


SCS CMU Civil and Environmental ENGINEERING

Sensor Mining at work: Principles and a Water Quality Case-Study


Jeanne M. VanBriesen (CIT, CMU)
Christos Faloutsos (SCS, CMU)

SCS CMU Civil and Environmental ENGINEERING


Thanks




Dr. Jia-Yu (Tim) Pan (CMU)



Dr. Spiros Papadimitriou (IBM)



Dr. Yasushi Sakurai (NTT)




Prof. Byoung-Kee Yi (Pohang U.)

KDD 2006
J. VanBriesen, C. Faloutsos
2

SCS CMU Civil and Environmental ENGINEERING

Acknowledgements



National Science Foundation
WHERE DISCOVERIES BEGIN

Dept. of Homeland Security

KDD 2006
J. VanBriesen, C. Faloutsos
3

SCS CMU Civil and Environmental ENGINEERING

Outline

- ➔ Motivation
- Traditional tools (wavelets, etc)
- Recent streaming tools
- Intro to water quality [Jeanne]
- Conclusions

KDD 2006
J. VanBriesen, C. Faloutsos
4

SCS CMU Civil and Environmental ENGINEERING

Problem definition

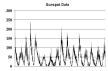
- Given: one or more sequences
 $x_1, x_2, \dots, x_t, \dots$
 $(y_1, y_2, \dots, y_p, \dots)$
 ...)
- Find
 - similar sequences; forecasts
 - patterns; clusters; outliers

KDD 2006
J. VanBriesen, C. Faloutsos
5

SCS CMU Civil and Environmental ENGINEERING

Motivation - Applications

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



KDD 2006
J. VanBriesen, C. Faloutsos
6

SCS CMU Civil and Environmental ENGINEERING

Motivation – Applications (cont'd)

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

KDD 2006 J. VanBriesen, C. Faloutsos 7

SCS CMU Civil and Environmental ENGINEERING

Motivation - Applications (cont'd)

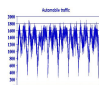
- ‘Smart house’
 - sensors monitor temperature, humidity, air quality
- video surveillance

KDD 2006 J. VanBriesen, C. Faloutsos 8

SCS CMU Civil and Environmental ENGINEERING

Motivation - Applications (cont'd)

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




KDD 2006 J. VanBriesen, C. Faloutsos 9

SCS CMU Civil and Environmental ENGINEERING

Motivation - Applications (cont'd)

- Computer systems
 - data centers (‘self-*)’)
 - web servers
 - network traffic monitoring
 - ...

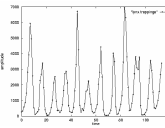


KDD 2006 J. VanBriesen, C. Faloutsos 10

SCS CMU Civil and Environmental ENGINEERING

Problem #1:

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



count lynx caught per year
(chlorine per minute;
temperature per day)

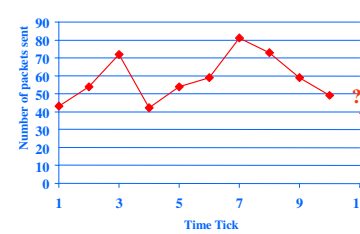
year

KDD 2006 J. VanBriesen, C. Faloutsos 11

SCS CMU Civil and Environmental ENGINEERING

Problem#2: Forecast

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



Number of packets sent

Time Tick

KDD 2006 J. VanBriesen, C. Faloutsos 12

SCS CMU Civil and Environmental ENGINEERING

Problem#2': Similarity search

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

KDD 2006 J. VanBriesen, C. Faloutsos 13

SCS CMU Civil and Environmental ENGINEERING

Problem #3:

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

KDD 2006 J. VanBriesen, C. Faloutsos 14

SCS CMU Civil and Environmental ENGINEERING

Important observations

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

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

KDD 2006 J. VanBriesen, C. Faloutsos 16

SCS CMU Civil and Environmental ENGINEERING

Important topics NOT in this tutorial:

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

KDD 2006 J. VanBriesen, C. Faloutsos 17

SCS CMU Civil and Environmental ENGINEERING

Outline

- Motivation
- Traditional tools
 - Similarity Search and Indexing
 - DSP (Digital Signal Processing)
 - Linear Forecasting
 - ICA
- Recent streaming tools
- Intro to water quality [Jeanne]
- Conclusions

KDD 2006 J. VanBriesen, C. Faloutsos 18

SCS CMU Civil and Environmental ENGINEERING

Importance of distance functions

Subtle, but **absolutely necessary**:

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

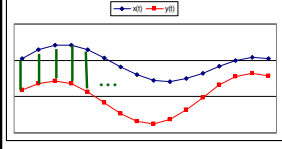
Two major families

- Euclidean and Lp norms
- Time warping and variations

KDD 2006 J. VanBriesen, C. Faloutsos 19

SCS CMU Civil and Environmental Engineering

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_∞

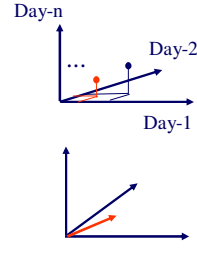
KDD 2006 J. VanBriesen, C. Faloutsos 20

SCS CMU Civil and Environmental Engineering

Observation

Euclidean distance is closely related to

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



KDD 2006 J. VanBriesen, C. Faloutsos 21

SCS CMU Civil and Environmental Engineering

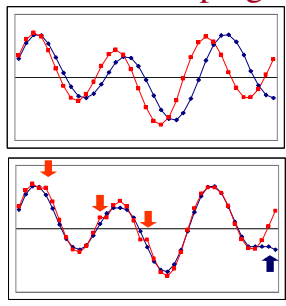
Time Warping

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

KDD 2006 J. VanBriesen, C. Faloutsos 22

SCS CMU Civil and Environmental Engineering

Time Warping



'stutters':

KDD 2006 J. VanBriesen, C. Faloutsos 23

SCS CMU Civil and Environmental Engineering

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]

KDD 2006 J. VanBriesen, C. Faloutsos 24

SCS CMU Civil and Environmental Engineering

Conclusions

Prevailing distances:

- Euclidean and
- time-warping

KDD 2006 J. VanBriesen, C. Faloutsos 25

SCS CMU Civil and Environmental ENGINEERING

Conclusions

Prevailing distances:

- Euclidean and ← probably most suitable for our setting
- time-warping

KDD 2006 J. VanBriesen, C. Faloutsos 26

SCS CMU Civil and Environmental ENGINEERING

Outline

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

KDD 2006 J. VanBriesen, C. Faloutsos 27

SCS CMU Civil and Environmental ENGINEERING

Indexing

Problem:

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

KDD 2006 J. VanBriesen, C. Faloutsos 28

SCS CMU Civil and Environmental ENGINEERING

Price vs. day (1 to 365)

distance function: by expert

KDD 2006 J. VanBriesen, C. Faloutsos 29

SCS CMU Civil and Environmental ENGINEERING

Idea: 'GEMINI'

Eg., 'find stocks similar to MSFT'

Seq. scanning: too slow

How to accelerate the search?

[Faloutsos96]

KDD 2006 J. VanBriesen, C. Faloutsos 30

SCS CMU Civil and Environmental ENGINEERING

'GEMINI' - Pictorially

S1 vs. day (1 to 365)

Sn vs. day (1 to 365)

eg., std

eg., avg

$F(S1)$

$F(Sn)$

KDD 2006 J. VanBriesen, C. Faloutsos 31

SCS CMU Civil and Environmental Engineering

GEMINI

Solution: Quick-and-dirty' filter:

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

KDD 2006 J. VanBriesen, C. Faloutsos 32

SCS CMU Civil and Environmental Engineering

Examples of GEMINI

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

KDD 2006 J. VanBriesen, C. Faloutsos 33

SCS CMU Civil and Environmental Engineering

Examples of GEMINI

Even on other-than-sequence data:

- Images (QBIC) [JIIS94]
- tumor-like shapes [VLDB96]
- video [Informedia + S-R-trees]
- automobile part shapes [Kriegel+97]

KDD 2006 J. VanBriesen, C. Faloutsos 34

SCS CMU Civil and Environmental Engineering

Indexing - SAMs

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

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

For example:

KDD 2006 J. VanBriesen, C. Faloutsos 35

SCS CMU Civil and Environmental Engineering Skip

R-trees

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

KDD 2006 J. VanBriesen, C. Faloutsos 36

SCS CMU Civil and Environmental Engineering Skip

R-trees

- eg., w/ fanout 4:

KDD 2006 J. VanBriesen, C. Faloutsos 37

SCS CMU Civil and Environmental Engineering

R-trees

Skip

- eg., w/ fanout 4:

KDD 2006 J. VanBriesen, C. Faloutsos 38

SCS CMU Civil and Environmental Engineering

R-trees - range search?

Skip

KDD 2006 J. VanBriesen, C. Faloutsos 39

SCS CMU Civil and Environmental Engineering

R-trees - range search?

Skip

KDD 2006 J. VanBriesen, C. Faloutsos 40

SCS CMU Civil and Environmental Engineering

Conclusions

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

KDD 2006 J. VanBriesen, C. Faloutsos 41

SCS CMU Civil and Environmental Engineering

Outline

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

KDD 2006 J. VanBriesen, C. Faloutsos 42

SCS CMU Civil and Environmental Engineering

Outline

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

KDD 2006 J. VanBriesen, C. Faloutsos 43

SCS CMU Civil and Environmental ENGINEERING

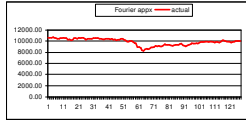
DFT and cousins

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

KDD 2006 J. VanBriesen, C. Faloutsos 44

SCS CMU Civil and Environmental ENGINEERING

DFT and stocks

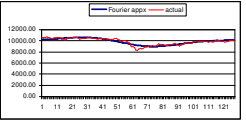


- Dow Jones Industrial index, 6/18/2001-12/21/2001

KDD 2006 J. VanBriesen, C. Faloutsos 45

SCS CMU Civil and Environmental ENGINEERING

DFT and stocks



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

KDD 2006 J. VanBriesen, C. Faloutsos 46

SCS CMU Civil and Environmental ENGINEERING

Outline

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

KDD 2006 J. VanBriesen, C. Faloutsos 47

SCS CMU Civil and Environmental ENGINEERING

SVD

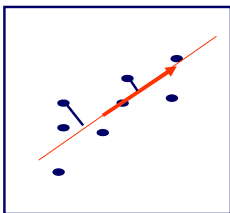
- THE optimal method for dimensionality reduction
 - (under the Euclidean metric)

KDD 2006 J. VanBriesen, C. Faloutsos 48

SCS CMU Civil and Environmental ENGINEERING

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]

KDD 2006 J. VanBriesen, C. Faloutsos 49

SCS CMU Civil and Environmental ENGINEERING

SVD

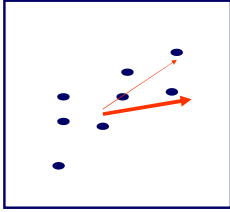
- **Extremely** useful tool
 - (also behind PageRank/google and Kleinberg’s algorithm for hubs and authorities)
- But may be slow: $O(N * M * M)$ if $N > M$
- any approximate, faster method?

KDD 2006 J. VanBriesen, C. Faloutsos 50

SCS CMU Civil and Environmental ENGINEERING

SVD shortcuts

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



KDD 2006 J. VanBriesen, C. Faloutsos 51

SCS CMU Civil and Environmental ENGINEERING

Random projections

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

KDD 2006 J. VanBriesen, C. Faloutsos 52

SCS CMU Civil and Environmental ENGINEERING

Variations of R.P.:

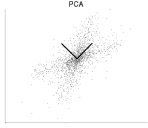
- ‘Sketches’ in time series analysis [Indyk+, VLDB 2000]

KDD 2006 J. VanBriesen, C. Faloutsos 53

SCS CMU Civil and Environmental ENGINEERING

SVD - quality

- Q: can we find better ‘hidden variables’?
- A: yes – with Independent Component Analysis (ICA) – see later

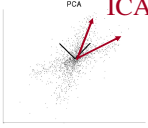


KDD 2006 J. VanBriesen, C. Faloutsos 54

SCS CMU Civil and Environmental ENGINEERING

SVD - quality

- Q: can we find better ‘hidden variables’?
- A: yes – with Independent Component Analysis (ICA) – see later



KDD 2006 J. VanBriesen, C. Faloutsos 55

SCS CMU Civil and Environmental ENGINEERING

Conclusions - Practitioner's guide

Similarity search in time sequences

- 1) establish/choose distance (Euclidean, time-warping,...)
- 2) extract features (SVD, DWT), and use an SAM (R-tree/variant) or a Metric Tree (M-tree)
- 2') for high intrinsic dimensionalities, consider sequential scan (it might win...)

KDD 2006 J. VanBriesen, C. Faloutsos 56

SCS CMU Civil and Environmental ENGINEERING

Books

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

KDD 2006 J. VanBriesen, C. Faloutsos 57

SCS CMU Civil and Environmental ENGINEERING

References

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

KDD 2006 J. VanBriesen, C. Faloutsos 58

SCS CMU Civil and Environmental ENGINEERING

References

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

KDD 2006 J. VanBriesen, C. Faloutsos 59

SCS CMU Civil and Environmental ENGINEERING

References

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

KDD 2006 J. VanBriesen, C. Faloutsos 60

SCS CMU Civil and Environmental ENGINEERING

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.

KDD 2006 J. VanBriesen, C. Faloutsos 61

SCS CMU Civil and Environmental ENGINEERING

References

- P. Indyk, N. Koudas, and S. Muthukrishnan. Identifying Representative Trends in Massive Time Series Data Sets Using Sketches. In Proc. of the 26th Int. Conf. on Very Large Data Bases, Cairo, Egypt, September 2000.
- Jolliffe, I. T. (1986). *Principal Component Analysis*, Springer Verlag.

KDD 2006 J. VanBriesen, C. Faloutsos 62

SCS CMU Civil and Environmental ENGINEERING

References

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

KDD 2006 J. VanBriesen, C. Faloutsos 63

SCS CMU Civil and Environmental ENGINEERING

References

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

KDD 2006 J. VanBriesen, C. Faloutsos 64

SCS CMU Civil and Environmental ENGINEERING

References

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

KDD 2006 J. VanBriesen, C. Faloutsos 65

SCS CMU Civil and Environmental ENGINEERING

References

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

KDD 2006 J. VanBriesen, C. Faloutsos 66

SCS CMU Civil and Environmental ENGINEERING

Traditional tools: DSP (Digital Signal Processing)

KDD 2006 J. VanBriesen, C. Faloutsos 67

SCS CMU Civil and Environmental Engineering

Outline

- Motivation
- Traditional tools
 - Similarity Search and Indexing
 - DSP (Digital Signal Processing)
 - Linear Forecasting
 - ICA
- Recent streaming tools
- Intro to water quality [Jeanne]
- Conclusions

KDD 2006 J. VanBriesen, C. Faloutsos 68

SCS CMU Civil and Environmental Engineering

Outline

- ➔ • DFT
 - Definition of DFT and properties
 - how to read the DFT spectrum
- DWT
 - Definition of DWT and properties
 - how to read the DWT scalogram

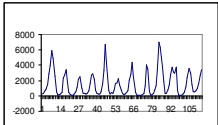
KDD 2006 J. VanBriesen, C. Faloutsos 69

SCS CMU Civil and Environmental Engineering

Introduction - Problem#1

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

count



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

year

KDD 2006 J. VanBriesen, C. Faloutsos 70

SCS CMU Civil and Environmental Engineering

What does DFT do?

A: highlights the periodicities
Powerful tool: Amplitude spectrum

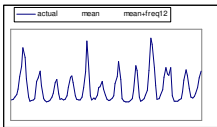
KDD 2006 J. VanBriesen, C. Faloutsos 71

SCS CMU Civil and Environmental Engineering

DFT: Amplitude spectrum

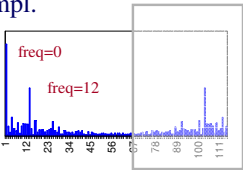
Amplitude: $A_f^2 = \text{Re}^2(X_f) + \text{Im}^2(X_f)$

count



year

Ampl.



1 12 23 34 45 56 67 78 89 100 111

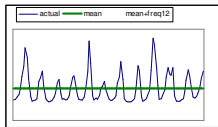
Freq.

KDD 2006 J. VanBriesen, C. Faloutsos 72

SCS CMU Civil and Environmental Engineering

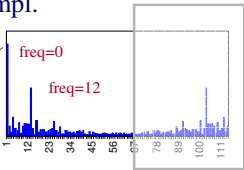
DFT: Amplitude spectrum

count



year

Ampl.



1 12 23 34 45 56 67 78 89 100 111

Freq.

KDD 2006 J. VanBriesen, C. Faloutsos 73

SCS CMU Civil and Environmental ENGINEERING

DFT: Amplitude spectrum

count

year

Ampl.

Freq.

KDD 2006 J. VanBriesen, C. Faloutsos 74

SCS CMU Civil and Environmental ENGINEERING

DFT: Amplitude spectrum

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

KDD 2006 J. VanBriesen, C. Faloutsos 75

SCS CMU Civil and Environmental ENGINEERING

DFT: Amplitude spectrum

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

KDD 2006 J. VanBriesen, C. Faloutsos 76

SCS CMU Civil and Environmental ENGINEERING

DFT: Amplitude spectrum

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

KDD 2006 J. VanBriesen, C. Faloutsos 77

SCS CMU Civil and Environmental ENGINEERING

Outline

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

KDD 2006 J. VanBriesen, C. Faloutsos 78

SCS CMU Civil and Environmental ENGINEERING

Problem #1:

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

count

year

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

KDD 2006 J. VanBriesen, C. Faloutsos 79

SCS CMU Civil and Environmental ENGINEERING

Wavelets - DWT

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

time

KDD 2006 J. VanBriesen, C. Faloutsos 80

SCS CMU Civil and Environmental ENGINEERING

Wavelets - DWT

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

time Ampl

Freq. 81

KDD 2006 J. VanBriesen, C. Faloutsos

SCS CMU Civil and Environmental ENGINEERING

Wavelets - DWT

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

time

KDD 2006 J. VanBriesen, C. Faloutsos 82

SCS CMU Civil and Environmental ENGINEERING

Wavelets - DWT

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

time

KDD 2006 J. VanBriesen, C. Faloutsos 83

SCS CMU Civil and Environmental ENGINEERING

Wavelets - DWT

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

time

KDD 2006 J. VanBriesen, C. Faloutsos 84

SCS CMU Civil and Environmental ENGINEERING

Wavelets - DWT

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

Time domain

	DFT	SWFT	DWT
freq			
time			

time

KDD 2006 J. VanBriesen, C. Faloutsos 85

SCS CMU Civil and Environmental Engineering

Haar Wavelets

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

KDD 2006J. VanBriesen, C. Faloutsos86

SCS CMU Civil and Environmental Engineering

Wavelets - construction

Skip

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

KDD 2006J. VanBriesen, C. Faloutsos87

SCS CMU Civil and Environmental Engineering

Wavelets - construction

Skip

KDD 2006J. VanBriesen, C. Faloutsos88

SCS CMU Civil and Environmental Engineering

Wavelets - construction

Skip

KDD 2006J. VanBriesen, C. Faloutsos89

SCS CMU Civil and Environmental Engineering

Wavelets - construction

Skip

KDD 2006J. VanBriesen, C. Faloutsos90

SCS CMU Civil and Environmental Engineering

Wavelets - construction

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

KDD 2006J. VanBriesen, C. Faloutsos91

SCS CMU Civil and Environmental Engineering

Skip

Wavelets - construction

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

KDD 2006 J. VanBriesen, C. Faloutsos 92

SCS CMU Civil and Environmental Engineering

Haar wavelets - code

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

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

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

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

KDD 2006 J. VanBriesen, C. Faloutsos 93

SCS CMU Civil and Environmental Engineering

Skip

Wavelets - construction

Observation1:
 '+' can be some weighted addition
 '-' is the corresponding weighted difference ('Quadrature mirror filters')

Observation2: unlike DFT/DCT, there are *many* wavelet bases: Haar, Daubechies-4, Daubechies-6, Coifman, Morlet, Gabor, ...

KDD 2006 J. VanBriesen, C. Faloutsos 94

SCS CMU Civil and Environmental Engineering

Skip

Wavelets - how do they look like?

- E.g., Daubechies-4

KDD 2006 J. VanBriesen, C. Faloutsos 95

SCS CMU Civil and Environmental Engineering

Skip

Wavelets - how do they look like?

- E.g., Daubechies-4

KDD 2006 J. VanBriesen, C. Faloutsos 96

SCS CMU Civil and Environmental Engineering

Skip

Wavelets - how do they look like?

- E.g., Daubechies-4

KDD 2006 J. VanBriesen, C. Faloutsos 97

SCS CMU Civil and Environmental ENGINEERING

Outline

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

KDD 2006 J. VanBriesen, C. Faloutsos 98

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#1:

- Q: baritone/silence/soprano - DWT?

KDD 2006 J. VanBriesen, C. Faloutsos 99

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#1:

- Q: baritone/soprano - DWT?

KDD 2006 J. VanBriesen, C. Faloutsos 100

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#2:

- Q: spike - DWT?

KDD 2006 J. VanBriesen, C. Faloutsos 101

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#2:

- Q: spike - DWT?

0.00	0.00	0.71	0.00
0.00	0.50		
	-0.35		
	0.35		

KDD 2006 J. VanBriesen, C. Faloutsos 102

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#3:

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

KDD 2006 J. VanBriesen, C. Faloutsos 103

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#3:

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

f t

KDD 2006 J. VanBriesen, C. Faloutsos 104

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#3:

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

f t

KDD 2006 J. VanBriesen, C. Faloutsos 105

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#3:

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

f t

KDD 2006 J. VanBriesen, C. Faloutsos 106

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#3:

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

f t

KDD 2006 J. VanBriesen, C. Faloutsos 107

SCS CMU Civil and Environmental ENGINEERING

Wavelets - Drill#3:

- Q: DFT?

DWT

f t

DFT

f t


KDD 2006 J. VanBriesen, C. Faloutsos 108

SCS CMU Civil and Environmental ENGINEERING


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

KDD 2006 J. VanBriesen, C. Faloutsos 109




SCS CMU




Overall Conclusions

- DFT (and 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!)

KDD 2006 J. VanBriesen, C. Faloutsos 110




SCS CMU




Overall Conclusions

- DWT: used for summarization of streams [Gilbert+01], db histograms etc

KDD 2006 J. VanBriesen, C. Faloutsos 111




SCS CMU




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)

KDD 2006 J. VanBriesen, C. Faloutsos 112




SCS CMU




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

KDD 2006 J. VanBriesen, C. Faloutsos 113




SCS CMU




Books

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

KDD 2006 J. VanBriesen, C. Faloutsos 114



SCS CMU



Additional Reading

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

KDD 2006 J. VanBriesen, C. Faloutsos 115

SCS CMU Civil and Environmental ENGINEERING

Traditional tools: Linear Forecasting

KDD 2006 J. VanBriesen, C. Faloutsos 116

SCS CMU Civil and Environmental ENGINEERING

Forecasting

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

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

KDD 2006 J. VanBriesen, C. Faloutsos 117

SCS CMU Civil and Environmental ENGINEERING

Outline

- Motivation
- Traditional tools
 - Similarity Search and Indexing
 - DSP (Digital Signal Processing)
 - Linear Forecasting
 - ICA
- Recent streaming tools
- Intro to water quality [Jeanne]
- Conclusions

KDD 2006 J. VanBriesen, C. Faloutsos 118

SCS CMU Civil and Environmental ENGINEERING

Problem#2: Forecast

- Example: give x_{t-1}, x_{t-2}, \dots , forecast x_t

KDD 2006 J. VanBriesen, C. Faloutsos 119

SCS CMU Civil and Environmental ENGINEERING

Forecasting: Preprocessing

MANUALLY:

- remove trends
- spot periodicities

KDD 2006 J. VanBriesen, C. Faloutsos 120

SCS CMU Civil and Environmental ENGINEERING

Problem#2: Forecast

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

Formally:

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

KDD 2006 J. VanBriesen, C. Faloutsos 121

SCS CMU Civil and Environmental ENGINEERING

Skip

(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

KDD 2006 J. VanBriesen, C. Faloutsos 122

SCS CMU Civil and Environmental ENGINEERING

Linear Auto Regression:

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

KDD 2006 J. VanBriesen, C. Faloutsos 123

SCS CMU Civil and Environmental ENGINEERING

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 ($S[t-1]$)

KDD 2006 J. VanBriesen, C. Faloutsos 124

SCS CMU Civil and Environmental ENGINEERING

Outline

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

KDD 2006 J. VanBriesen, C. Faloutsos 125

SCS CMU Civil and Environmental ENGINEERING

More details:

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

KDD 2006 J. VanBriesen, C. Faloutsos 126

SCS CMU Civil and Environmental ENGINEERING

More details:

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

KDD 2006 J. VanBriesen, C. Faloutsos 127

SCS CMU Civil and Environmental ENGINEERING

More details:

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

KDD 2006 J. VanBriesen, C. Faloutsos 128

SCS CMU Civil and Environmental ENGINEERING

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
 - \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

KDD 2006 J. VanBriesen, C. Faloutsos 129

SCS CMU Civil and Environmental ENGINEERING

Skip

More details:

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

Ind-var1 Ind-var-w

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

KDD 2006 J. VanBriesen, C. Faloutsos 130

SCS CMU Civil and Environmental ENGINEERING

Skip

More details:

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

Ind-var1 Ind-var-w

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

KDD 2006 J. VanBriesen, C. Faloutsos 131

SCS CMU Civil and Environmental ENGINEERING

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}

KDD 2006 J. VanBriesen, C. Faloutsos 132

SCS CMU Civil and Environmental ENGINEERING

Even more details

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

KDD 2006 J. VanBriesen, C. Faloutsos 133

SCS CMU Civil and Environmental ENGINEERING

Even more details

- Given:

Dependent Variable

Independent Variable

KDD 2006 J. VanBriesen, C. Faloutsos 134

SCS CMU Civil and Environmental ENGINEERING

Even more details

Dependent Variable

Independent Variable

new point

KDD 2006 J. VanBriesen, C. Faloutsos 135

SCS CMU Civil and Environmental ENGINEERING

Even more details

RLS: quickly compute new best fit

Dependent Variable

Independent Variable

new point

KDD 2006 J. VanBriesen, C. Faloutsos 136

SCS CMU Civil and Environmental ENGINEERING

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 $O(1 \times w^2)$

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

KDD 2006 J. VanBriesen, C. Faloutsos 137

SCS CMU Civil and Environmental ENGINEERING

Even more details

Skip

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

KDD 2006 J. VanBriesen, C. Faloutsos 138

SCS CMU Civil and Environmental ENGINEERING

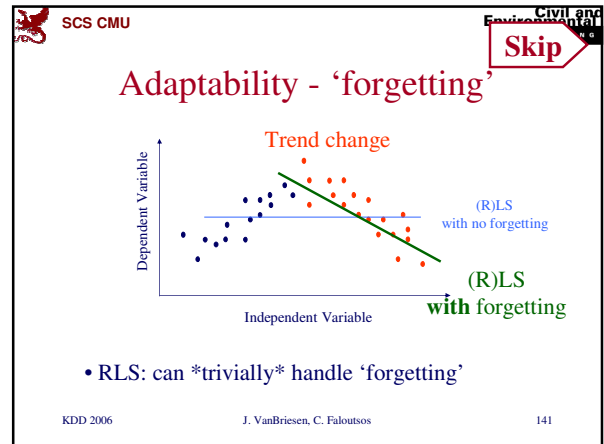
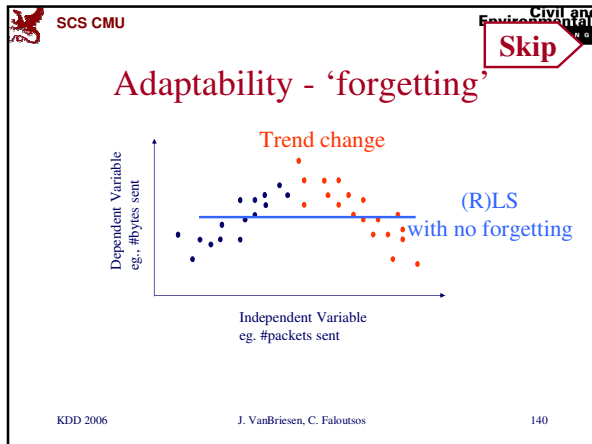
Adaptability - 'forgetting'

Skip

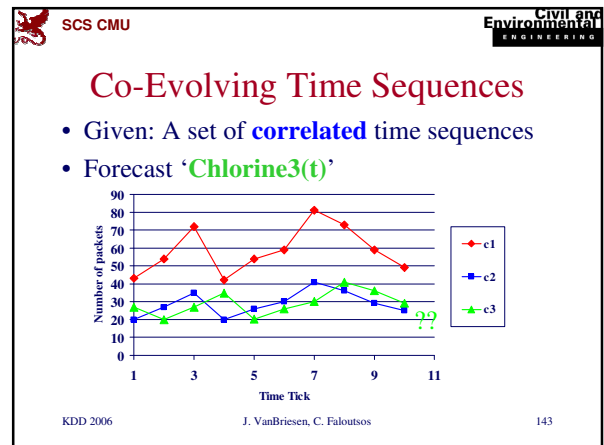
Dependent Variable
eg., #bytes sent

Independent Variable
eg., #packets sent

KDD 2006 J. VanBriesen, C. Faloutsos 139



-
- SCS CMU Civil and Environmental ENGINEERING
- ## Outline
- Motivation
 - ...
 - Linear Forecasting
 - Auto-regression: Least Squares; RLS
 - Co-evolving time sequences
 - Conclusions
- KDD 2006 J. VanBriesen, C. Faloutsos 142




SCS CMU Civil and Environmental ENGINEERING

Solution:


Q: what should we do?

KDD 2006 J. VanBriesen, C. Faloutsos 144

-
- SCS CMU Civil and Environmental ENGINEERING
- ## Solution:
- Least Squares, with
- Dep. Variable: **c3(t)**
 - Indep. Variables:
 - **c3(t-1) ... c3(t-w)**;
 - **c2(t-1) ... c2(t-w)**;
 - **c1(t-1) ,...**
 - (named: 'MUSCLES' [Yi+00])
- KDD 2006 J. VanBriesen, C. Faloutsos 145




SCS CMU




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]

KDD 2006 J. VanBriesen, C. Faloutsos 146




SCS CMU




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/>

KDD 2006 J. VanBriesen, C. Faloutsos 147




SCS CMU




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.

KDD 2006 J. VanBriesen, C. Faloutsos 148



SCS CMU



Additional Reading

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

KDD 2006 J. VanBriesen, C. Faloutsos 149