2004          83

**CMU SCS**

# Outline

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

CIKM 04          (c) C. Faloutsos, 2004          84
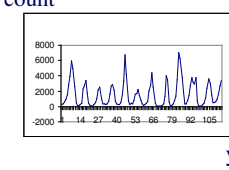
---

**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

CIKM 04          (c) C. Faloutsos, 2004          85

---

**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

CIKM 04          (c) C. Faloutsos, 2004          86

---

**CMU SCS**

## What does DFT do?

A: highlights the periodicities

CIKM 04          (c) C. Faloutsos, 2004          87

---

**CMU SCS**

Skip

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

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

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

inverse DFT

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

CIKM 04          (c) C. Faloutsos, 2004          88

---

**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] ];

CIKM 04          (c) C. Faloutsos, 2004          89

---

**CMU SCS**

## DFT: Amplitude spectrum

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

count

Ampl.

freq=0

freq=12

year

Freq.

CIKM 04          (c) C. Faloutsos, 2004          90

---

**CMU SCS**

Skip

# DFT: examples

flat

Amplitude



time

freq

CIKM 04      (c) C. Faloutsos, 2004      91

**CMU SCS**

Skip

# DFT: examples

Low frequency sinusoid



time

freq

CIKM 04      (c) C. Faloutsos, 2004      92

**CMU SCS**

Skip

# DFT: examples

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



time

freq

CIKM 04      (c) C. Faloutsos, 2004      93

**CMU SCS**

Skip

# DFT: examples

• Higher freq. sinusoid



time

freq

CIKM 04      (c) C. Faloutsos, 2004      94

**CMU SCS**

Skip

# DFT: examples

examples



CIKM 04      (c) C. Faloutsos, 2004      95

**CMU SCS**

Skip

# DFT: examples

examples

Ampl.



Freq.

CIKM 04      (c) C. Faloutsos, 2004      96

**CMU SCS**

### Outline

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

CIKM 04        (c) C. Faloutsos, 2004        97

**CMU SCS**

### Outline

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

CIKM 04        (c) C. Faloutsos, 2004        98

**CMU SCS**

### DFT: Amplitude spectrum

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



count

Ampl.

freq=0

freq=12

year

Freq.

CIKM 04        (c) C. Faloutsos, 2004        99

**CMU SCS**

### DFT: Amplitude spectrum



count

Ampl.

freq=0

freq=12

year

Freq.

CIKM 04        (c) C. Faloutsos, 2004        100

**CMU SCS**

### DFT: Amplitude spectrum



count

Ampl.

freq=0

freq=12

year

Freq.

CIKM 04        (c) C. Faloutsos, 2004        101

**CMU SCS**

### DFT: Amplitude spectrum

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



Freq.

CIKM 04        (c) C. Faloutsos, 2004        102

**CMU SCS**

# DFT: Amplitude spectrum

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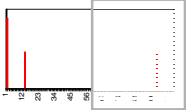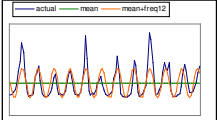- A2: pattern discovery

CIKM 04 (c) C. Faloutsos, 2004 103

**CMU SCS**

# DFT: Amplitude spectrum

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

CIKM 04 (c) C. Faloutsos, 2004 104

**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)

CIKM 04 (c) C. Faloutsos, 2004 105

**CMU SCS**

# Outline

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

CIKM 04 (c) C. Faloutsos, 2004 106

**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

CIKM 04 (c) C. Faloutsos, 2004 107

**CMU SCS**

# Wavelets - DWT

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

value

time

CIKM 04 (c) C. Faloutsos, 2004 108

## Wavelets - DWT

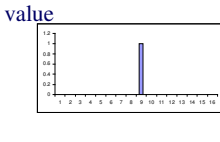- DFT is great - but, how about compressing a spike?
- A: Terrible - all DFT coefficients needed!

value

Ampl

time

Freq

CIKM 04      (c) C. Faloutsos, 2004      109

## Wavelets - DWT

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

value

time

CIKM 04      (c) C. Faloutsos, 2004      110

## Wavelets - DWT

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

value

time

CIKM 04      (c) C. Faloutsos, 2004      111

## Wavelets - DWT

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

freq

value

time

time

CIKM 04      (c) C. Faloutsos, 2004      112

## Wavelets - DWT

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

Time domain    DFT    SWFT    DWT

freq

time

CIKM 04      (c) C. Faloutsos, 2004      113

## Haar Wavelets

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

CIKM 04      (c) C. Faloutsos, 2004      114

**CMU SCS**

# Wavelets - construction

Skip

x0  x1  x2  x3  x4  x5  x6  x7

CIKM 04 — (c) C. Faloutsos, 2004 — 115

---

**CMU SCS**

# Wavelets - construction

Skip

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

−   +

x0  x1  x2  x3  x4  x5  x6  x7

CIKM 04 — (c) C. Faloutsos, 2004 — 116

---

**CMU SCS**

# Wavelets - construction

Skip

level 2   d2,0    s2,0

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

−   +

x0  x1  x2  x3  x4  x5  x6  x7

CIKM 04 — (c) C. Faloutsos, 2004 — 117

---

**CMU SCS**

# Wavelets - construction

Skip

etc ...

d2,0    s2,0

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

−   +

x0  x1  x2  x3  x4  x5  x6  x7

CIKM 04 — (c) C. Faloutsos, 2004 — 118

---

**CMU SCS**

# Wavelets - construction

Skip

Q: map each coefficient

on the time-freq. plane

f

d2,0    s2,0

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

−   +

t

x0  x1  x2  x3  x4  x5  x6  x7

CIKM 04 — (c) C. Faloutsos, 2004 — 119

---

**CMU SCS**

# Wavelets - construction

Skip

Q: map each coefficient

on the time-freq. plane

f

d2,0    s2,0

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

−   +

t

x0  x1  x2  x3  x4  x5  x6  x7

CIKM 04 — (c) C. Faloutsos, 2004 — 120

---

## Slide 1

# Haar wavelets - code

```perl
#!/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 );
}
```

```perl
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
```

CIKM 04          (c) C. Faloutsos, 2004          121

## Slide 2

# 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, ...

CIKM 04          (c) C. Faloutsos, 2004          122

## Slide 3

# Wavelets - how do they look like?



- E.g., Daubechies-4

CIKM 04          (c) C. Faloutsos, 2004          123

## Slide 4

# Wavelets - how do they look like?



- E.g., Daubechies-4

?

?

CIKM 04          (c) C. Faloutsos, 2004          124

## Slide 5

# Wavelets - how do they look like?



- E.g., Daubechies-4

CIKM 04          (c) C. Faloutsos, 2004          125

## Slide 6

# Outline

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

CIKM 04          (c) C. Faloutsos, 2004          126

C. Faloutsos



**Wavelets - Drill#1:**

- Q: baritone/silence/soprano - DWT?

f

t

value

time

CIKM 04    (c) C. Faloutsos, 2004    127

**Wavelets - Drill#1:**

- Q: baritone/soprano - DWT?

f

t

value

time

CIKM 04    (c) C. Faloutsos, 2004    128

**Wavelets - Drill#2:**

- Q: spike - DWT?

f

t

CIKM 04    (c) C. Faloutsos, 2004    129

**Wavelets - Drill#2:**

- Q: spike - DWT?

f

t

0.00    0.00    **0.71**    0.00

0.00    **0.50**

**-0.35**

**0.35**

CIKM 04    (c) C. Faloutsos, 2004    130

**Wavelets - Drill#3:**

- Q: weekly + daily periodicity, + spike - DWT?

f

t

CIKM 04    (c) C. Faloutsos, 2004    131

**Wavelets - Drill#3:**

- Q: **weekly** + daily periodicity, + spike - DWT?

f

t

CIKM 04    (c) C. Faloutsos, 2004    132

**CMU SCS**

# Wavelets - Drill#3:

- Q: weekly + **daily** periodicity, + spike - DWT?

f

t

**CMU SCS**

# Wavelets - Drill#3:

- Q: weekly + daily periodicity, + **spike** - DWT?

f

t

**CMU SCS**

# Wavelets - Drill#3:

- Q: weekly + daily periodicity, + spike - DWT?

f

t

**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

CIKM 04     (c) C. Faloutsos, 2004     139

**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)

CIKM 04     (c) C. Faloutsos, 2004     140

**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

CIKM 04     (c) C. Faloutsos, 2004     141

**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 DFT, DWT)
- C. Faloutsos: *Searching Multimedia Databases by Content*, Kluwer Academic Press, 1996 (introduction to DFT, DWT)

CIKM 04     (c) C. Faloutsos, 2004     142

**CMU SCS**

## 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

CIKM 04     (c) C. Faloutsos, 2004     143

**CMU SCS**

## BREAK!

CIKM 04     (c) C. Faloutsos, 2004     144

**CMU SCS**

## Sensor Data Mining: Similarity Search and Pattern Analysis

*Christos Faloutsos*
CMU

---

**CMU SCS**

# Part 3: Linear Forecasting

CIKM 04          (c) C. Faloutsos, 2004          146

---

**CMU SCS**

## Outline

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

CIKM 04          (c) C. Faloutsos, 2004          147

---

**CMU SCS**

## Forecasting

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

**http://www.hfac.uh.edu/MediaFutures/t houghts.html**

CIKM 04          (c) C. Faloutsos, 2004          148

---

**CMU SCS**

## Outline

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

CIKM 04          (c) C. Faloutsos, 2004          149

---

**CMU SCS**

## Problem#2: Forecast

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

??

CIKM 04          (c) C. Faloutsos, 2004          150

---

**CMU SCS**

## Forecasting: Preprocessing

MANUALLY:

remove trends        spot periodicities

7 days

time                time

CIKM 04        (c) C. Faloutsos, 2004        151

---

**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:

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

**Time Tick**

CIKM 04        (c) C. Faloutsos, 2004        152

---

**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

**Time Tick**

CIKM 04        (c) C. Faloutsos, 2004        153

---

**CMU SCS**

## Linear Regression: idea

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

**Body height**

**Body weight**

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

CIKM 04        (c) C. Faloutsos, 2004        154

---

**CMU SCS**

## Linear Auto Regression:

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

CIKM 04        (c) C. Faloutsos, 2004        155

---

**CMU SCS**

## Linear Auto Regression:

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

**Number of packets sent (t)**

'lag-plot'

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

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

CIKM 04        (c) C. Faloutsos, 2004        156

---

**CMU SCS**

## Outline

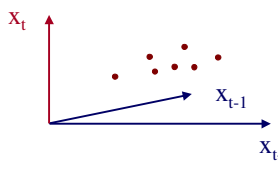- Motivation
- ...
- Linear Forecasting
  - Auto-regression: **Least Squares; RLS**
  - Co-evolving time sequences
  - Examples
  - Conclusions

CIKM 04 (c) C. Faloutsos, 2004 157

---

**CMU SCS**

## More details:

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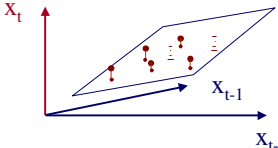- A1: YES!



CIKM 04 (c) C. Faloutsos, 2004 158

---

**CMU SCS**

## More details:

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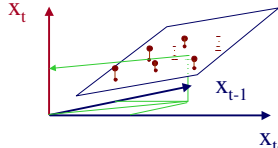- A1: YES! (we'll fit a hyper-plane, then!)



CIKM 04 (c) C. Faloutsos, 2004 159

---

**CMU SCS**

## More details:

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



CIKM 04 (c) C. Faloutsos, 2004 160

---

**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

CIKM 04 (c) C. Faloutsos, 2004 161

---

**CMU SCS**

**Skip**

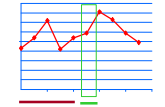## More details:

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

Ind-var1    Ind-var-w

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

CIKM 04 (c) C. Faloutsos, 2004 162

---

**Slide 163:**

CMU SCS

Skip

## More details:

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

Ind-var1    Ind-var-w

time

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

CIKM 04       (c) C. Faloutsos, 2004       163

---

**Slide 164:**

CMU SCS

Skip

## 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 $\mathbf{y}$

CIKM 04       (c) C. Faloutsos, 2004       164

---

**Slide 165:**

CMU SCS

Skip

## 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:

CIKM 04       (c) C. Faloutsos, 2004       165

---

**Slide 166:**

CMU SCS

Skip

## Even more details

- Given:

Dependent Variable

Independent Variable

CIKM 04       (c) C. Faloutsos, 2004       166

---

**Slide 167:**

CMU SCS

Skip

## Even more details

Dependent Variable

← new point

Independent Variable

CIKM 04       (c) C. Faloutsos, 2004       167

---

**Slide 168:**

CMU SCS

Skip

## Even more details

### RLS: quickly compute new best fit

Dependent Variable

← new point

Independent Variable

CIKM 04       (c) C. Faloutsos, 2004       168

---

**CMU SCS**

**Skip**

# 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 $O(1 \times w^2)$

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

CIKM 04     (c) C. Faloutsos, 2004     169

---

**CMU SCS**

**Skip**

# Even more details

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

CIKM 04     (c) C. Faloutsos, 2004     170

---

**CMU SCS**

**Skip**

# Adaptability - 'forgetting'



CIKM 04     (c) C. Faloutsos, 2004     171

---

**CMU SCS**

**Skip**

# Adaptability - 'forgetting'



CIKM 04     (c) C. Faloutsos, 2004     172

---

**CMU SCS**

**Skip**

# Adaptability - 'forgetting'



- RLS: can *trivially* handle 'forgetting'

CIKM 04     (c) C. Faloutsos, 2004     173

---

**CMU SCS**

# How to choose '$w$'?

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

CIKM 04     (c) C. Faloutsos, 2004     174

**Slide 175:**

## Answer:

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

CIKM 04      (c) C. Faloutsos, 2004      175

**Slide 176:**

## AWSOM



CIKM 04      (c) C. Faloutsos, 2004      176

**Slide 177:**

## AWSOM



CIKM 04      (c) C. Faloutsos, 2004      177

**Slide 178:**

## AWSOM - idea

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

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

CIKM 04      (c) C. Faloutsos, 2004      178

**Slide 179:**

## 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")

CIKM 04      (c) C. Faloutsos, 2004      179

**Slide 180:**

## Results - Synthetic data



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

CIKM 04      (c) C. Faloutsos, 2004      180

**CMU SCS**

## Results - Real data



Automobile – Original | Automobile – AWSOM (3) | Automobile – AR (36) | Automobile – SAR

FAILED

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

CIKM 04      (c) C. Faloutsos, 2004      181

---

**CMU SCS**

## Results - real data



Sunspot – Original | Sunspot – AWSOM (6,1) | Sunspot – AR (60) | nspot – SARIMA (2,1,0) x (1,1,0

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

CIKM 04      (c) C. Faloutsos, 2004      182

---

**CMU SCS**

## 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; $O(lgN)$

CIKM 04      (c) C. Faloutsos, 2004      183

---

**CMU SCS**

## Outline

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

CIKM 04      (c) C. Faloutsos, 2004      184

---

**CMU SCS**

## Co-Evolving Time Sequences

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



CIKM 04      (c) C. Faloutsos, 2004      185

---

**CMU SCS**

## Solution:

Q: what should we do?

CIKM 04      (c) C. Faloutsos, 2004      186

**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])

CIKM 04      (c) C. Faloutsos, 2004      187

---

**CMU SCS**

**Skip**

## B.II - Time Series Analysis Outline

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

CIKM 04      (c) C. Faloutsos, 2004      188

---

**CMU SCS**

**Skip**

## Examples - Experiments

- Datasets
  - Modem pool traffic (14 modems, 1500 time-ticks; #packets per time unit)
  - AT&T WorldNet internet usage (several data streams; 980 time-ticks)
- Measures of success
  - Accuracy : Root Mean Square Error (RMSE)

CIKM 04      (c) C. Faloutsos, 2004      189

---

**CMU SCS**

**Skip**

## Accuracy - "Modem"



MUSCLES outperforms AR & "yesterday"

CIKM 04      (c) C. Faloutsos, 2004      190

---

**CMU SCS**

**Skip**

## Accuracy - "Internet"



MUSCLES consistently outperforms AR & "yesterday"

CIKM 04      (c) C. Faloutsos, 2004      191

---

**CMU SCS**

**Skip**

## B.II - Time Series Analysis Outline

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

CIKM 04      (c) C. Faloutsos, 2004      192

---

**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)

CIKM 04        (c) C. Faloutsos, 2004        193

**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/`

CIKM 04        (c) C. Faloutsos, 2004        194

**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.

CIKM 04        (c) C. Faloutsos, 2004        195

**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)

CIKM 04        (c) C. Faloutsos, 2004        196

**CMU SCS**

# Part 4: Bursty traffic and multifractals

CIKM 04        (c) C. Faloutsos, 2004        197

**CMU SCS**

## Outline

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

CIKM 04        (c) C. Faloutsos, 2004        198

**CMU SCS**

## Outline

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

CIKM 04 (c) C. Faloutsos, 2004 199

---

**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

CIKM 04 (c) C. Faloutsos, 2004 200

---

**CMU SCS**

## Problem #1

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

# bytes



Poisson →

time

CIKM 04 (c) C. Faloutsos, 2004 201

---

**CMU SCS**

## Motivation

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

CIKM 04 (c) C. Faloutsos, 2004 202

---

**CMU SCS**

## Q: any 'pattern'?

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

# bytes



time

CIKM 04 (c) C. Faloutsos, 2004 203

---

**CMU SCS**

## solution: self-similarity

# bytes

# bytes



time

time

CIKM 04 (c) C. Faloutsos, 2004 204

---

**CMU SCS**

## But:

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

CIKM 04 (c) C. Faloutsos, 2004 205

**CMU SCS**

## Outline

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

CIKM 04 (c) C. Faloutsos, 2004 206

**CMU SCS**

## Approach

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

CIKM 04 (c) C. Faloutsos, 2004 207

**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%)

CIKM 04 (c) C. Faloutsos, 2004 208

**CMU SCS**

## Binary multifractals

**20** **80**



CIK 209

**CMU SCS**

## Binary multifractals

**20** **80**



CIK 210

---

**CMU SCS**

# Parameter estimation

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

CIKM 04          (c) C. Faloutsos, 2004          211

---

**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]

CIKM 04          (c) C. Faloutsos, 2004          212

---

**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.

CIKM 04          (c) C. Faloutsos, 2004          213

---

**CMU SCS**

# Entropy plot

p1          p2

% of bytes here



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

CIKM 04          (c) C. Faloutsos, 2004          214

---

**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: $E(2) = -\Sigma_t \, p_{2,i} * \log_2 (p_{2,i})$

CIKM 04          (c) C. Faloutsos, 2004          215

---

**CMU SCS**

# Real traffic

Entropy *E(n)*



0.73

Real
0.73*x-0.24

# of levels (*n*)

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

CIKM 04          (c) C. Faloutsos, 2004          216

---

**CMU SCS**

## Observation - intuition:

Entropy

$E(n)$

intuition: slope =

intrinsic dimensionality =

info-bits per coordinate-bit

– unif. Dataset: slope =1

– multi-point: slope = 0

# of levels ($n$)

CIKM 04      (c) C. Faloutsos, 2004      217

---

**CMU SCS**

## 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?

CIKM 04      (c) C. Faloutsos, 2004      218

---

**CMU SCS**

## 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

CIKM 04      (c) C. Faloutsos, 2004      219

---

**CMU SCS**

## 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?

CIKM 04      (c) C. Faloutsos, 2004      220

---

**CMU SCS**

## 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?

CIKM 04      (c) C. Faloutsos, 2004      221

---

**CMU SCS**

## Entropy plot

- Repeat, for all points at same position:

Dim=0

CIKM 04      (c) C. Faloutsos, 2004      222

C. Faloutsos

---

**CMU SCS**

**Skip**

## Entropy plot

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

Dim=0

CIKM 04    (c) C. Faloutsos, 2004    223

---

**CMU SCS**

**Skip**

## 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

CIKM 04    (c) C. Faloutsos, 2004    224

---

**CMU SCS**

## (Fractals, again)

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

CIKM 04    (c) C. Faloutsos, 2004    225

---

**CMU SCS**

## Cantor dust

- Eliminate the middle third
- Recursively!

CIKM 04    (c) C. Faloutsos, 2004    226

---

**CMU SCS**

## Cantor dust

CIKM 04    (c) C. Faloutsos, 2004    227

---

**CMU SCS**

## Cantor dust

CIKM 04    (c) C. Faloutsos, 2004    228

---

CMU

**CMU SCS**

# Cantor dust



CIKM 04     (c) C. Faloutsos, 2004     229

---

**CMU SCS**

# Cantor dust



CIKM 04     (c) C. Faloutsos, 2004     230

---

**CMU SCS**

# Cantor dust



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

CIKM 04     (c) C. Faloutsos, 2004     231

---

**CMU SCS**

# Some more entropy plots:

- **Poisson** vs real



Poisson: slope = ~1 -> uniformly distributed

CIKM 04     (c) C. Faloutsos, 2004     232

---

**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*

CIKM 04     (c) C. Faloutsos, 2004     233

---

**CMU SCS**

# Outline

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

CIKM 04     (c) C. Faloutsos, 2004     234

---

**CMU SCS**

# Experimental setup

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

CIKM 04 (c) C. Faloutsos, 2004 235

---

**CMU SCS**

# Model validation

- Linear entropy plots



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

CIKM 04 (c) C. Faloutsos, 2004 236

---

**CMU SCS**

# Web traffic - results

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

Prob( $>l$)



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

CIKM 04 (c) C. Faloutsos, 2004 237

---

**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 $l$)

CIKM 04 (c) C. Faloutsos, 2004 238

---

**CMU SCS**

# Conclusions

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

CIKM 04 (c) C. Faloutsos, 2004 239

---

**CMU SCS**

# 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!)

CIKM 04 (c) C. Faloutsos, 2004 240

---

**CMU SCS**

## 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.

CIKM 04 (c) C. Faloutsos, 2004 241

---

**CMU SCS**

## 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.

CIKM 04 (c) C. Faloutsos, 2004 242

---

**CMU SCS**

## Part 5:
## chaos and
## non-linear forecasting

CIKM 04 (c) C. Faloutsos, 2004 243

---

**CMU SCS**

## Outline

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

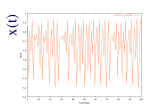CIKM 04 (c) C. Faloutsos, 2004 244

---

**CMU SCS**

## Detailed Outline

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

CIKM 04 (c) C. Faloutsos, 2004 245

---

**CMU SCS**

## Recall: Problem #1

Value



Time

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

CIKM 04 (c) C. Faloutsos, 2004 246

### How to forecast?

- ARIMA - but: linearity assumption

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

CIKM 04      (c) C. Faloutsos, 2004      247

---

### General Intuition (Lag Plot)

CIKM 04      (c) C. Faloutsos, 2004      248

---

### 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?

CIKM 04      (c) C. Faloutsos, 2004      249

---

### Q1: Choosing lag $L$

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

CIKM 04      (c) C. Faloutsos, 2004      250

---

### Q2: Choosing number of neighbors $k$

- Manually (typically ~ 1-10)

CIKM 04      (c) C. Faloutsos, 2004      251

---

### 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?)

CIKM 04      (c) C. Faloutsos, 2004      252

**CMU SCS**

## Q3: How to interpolate?

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



CIKM 04     (c) C. Faloutsos, 2004     253

**CMU SCS**

## Q4: Any theory behind it?

A4: YES!

CIKM 04     (c) C. Faloutsos, 2004     254

**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)

CIKM 04     (c) C. Faloutsos, 2004     255

**CMU SCS**

**Skip**

## Theoretical foundation

Example: Lotka-Volterra equations

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



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, …).

CIKM 04     (c) C. Faloutsos, 2004     256

**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**



CIKM 04     (c) C. Faloutsos, 2004     257

**CMU SCS**

## Detailed Outline

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

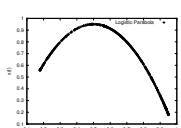CIKM 04     (c) C. Faloutsos, 2004     258

# Datasets

x(t)

time

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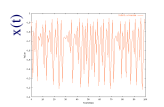
Lag-plot

CIKM 04          (c) C. Faloutsos, 2004          259
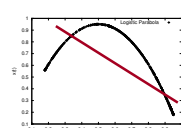
---

# Datasets

x(t)

time

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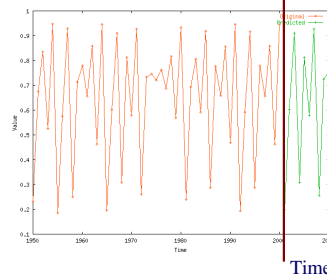Models population of flies [R. May/1976]

Lag-plot

ARIMA: fails

CIKM 04          (c) C. Faloutsos, 2004          260
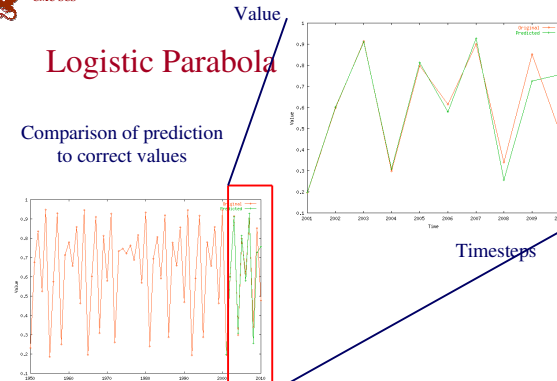
---

# Logistic Parabola

Our Prediction from here

Value

Timesteps

CIKM 04          (c) C. Faloutsos, 2004          261

---

Value

# Logistic Parabola

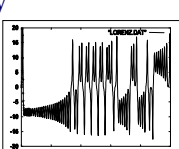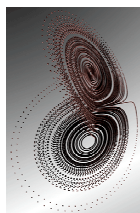Comparison of prediction to correct values

Timesteps

CIKM 04          (c) C. Faloutsos, 2004          262

---

Value

# Datasets

Skip

LORENZ: Models convection currents in the air
$dx / dt = a (y - x)$
$dy / dt = x (b - z) - y$
$dz / dt = xy - c z$
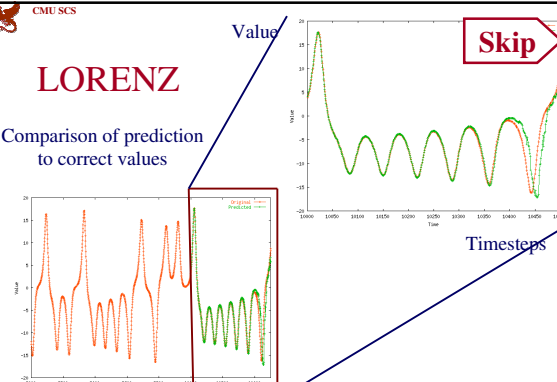
CIKM 04          (c) C. Faloutsos, 2004          263

---

Value

Skip

# LORENZ

Comparison of prediction to correct values

Timesteps

CIKM 04          (c) C. Faloutsos, 2004          264

**CMU SCS**

Value

## Datasets

**Skip**

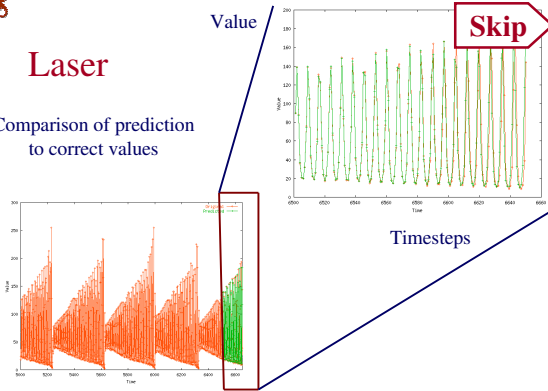- LASER: fluctuations in a Laser over time (used in Santa Fe competition)

Time

CIKM 04          (c) C. Faloutsos, 2004          265

---

**CMU SCS**

Value

## Laser

**Skip**

Comparison of prediction to correct values

Timesteps

CIKM 04          (c) C. Faloutsos, 2004          266

---

**CMU SCS**

## Conclusions

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

CIKM 04          (c) C. Faloutsos, 2004          267

---

**CMU SCS**

## 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.

CIKM 04          (c) C. Faloutsos, 2004          268

---

**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.)

CIKM 04          (c) C. Faloutsos, 2004          269

---

**CMU SCS**

## Overall conclusions

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

CIKM 04          (c) C. Faloutsos, 2004          270

---

**CMU SCS**

## Overall conclusions

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

**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) methodology

**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) methodology; AWSOM
- Bursty traffic: **multifractals** (80-20 'law')

**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) methodology
- Bursty traffic: **multifractals** (80-20 'law')
- Non-linear forecasting: **lag-plots** (Takens)

**CMU SCS**

## 'Take home' messages

- Hard, but desirable query for sensor data: '*find patterns / outliers*'
- We need **fast**, **automated** such tools
  - Many great tools exist (DWT, ARIMA, …)
  - some are readily usable; others need to be made scalable / single pass/ automatic

**CMU SCS**

# THANK YOU!

`christos@cs.cmu.edu`
`www.cs.cmu.edu/~christos`