

CMU SCS

Indexing and Mining Streams

Christos Faloutsos
CMU

SIGMOD 04 Copyright: C. Faloutsos, 2004

CMU SCS

Thanks


 Deepay Chakrabarti (CMU)


 Prof. Dimitris Gunopulos (UCR)


 Spiros Papadimitriou (CMU)


 Mengzhi Wang (CMU)


 Prof. Byoung-Kee Yi (Pohang U.)

SIGMOD 04 Copyright: 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

SIGMOD 04 Copyright: C. Faloutsos, 2004 3

CMU SCS

Problem definition

- Given: one or more sequences
 $x_1, x_2, \dots, x_t, \dots$
 $(y_1, y_2, \dots, y_p, \dots$
 $\dots)$
- Find
 - similar sequences; forecasts
 - patterns; clusters; outliers

SIGMOD 04 Copyright: C. Faloutsos, 2004 4

CMU SCS

Motivation - Applications

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 5

CMU SCS

Motivation - Applications (cont'd)

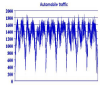
- 'Smart house'
 - sensors monitor temperature, humidity, air quality
- video surveillance

SIGMOD 04 Copyright: C. Faloutsos, 2004 6

CMU SCS

Motivation - Applications (cont'd)

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

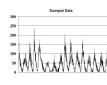


SIGMOD 04 Copyright: C. Faloutsos, 2004 7

CMU SCS

Motivation - Applications (cont'd)

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



SIGMOD 04 Copyright: C. Faloutsos, 2004 8

CMU SCS

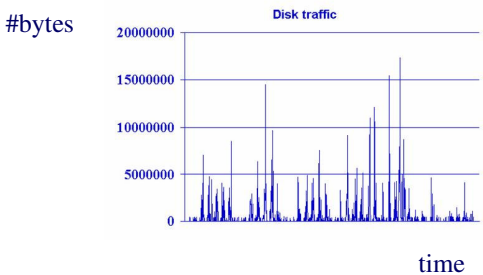
Motivation - Applications (cont'd)

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 9

CMU SCS

Stream Data: Disk accesses



SIGMOD 04 Copyright: C. Faloutsos, 2004 10

CMU SCS

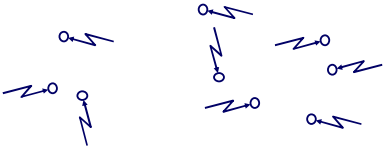
Settings & Applications

- One or more sensors, collecting time-series data

SIGMOD 04 Copyright: C. Faloutsos, 2004 11

CMU SCS

Settings & Applications



Each sensor collects data $(x_1, x_2, \dots, x_p, \dots)$

SIGMOD 04 Copyright: C. Faloutsos, 2004 12

CMU SCS

Settings & Applications

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 13

CMU SCS

Settings & Applications

Goal #1:
Finding patterns
in a single time sequence

SIGMOD 04 Copyright: C. Faloutsos, 2004 14

CMU SCS

Settings & Applications

Goal #2:
Finding patterns
in many time sequences

SIGMOD 04 Copyright: C. Faloutsos, 2004 15

CMU SCS

Problem #1:

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

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 16

CMU SCS

Problem#2: Forecast

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

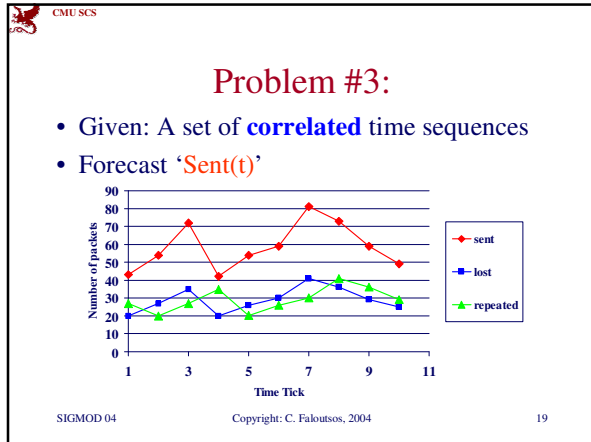
SIGMOD 04 Copyright: C. Faloutsos, 2004 17

CMU SCS

Problem#2': Similarity search

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 18



- 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
- SIGMOD 04 Copyright: C. Faloutsos, 2004 20

- 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)
- SIGMOD 04 Copyright: C. Faloutsos, 2004 21

- 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)
- SIGMOD 04 Copyright: C. Faloutsos, 2004 22

- CMU SCS
- ### Outline
- Motivation
 - ➡ Similarity Search and Indexing
 - DSP
 - Linear Forecasting
 - Bursty traffic - fractals and multifractals
 - Non-linear forecasting
 - Conclusions
- SIGMOD 04 Copyright: C. Faloutsos, 2004 23

- CMU SCS
- ### Outline
- Motivation
 - ➡ Similarity Search and Indexing
 - distance functions: Euclidean; Time-warping
 - indexing
 - feature extraction
 - DSP
 - ...
- SIGMOD 04 Copyright: C. Faloutsos, 2004 24

CMU SCS

Importance of distance functions

Subtle, but **absolutely necessary**:

- A 'must' for similarity indexing (-> forecasting)
- A 'must' for clustering

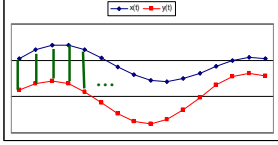
Two major families

- Euclidean and L_p norms
- Time warping and variations

SIGMOD 04 Copyright: C. Faloutsos, 2004 25

CMU SCS

Euclidean and L_p



$$D(\bar{x}, \bar{y}) = \sum_{i=1}^n (x_i - y_i)^2$$

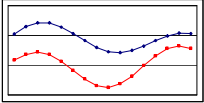
$$L_p(\bar{x}, \bar{y}) = \sum_{i=1}^n |x_i - y_i|^p$$

- L₁: city-block = Manhattan
- L₂ = Euclidean
- L_∞

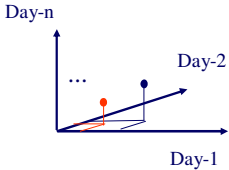
SIGMOD 04 Copyright: C. Faloutsos, 2004 26

CMU SCS

Observation #1



- Time sequence -> n-d vector



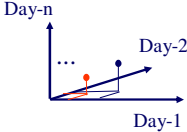
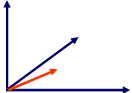
SIGMOD 04 Copyright: C. Faloutsos, 2004 27

CMU SCS

Observation #2

Euclidean distance is closely related to

- cosine similarity
- dot product
- 'cross-correlation' function

SIGMOD 04 Copyright: C. Faloutsos, 2004 28

CMU SCS

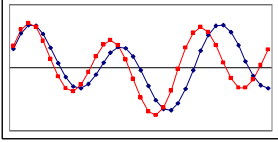
Time Warping

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

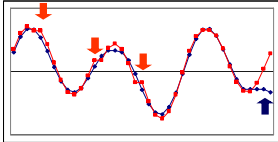
SIGMOD 04 Copyright: C. Faloutsos, 2004 29

CMU SCS

Time Warping



'stutters':



SIGMOD 04 Copyright: 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

SIGMOD 04 Copyright: C. Faloutsos, 2004 31

CMU SCS Skip

Time warping **Time warping**

Thus, with no penalty for stutter, for sequences
 $x_1, x_2, \dots, x_i; \quad y_1, y_2, \dots, y_j$

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

SIGMOD 04 Copyright: 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]

SIGMOD 04 Copyright: 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]

SIGMOD 04 Copyright: C. Faloutsos, 2004 34

CMU SCS

Conclusions

Prevailing distances:

- Euclidean and
- time-warping

SIGMOD 04 Copyright: C. Faloutsos, 2004 35

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 36

CMU SCS

Indexing

Problem:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 37

CMU SCS

Price

1 365 day

Price

1 365 day

Price

1 365 day

Price

1 365 day

distance function: by expert

SIGMOD 04 Copyright: C. Faloutsos, 2004 38

CMU SCS

Idea: 'GEMINI'

Eg., 'find stocks similar to MSFT'

Seq. scanning: too slow

How to accelerate the search?

[Faloutsos96]

SIGMOD 04 Copyright: C. Faloutsos, 2004 39

CMU SCS

'GEMINI' - Pictorially

S1

1 365 day

Sn

1 365 day

eg., std

• F(S1)

• F(Sn)

eg., avg

SIGMOD 04 Copyright: C. Faloutsos, 2004 40

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 41

CMU SCS

Examples of GEMINI

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 42

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]

SIGMOD 04 Copyright: C. Faloutsos, 2004 43

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:

SIGMOD 04 Copyright: C. Faloutsos, 2004 44

CMU SCS

R-trees

Skip

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 45

CMU SCS

R-trees

Skip

- eg., w/ fanout 4:

SIGMOD 04 Copyright: C. Faloutsos, 2004 46

CMU SCS

R-trees

Skip

- eg., w/ fanout 4:

SIGMOD 04 Copyright: C. Faloutsos, 2004 47

CMU SCS

R-trees - range search?

Skip

SIGMOD 04 Copyright: C. Faloutsos, 2004 48

CMU SCS

Skip

R-trees - range search?

SIGMOD 04 Copyright: C. Faloutsos, 2004 49

CMU SCS

Conclusions

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 50

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 51

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 52

CMU SCS

DFT and cousins

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 53

CMU SCS

DFT and stocks

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 54

CMU SCS

DFT and stocks

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 55

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 56

CMU SCS

SVD

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 57

CMU SCS

Singular Value Decomposition (SVD)

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

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 58

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?

SIGMOD 04 Copyright: C. Faloutsos, 2004 59

CMU SCS

SVD shortcuts

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 60

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])

SIGMOD 04 Copyright: C. Faloutsos, 2004 61

CMU SCS

Skip

Feature extraction - w/ fractals

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 62

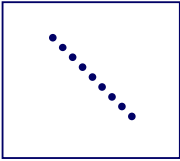
CMU SCS

Skip

Fractals

Fractal dimension
= intrinsic dimension
~ degrees of freedom

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



SIGMOD 04 Copyright: C. Faloutsos, 2004 63

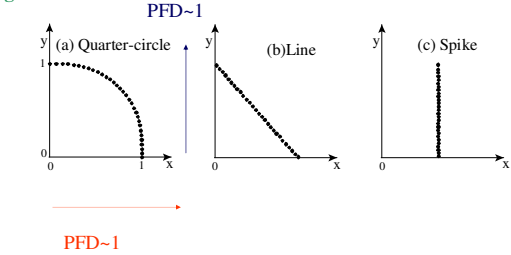
CMU SCS

Skip

Feature extraction - w/ fractals

global FD=1

PFD~1



SIGMOD 04 Copyright: C. Faloutsos, 2004 64

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 65

CMU SCS

MDS / FastMap

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 66

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 67

CMU SCS

MDS

Multi Dimensional Scaling

SIGMOD 04 Copyright: C. Faloutsos, 2004 68

CMU SCS

FastMap

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 69

CMU SCS

FastMap: Application

VideoTrails [Kobla+97]

scene-cut detection (about 10% errors)

SIGMOD 04 Copyright: C. Faloutsos, 2004 70

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 71


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...)

SIGMOD 04 Copyright: C. Faloutsos, 2004 72




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)

SIGMOD 04 Copyright: C. Faloutsos, 2004 73




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/>

SIGMOD 04 Copyright: C. Faloutsos, 2004 74




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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 75




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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 76




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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 77



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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 78

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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 79

CMU SCS

References

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 80

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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 81

CMU SCS

Part 2: DSP (Digital Signal Processing)

SIGMOD 04 Copyright: C. Faloutsos, 2004 82

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 83

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 84

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 85

CMU SCS

What does DFT do?

A: highlights the periodicities

SIGMOD 04 Copyright: C. Faloutsos, 2004 86

CMU SCS

DFT: definition

Skip

- 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} :

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

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

$$x_t = 1/\sqrt{n} \sum_{f=0}^{n-1} X_f * \exp(+j2\pi tf/n) \quad \swarrow \text{inverse DFT}$$

SIGMOD 04 Copyright: C. Faloutsos, 2004 87

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] ];
```

SIGMOD 04 Copyright: C. Faloutsos, 2004 88

CMU SCS

DFT: Amplitude spectrum

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

count

year

Ampl.

freq=0
freq=12

Freq.

SIGMOD 04 Copyright: C. Faloutsos, 2004 89

CMU SCS

DFT: examples

Skip

flat

time

Amplitude

freq

SIGMOD 04 Copyright: C. Faloutsos, 2004 90

CMU SCS Skip

DFT: examples

Low frequency sinusoid

time freq

SIGMOD 04 Copyright: C. Faloutsos, 2004 91

CMU SCS Skip

DFT: examples

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

time freq

SIGMOD 04 Copyright: C. Faloutsos, 2004 92

CMU SCS Skip

DFT: examples

- Higher freq. sinusoid

time freq

SIGMOD 04 Copyright: C. Faloutsos, 2004 93

CMU SCS Skip

DFT: examples

examples

SIGMOD 04 Copyright: C. Faloutsos, 2004 94

CMU SCS Skip

DFT: examples

examples

Ampl. Freq.

SIGMOD 04 Copyright: C. Faloutsos, 2004 95

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 96

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 97

CMU SCS

DFT: Amplitude spectrum

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

count

year

Ampl.

Freq.

SIGMOD 04 Copyright: C. Faloutsos, 2004 98

CMU SCS

DFT: Amplitude spectrum

count

year

Ampl.

Freq.

SIGMOD 04 Copyright: C. Faloutsos, 2004 99

CMU SCS

DFT: Amplitude spectrum

count

year

Ampl.

Freq.

SIGMOD 04 Copyright: C. Faloutsos, 2004 100

CMU SCS

DFT: Amplitude spectrum

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

Freq.

SIGMOD 04 Copyright: C. Faloutsos, 2004 101

CMU SCS

DFT: Amplitude spectrum

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 102

CMU SCS

DFT: Amplitude spectrum

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 103

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)

SIGMOD 04 Copyright: C. Faloutsos, 2004 104

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 105

CMU SCS

Problem #1:

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

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 106

CMU SCS

Wavelets - DWT

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 107

CMU SCS

Wavelets - DWT

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 108

CMU SCS

Wavelets - DWT

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

value

time

SIGMOD 04 Copyright: C. Faloutsos, 2004 109

CMU SCS

Wavelets - DWT

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

value

time

SIGMOD 04 Copyright: C. Faloutsos, 2004 110

CMU SCS

Wavelets - DWT

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

freq

value

time

SIGMOD 04 Copyright: C. Faloutsos, 2004 111

CMU SCS

Wavelets - DWT

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

Time domain

freq

DFT

SWFT

DWT

time

SIGMOD 04 Copyright: C. Faloutsos, 2004 112

CMU SCS

Haar Wavelets

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 113

CMU SCS

Wavelets - construction

Skip

$x_0 \ x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ x_6 \ x_7$

SIGMOD 04 Copyright: C. Faloutsos, 2004 114

CMU SCS

Skip

Wavelets - construction

level 1 $d_{1,0}$ $s_{1,0}$ $d_{1,1}$ $s_{1,1}$

x_0 x_1 x_2 x_3 x_4 x_5 x_6 x_7

SIGMOD 04 Copyright: C. Faloutsos, 2004 115

CMU SCS

Skip

Wavelets - construction

level 2 $d_{2,0}$ $s_{2,0}$

$d_{1,0}$ $s_{1,0}$ $d_{1,1}$ $s_{1,1}$

x_0 x_1 x_2 x_3 x_4 x_5 x_6 x_7

SIGMOD 04 Copyright: C. Faloutsos, 2004 116

CMU SCS

Skip

Wavelets - construction

etc ...

$d_{2,0}$ $s_{2,0}$

$d_{1,0}$ $s_{1,0}$ $d_{1,1}$ $s_{1,1}$

x_0 x_1 x_2 x_3 x_4 x_5 x_6 x_7

SIGMOD 04 Copyright: C. Faloutsos, 2004 117

CMU SCS

Skip

Wavelets - construction

Q: map each coefficient on the time-freq. plane

$d_{2,0}$ $s_{2,0}$

$d_{1,0}$ $s_{1,0}$ $d_{1,1}$ $s_{1,1}$

x_0 x_1 x_2 x_3 x_4 x_5 x_6 x_7

f

t

SIGMOD 04 Copyright: C. Faloutsos, 2004 118

CMU SCS

Skip

Wavelets - construction

Q: map each coefficient on the time-freq. plane

$d_{2,0}$ $s_{2,0}$

$d_{1,0}$ $s_{1,0}$ $d_{1,1}$ $s_{1,1}$

x_0 x_1 x_2 x_3 x_4 x_5 x_6 x_7

f

t

SIGMOD 04 Copyright: C. Faloutsos, 2004 119

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(<stdin>){
    @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 "u", $diff[$i];
        $smooth[$i] = ($vals[2*$i] + $vals[2*$i + 1]) / sqrt(2);
    }
    print "\n";
    @vals = @smooth;
    $half = int($half/2);
}
print "v", $vals[0], " "; # the final, smooth component
    
```

SIGMOD 04 Copyright: C. Faloutsos, 2004 120

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, ...

SIGMOD 04 Copyright: C. Faloutsos, 2004 121

CMU SCS

Wavelets - how do they look like?

- E.g., Daubechies-4

SIGMOD 04 Copyright: C. Faloutsos, 2004 122

CMU SCS

Wavelets - how do they look like?

- E.g., Daubechies-4

SIGMOD 04 Copyright: C. Faloutsos, 2004 123

CMU SCS

Wavelets - how do they look like?

- E.g., Daubechies-4

SIGMOD 04 Copyright: C. Faloutsos, 2004 124

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 125

CMU SCS

Wavelets - Drill#1:

- Q: baritone/silence/soprano - DWT?

SIGMOD 04 Copyright: C. Faloutsos, 2004 126

CMU SCS

Wavelets - Drill#1:

- Q: baritone/soprano - DWT?

SIGMOD 04 Copyright: C. Faloutsos, 2004 127

CMU SCS

Wavelets - Drill#2:

- Q: spike - DWT?

SIGMOD 04 Copyright: C. Faloutsos, 2004 128

CMU SCS

Wavelets - Drill#2:

- Q: spike - DWT?

0.00	0.00	0.71	0.00
0.00	0.50	-0.35	0.35

SIGMOD 04 Copyright: C. Faloutsos, 2004 129

CMU SCS

Wavelets - Drill#3:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 130

CMU SCS

Wavelets - Drill#3:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 131

CMU SCS

Wavelets - Drill#3:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 132

CMU SCS

Wavelets - Drill#3:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 133

CMU SCS

Wavelets - Drill#3:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 134

CMU SCS

Wavelets - Drill#3:

- Q: DFT?

SIGMOD 04 Copyright: C. Faloutsos, 2004 135

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 136

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!)

SIGMOD 04 Copyright: C. Faloutsos, 2004 137

CMU SCS

Overall Conclusions

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 138

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)

SIGMOD 04 Copyright: C. Faloutsos, 2004 139

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 140

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)

SIGMOD 04 Copyright: C. Faloutsos, 2004 141

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 142

CMU SCS

BREAK!

SIGMOD 04 Copyright: C. Faloutsos, 2004 143

CMU SCS

Indexing and Mining Streams

Christos Faloutsos
CMU

CMU SCS

Part 3: Linear Forecasting

SIGMOD 04 Copyright: C. Faloutsos, 2004 145

- CMU SCS
- ## Outline
- Motivation
 - Similarity Search and Indexing
 - DSP
 - ➔ • Linear Forecasting
 - Bursty traffic - fractals and multifractals
 - Non-linear forecasting
 - Conclusions
- SIGMOD 04 Copyright: C. Faloutsos, 2004 146

CMU SCS

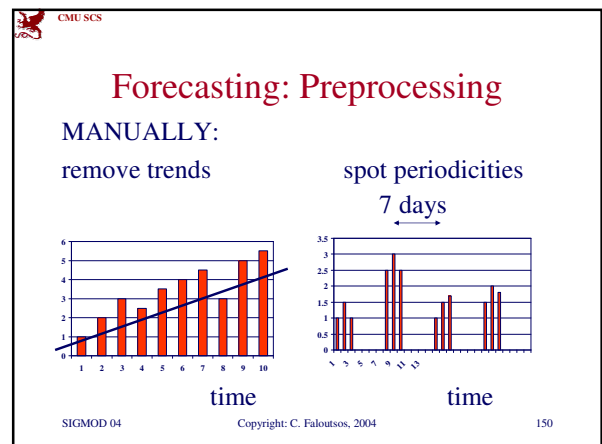
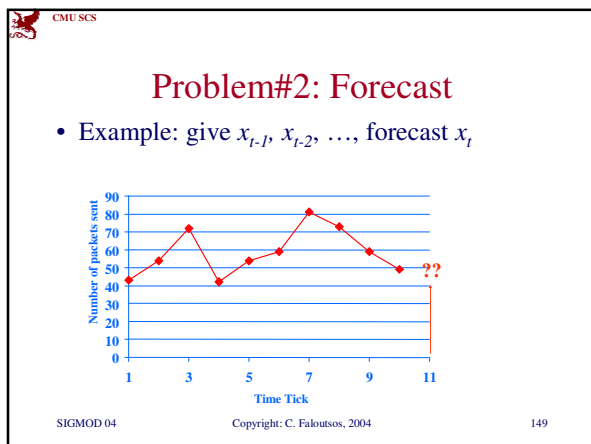
Forecasting

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

<http://www.hfac.uh.edu/MediaFutures/thoughts.html>

SIGMOD 04 Copyright: C. Faloutsos, 2004 147

- CMU SCS
- ## Outline
- Motivation
 - ...
 - ➔ • Linear Forecasting
 - Auto-regression: Least Squares; RLS
 - Co-evolving time sequences
 - Examples
 - Conclusions
- SIGMOD 04 Copyright: C. Faloutsos, 2004 148



CMU SCS

Problem#2: Forecast

- Solution: try to express x_t as a linear function of the past: x_{t-2}, x_{t-3}, \dots (up to a window of w)

Formally:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 151

CMU SCS

(Problem: Back-cast; interpolate)

- Solution - interpolate: try to express x_t as a linear function of the past AND the future: $x_{t+1}, x_{t+2}, \dots, x_{t+w_{future}}, x_{t-1}, \dots, x_{t-w_{past}}$ (up to windows of w_{past}, w_{future})

EXACTLY the same algo's

SIGMOD 04 Copyright: C. Faloutsos, 2004 152

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)')

SIGMOD 04 Copyright: C. Faloutsos, 2004 153

CMU SCS

Linear Auto Regression:

Time	Packets Sent(t)
1	43
2	54
3	72
...	...
N	??

SIGMOD 04 Copyright: C. Faloutsos, 2004 154

CMU SCS

Linear Auto Regression:

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

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 155

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 156

CMU SCS

More details:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 157

CMU SCS

More details:

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

SIGMOD 04 Copyright: 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!)

SIGMOD 04 Copyright: C. Faloutsos, 2004 159

CMU SCS

More details:

Skip

- 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
 - \mathbf{a} is the vector of the regression coefficients
 - \mathbf{X} has the N values of the w indep. variables
 - \mathbf{y} has the N values of the dependent variable

SIGMOD 04 Copyright: C. Faloutsos, 2004 160

CMU SCS

More details:

Skip

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

Ind-var1 Ind-var-w

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 161

CMU SCS

More details:

Skip

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

Ind-var1 Ind-var-w

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 162

CMU SCS

Skip

More details

- Q2: How to estimate $a_1, a_2, \dots, a_w = \mathbf{a}$?
- 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}

SIGMOD 04 Copyright: C. Faloutsos, 2004 163

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:

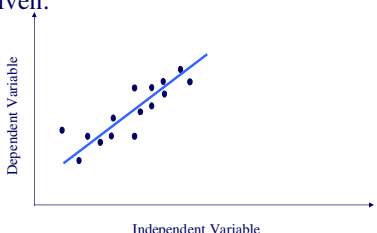
SIGMOD 04 Copyright: C. Faloutsos, 2004 164

CMU SCS

Skip

Even more details

- Given:

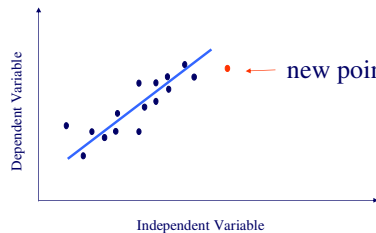


SIGMOD 04 Copyright: C. Faloutsos, 2004 165

CMU SCS

Skip

Even more details



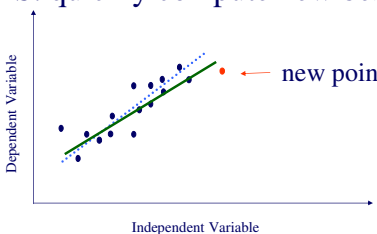
SIGMOD 04 Copyright: C. Faloutsos, 2004 166

CMU SCS

Skip

Even more details

RLS: quickly compute new best fit



SIGMOD 04 Copyright: C. Faloutsos, 2004 167

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 $O(1 \times w^2)$

$N = 10^6, w = 1-100$

SIGMOD 04 Copyright: C. Faloutsos, 2004 168

CMU SCS

Skip

Even more details

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 169

CMU SCS

Skip

Adaptability - ‘forgetting’

SIGMOD 04 Copyright: C. Faloutsos, 2004 170

CMU SCS

Skip

Adaptability - ‘forgetting’

SIGMOD 04 Copyright: C. Faloutsos, 2004 171

CMU SCS

Skip

Adaptability - ‘forgetting’

- RLS: can *trivially* handle ‘forgetting’

SIGMOD 04 Copyright: C. Faloutsos, 2004 172

CMU SCS

How to choose ‘w’?

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

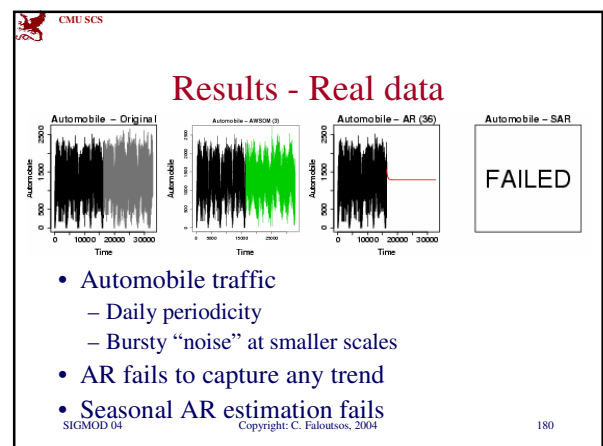
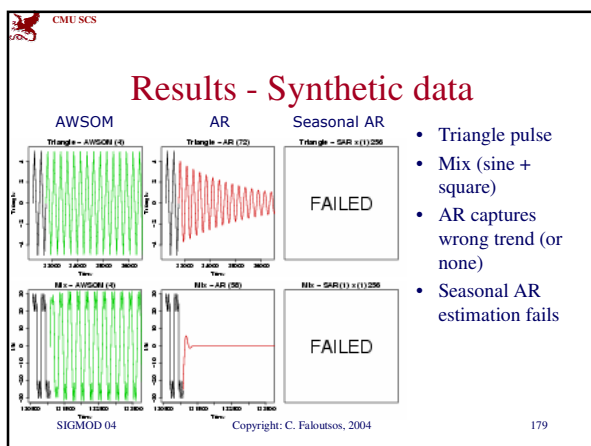
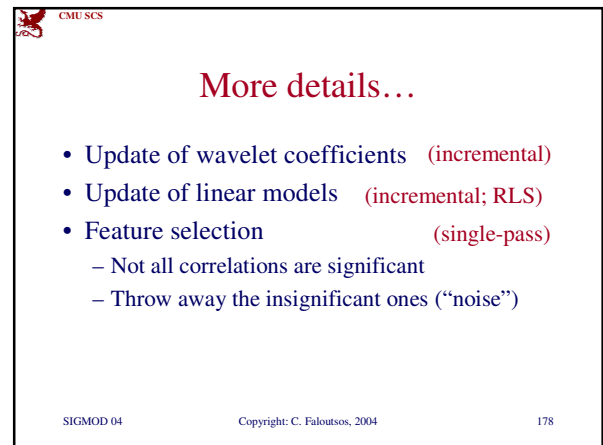
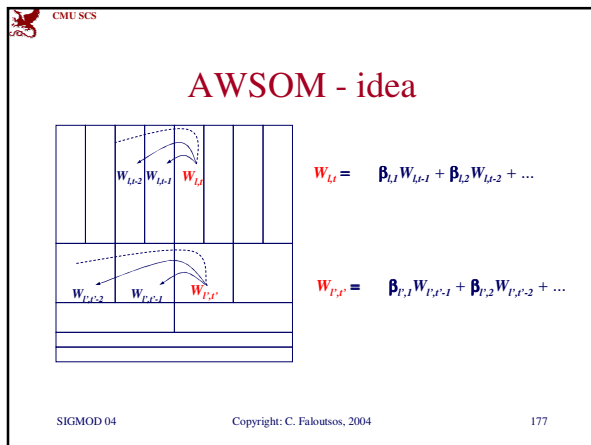
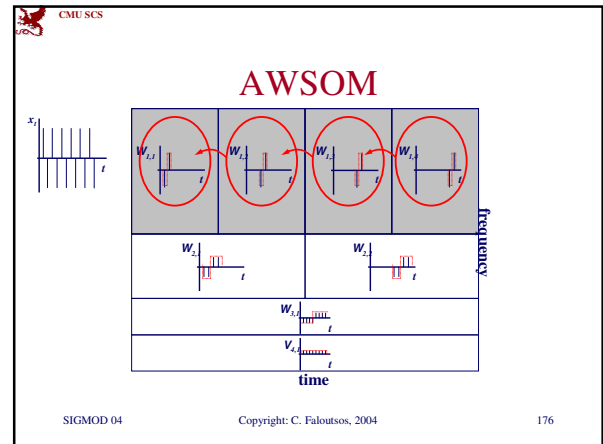
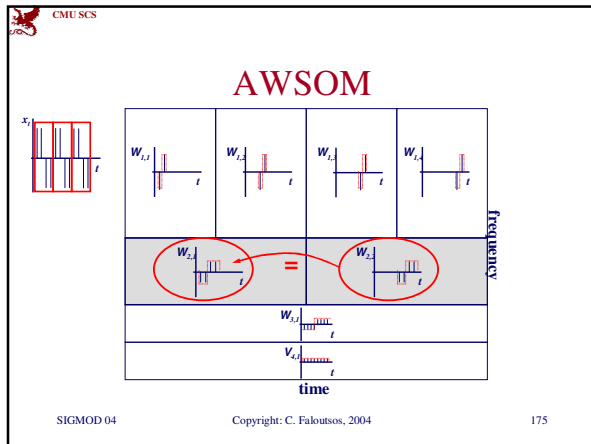
SIGMOD 04 Copyright: C. Faloutsos, 2004 173

CMU SCS

Answer:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 174



CMU SCS

Results - real data

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 181

CMU SCS

Skip

Complexity

- Model update
 - Space: $O(\lg N + mk^2) \approx O(\lg N)$
 - 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; $O(\lg N)$

SIGMOD 04 Copyright: C. Faloutsos, 2004 182

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 183

CMU SCS

Co-Evolving Time Sequences

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 184

CMU SCS

Solution:

Q: what should we do?

SIGMOD 04 Copyright: C. Faloutsos, 2004 185

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])

SIGMOD 04 Copyright: C. Faloutsos, 2004 186

CMU SCS

B.II - Time Series Analysis Outline

Skip

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 187

CMU SCS

Examples - Experiments

Skip

- 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)

SIGMOD 04 Copyright: C. Faloutsos, 2004 188

CMU SCS

Accuracy - "Modem"

Skip

Modem	AR	yesterday	MUSCLES
1	1.8	1.5	1.2
2	1.5	1.2	1.0
3	2.0	1.8	1.5
4	2.5	2.2	1.8
5	2.2	2.0	1.6
6	2.8	2.5	2.0
7	2.5	2.2	1.8
8	2.0	1.8	1.5
9	2.2	2.0	1.6
10	2.5	2.2	1.8
11	2.0	1.8	1.5
12	2.5	2.2	1.8
13	2.0	1.8	1.5
14	3.8	3.5	2.8

MUSCLES outperforms AR & "yesterday"

SIGMOD 04 Copyright: C. Faloutsos, 2004 189

CMU SCS

Accuracy - "Internet"

Skip

Streams	AR	yesterday	MUSCLES
1	0.7	0.6	0.5
2	0.7	0.6	0.5
3	0.7	0.6	0.5
4	0.7	0.6	0.5
5	0.7	0.6	0.5
6	0.7	0.6	0.5
7	0.7	0.6	0.5
8	0.7	0.6	0.5
9	0.7	0.6	0.5
10	0.7	0.6	0.5
11	0.7	0.6	0.5
12	0.7	0.6	0.5
13	1.3	1.2	1.0
14	1.3	1.2	1.0
15	1.3	1.2	1.0

MUSCLES consistently outperforms AR & "yesterday"

SIGMOD 04 Copyright: C. Faloutsos, 2004 190

CMU SCS

B.II - Time Series Analysis Outline

Skip

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 191

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)

SIGMOD 04 Copyright: C. Faloutsos, 2004 192

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/>

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.

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)

Part 4: Bursty traffic and multifractals

Outline

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

Outline

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

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 199

CMU SCS

Problem #1

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

bytes

time

SIGMOD 04 Copyright: C. Faloutsos, 2004 200

CMU SCS

Motivation

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 201

CMU SCS

Q: any ‘pattern’?

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

bytes

time

SIGMOD 04 Copyright: C. Faloutsos, 2004 202

CMU SCS

solution: self-similarity

bytes

bytes

time

time

SIGMOD 04 Copyright: C. Faloutsos, 2004 203

CMU SCS

But:

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 204

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 205

CMU SCS

Approach

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 206

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%)

SIGMOD 04 Copyright: C. Faloutsos, 2004 207

CMU SCS

binary multifractals

20 ↗ ↘ 80

SIG 208

CMU SCS

binary multifractals

20 ↗ ↘ 80

SIG 209

CMU SCS

Parameter estimation

- Q2: How to estimate the bias factor b ?

SIGMOD 04 Copyright: C. Faloutsos, 2004 210

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]

SIGMOD 04 Copyright: C. Faloutsos, 2004 211

CMU SCS

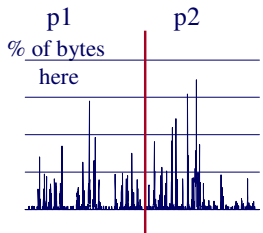
Entropy plot

- 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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 212

CMU SCS

Entropy plot

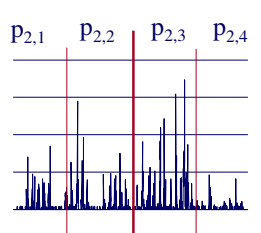


- Entropy $E(n)$ after n levels of splits
- $n=1$: $E(1) = -p_1 \log_2(p_1) - p_2 \log_2(p_2)$

SIGMOD 04 Copyright: C. Faloutsos, 2004 213

CMU SCS

Entropy plot

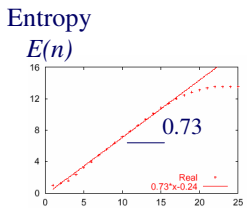


- Entropy $E(n)$ after n levels of splits
- $n=1$: $E(1) = -p_1 \log_2(p_1) - p_2 \log_2(p_2)$
- $n=2$: $E(2) = -\sum_i p_{2,i} \log_2(p_{2,i})$

SIGMOD 04 Copyright: C. Faloutsos, 2004 214

CMU SCS

Real traffic

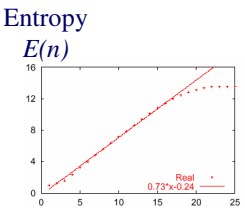


- Has linear entropy plot (\rightarrow self-similar)

SIGMOD 04 Copyright: C. Faloutsos, 2004 215

CMU SCS

Observation - intuition: Skip



intuition: slope = intrinsic dimensionality = info-bits per coordinate-bit

- unif. Dataset: slope = 1
- multi-point: slope = 0


SIGMOD 04 Copyright: C. Faloutsos, 2004 216

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?


SIGMOD 04 Copyright: C. Faloutsos, 2004 217

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


SIGMOD 04 Copyright: C. Faloutsos, 2004 218

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?


SIGMOD 04 Copyright: C. Faloutsos, 2004 219

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 = 1 bit
= E(2) - E(1) = slope!
↑ Is next MSB =0?


SIGMOD 04 Copyright: C. Faloutsos, 2004 220

CMU SCS

Skip

Entropy plot

- Repeat, for all points at same position:

Dim=0 


SIGMOD 04 Copyright: C. Faloutsos, 2004 221

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 


SIGMOD 04 Copyright: C. Faloutsos, 2004 222


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 

$0 < \text{Dim} < 1$ 

SIGMOD 04 Copyright: C. Faloutsos, 2004 223

CMU SCS

(Fractals, again)

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 224

CMU SCS


Cantor dust

- Eliminate the middle third
- Recursively!

SIGMOD 04 Copyright: C. Faloutsos, 2004 225

CMU SCS


Cantor dust



SIGMOD 04 Copyright: C. Faloutsos, 2004 226

CMU SCS


Cantor dust



SIGMOD 04 Copyright: C. Faloutsos, 2004 227

CMU SCS

Cantor dust



SIGMOD 04 Copyright: C. Faloutsos, 2004 228

CMU SCS

Cantor dust

SIGMOD 04 Copyright: C. Faloutsos, 2004 229

CMU SCS

Cantor dust

Dimensionality?
(no length; infinite # points!)
Answer: $\log_2 / \log_3 = 0.6$

SIGMOD 04 Copyright: C. Faloutsos, 2004 230

CMU SCS

Some more entropy plots:

- Poisson vs real

Poisson: slope = ~ 1 \rightarrow uniformly distributed

SIGMOD 04 Copyright: C. Faloutsos, 2004 231

CMU SCS

B-model

$E(n)$

n

- b-model traffic gives perfectly linear plot
- Lemma: its slope is $\text{slope} = -b \log_2 b - (1-b) \log_2 (1-b)$
- Fitting: do entropy plot; get slope; solve for b

SIGMOD 04 Copyright: C. Faloutsos, 2004 232

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 233

CMU SCS

Experimental setup

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 234

CMU SCS

Model validation

- Linear entropy plots

(a) Disk Traces

(b) Web Traces

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 235

CMU SCS

Web traffic - results

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

Prob($>l$)

(a) lbl-all

(b) lbl-nntp

(c) lbl-smtp

(d) lbl-ftp

Queue length distribution

How to give guarantees? (queue length l)

SIGMOD 04 Copyright: C. Faloutsos, 2004 236

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)

SIGMOD 04 Copyright: C. Faloutsos, 2004 237

CMU SCS

Conclusions

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 238

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!)

SIGMOD 04 Copyright: C. Faloutsos, 2004 239

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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 240

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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 241

CMU SCS

Part 5: chaos and non-linear forecasting

SIGMOD 04 Copyright: C. Faloutsos, 2004 242

CMU SCS

Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 243

CMU SCS

Detailed Outline

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 244

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} , ...

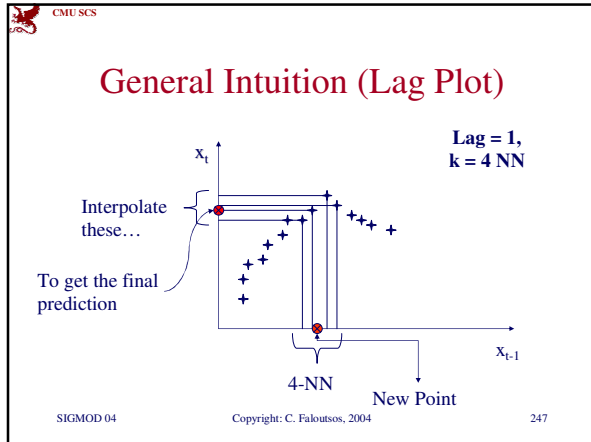
SIGMOD 04 Copyright: C. Faloutsos, 2004 245

CMU SCS

How to forecast?

- ARIMA - but: linearity assumption
- ANSWER: 'Delayed Coordinate Embedding' = Lag Plots [Sauer92]

SIGMOD 04 Copyright: C. Faloutsos, 2004 246

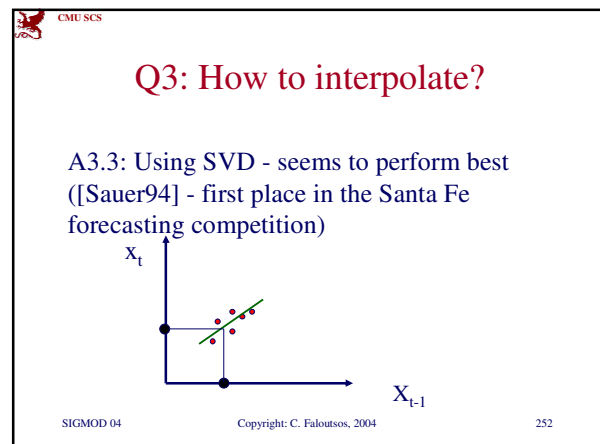


- CMU SCS
- ## Questions:
- Q1: How to choose lag L ?
 - Q2: How to choose k (the # of NN)?
 - Q3: How to interpolate?
 - Q4: why should this work at all?
- SIGMOD 04 Copyright: C. Faloutsos, 2004 248

- CMU SCS
- ## Q1: Choosing lag L
- Manually (16, in award winning system by [Sauer94])
- SIGMOD 04 Copyright: C. Faloutsos, 2004 249

- CMU SCS
- ## Q2: Choosing number of neighbors k
- Manually (typically ~ 1-10)
- SIGMOD 04 Copyright: C. Faloutsos, 2004 250

- 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?)
- SIGMOD 04 Copyright: C. Faloutsos, 2004 251



CMU SCS

Q4: Any theory behind it?

A4: YES!

SIGMOD 04 Copyright: C. Faloutsos, 2004 253

CMU SCS

Theoretical foundation

- Based on the “Takens’ Theorem” [Takens81]
- which says that **long enough** delay vectors **can do prediction**, even if there are unobserved variables in the dynamical system (= diff. equations)

SIGMOD 04 Copyright: C. Faloutsos, 2004 254

CMU SCS

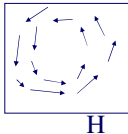
Theoretical foundation Skip

Example: Lotka-Volterra equations

$$\begin{aligned} dH/dt &= r H - a H * P \\ dP/dt &= b H * P - m P \end{aligned}$$

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, ...).

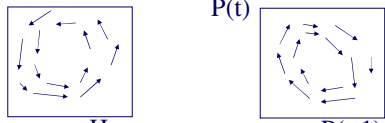


SIGMOD 04 Copyright: C. Faloutsos, 2004 255

CMU SCS

Theoretical foundation Skip

- 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**



SIGMOD 04 Copyright: C. Faloutsos, 2004 256

CMU SCS

Detailed Outline

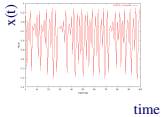
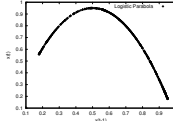
- Non-linear forecasting
 - Problem
 - Idea
 - How-to
 - ➔ – Experiments
 - Conclusions

SIGMOD 04 Copyright: C. Faloutsos, 2004 257

CMU SCS

Datasets

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

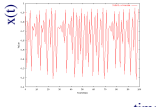
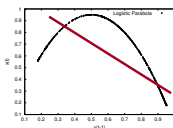
Lag-plot

SIGMOD 04 Copyright: C. Faloutsos, 2004 258

CMU SCS

Datasets

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

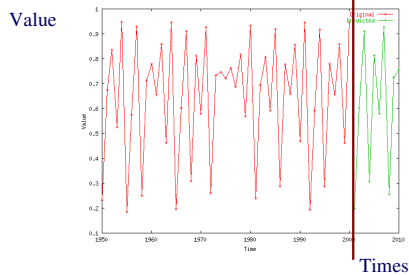
Lag-plot
ARIMA: fails

SIGMOD 04 Copyright: C. Faloutsos, 2004 259

CMU SCS

Logistic Parabola

Our Prediction from here

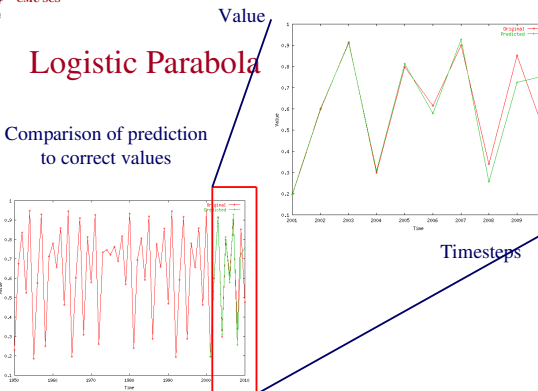


SIGMOD 04 Copyright: C. Faloutsos, 2004 260

CMU SCS

Logistic Parabola

Comparison of prediction to correct values

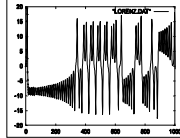
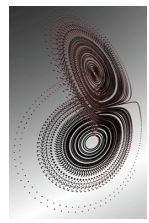


SIGMOD 04 Copyright: C. Faloutsos, 2004 261

CMU SCS

Datasets

LORENZ: Models convection currents in the air

$$\begin{aligned} dx / dt &= a(y - x) \\ dy / dt &= x(b - z) - y \\ dz / dt &= xy - cz \end{aligned}$$



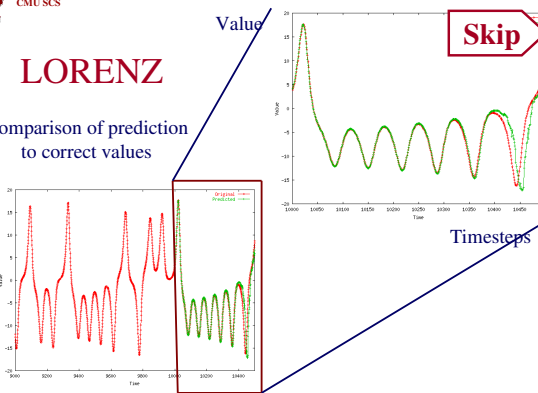
Skip

SIGMOD 04 Copyright: C. Faloutsos, 2004 262

CMU SCS

LORENZ

Comparison of prediction to correct values




Skip

SIGMOD 04 Copyright: C. Faloutsos, 2004 263

CMU SCS

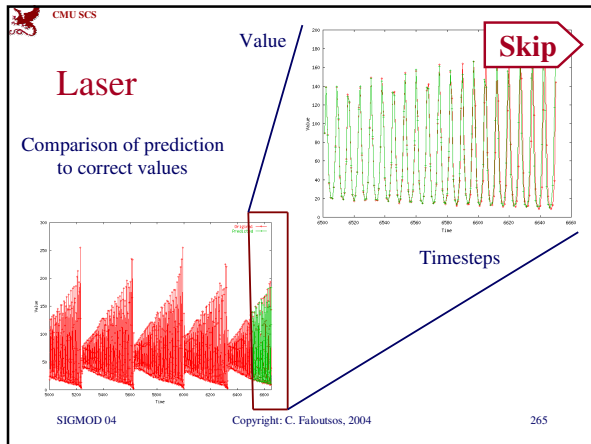
Datasets

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



Skip

SIGMOD 04 Copyright: C. Faloutsos, 2004 264



CMU SCS

Conclusions

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 266

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.

SIGMOD 04 Copyright: C. Faloutsos, 2004 267

CMU SCS

References

- Weigend, A. S. and N. A. Gershenfeld (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.)

SIGMOD 04 Copyright: C. Faloutsos, 2004 268

CMU SCS

Overall conclusions

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 269

CMU SCS

Overall conclusions

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 270

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 271

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')

SIGMOD 04 Copyright: C. Faloutsos, 2004 272

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)

SIGMOD 04 Copyright: C. Faloutsos, 2004 273

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

SIGMOD 04 Copyright: C. Faloutsos, 2004 274

CMU SCS

THANK YOU!



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

SIGMOD 04 Copyright: C. Faloutsos, 2004 275