




# Mining and Forecasting of Big Time-series Data

Yasushi Sakurai (Kumamoto University)  
 Yasuko Matsubara (Kumamoto University)  
 Christos Faloutsos (Carnegie Mellon University)


<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 1

# Roadmap

- Motivation
- **Similarity search, pattern discovery and summarization** **Part 1** 
- Non-linear modeling and forecasting **Part 2**
- Extension of time-series data: tensor analysis **Part 3**

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 2





Part 1

# Similarity search, pattern discovery and summarization

Yasushi Sakurai (Kumamoto University)  
 Yasuko Matsubara (Kumamoto University)  
 Christos Faloutsos (Carnegie Mellon University)



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 3

# Part 1 - Roadmap

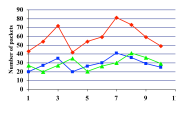
- ➔ Motivation
  - Similarity Search and Indexing
  - Feature extraction
  - Linear forecasting
  - Streaming pattern discovery
  - Automatic mining

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 4






# Problem definition

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




<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 5

# Motivation - Applications

- Financial, sales, economic series
- Medical
  - reactions to new drugs
  - elderly care
  - ECG ('physionet.org')



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 6

Kumamoto U CMU CS

## EEG - epilepsy

from wikipedia

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 7

Kumamoto U CMU CS

## Motivation - Applications (cont' d)


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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 8

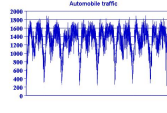
Kumamoto U CMU CS

## Motivation - Applications (cont' d)

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



Tokyo Gate Bridge



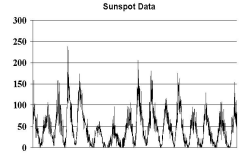
Automobile traffic

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 9

Kumamoto U CMU CS

## Motivation - Applications (cont' d)

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



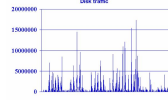
Sunspot Data

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 10

Kumamoto U CMU CS

## Motivation - Applications (cont' d)

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




Disk traffic

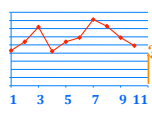
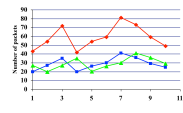
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 11

Kumamoto U CMU CS

## Wish list



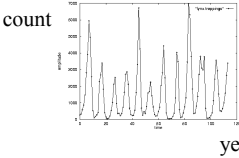
- Problem 1: find patterns/rules
- Problem 2: forecast
  - Problem 2': similarity search
- Problem 3: find patterns/rules/forecast, with many time sequences

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 12

**Problem #1:**

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

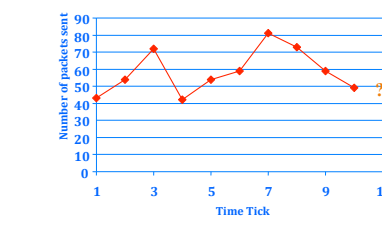


lynx caught per year  
 (packets per day;  
 temperature per day)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 13

**Problem#2: Forecast**

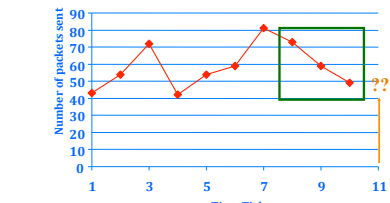
Given  $x_t, x_{t-1}, \dots$ , forecast  $x_{t+1}$



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 14

**Problem#2' : Similarity search**

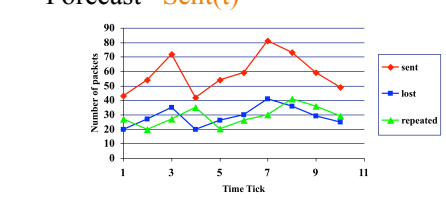
Eg., Find a 3-tick pattern, similar to the last one



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 15

**Problem #3:**

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

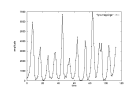


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 16

**Important observations**

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

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



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 17

**Important topics NOT in this tutorial:**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 18

Kumamoto U CMU CS

## Roadmap

- Motivation
- ➔ • Similarity Search and Indexing
- Feature extraction
- Linear forecasting
- Streaming pattern discovery
- Automatic mining

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 19

Kumamoto U CMU CS

## Roadmap

- Motivation
- Similarity Search and Indexing
  - ➔ – distance functions: Euclidean; Time-warping
  - indexing
- Feature extraction
- Linear forecasting
- Streaming pattern discovery
- Automatic mining

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 20

Kumamoto U CMU CS

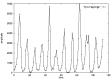
## Importance of distance functions

Subtle, but **absolutely necessary**:

- A ‘must’ for similarity indexing
  - (-> forecasting)
- A ‘must’ for clustering

Two major families

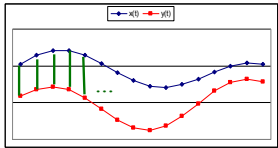
- Euclidean and Lp norms
- time warping and variations



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 21

Kumamoto U CMU CS

## Euclidean and Lp



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

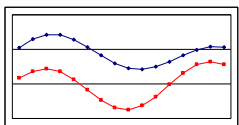
$$L_p(\vec{x}, \vec{y}) = \sum_{i=1}^n |x_i - y_i|^p$$

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

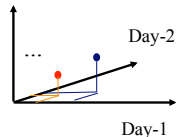
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 22

Kumamoto U CMU CS

## Observation #1



- Time sequence
  - > n-d vector

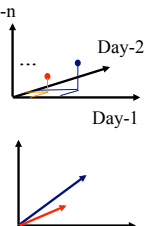


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 23

Kumamoto U CMU CS

## Observation #2

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



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 24

**Time Warping**

- allow accelerations - decelerations
  - (with or w/o penalty)
- THEN compute the (Euclidean) distance (+ penalty)
- related to the string-editing distance
- fast search methods [Yi+98] [Keogh+02] [Sakurai+05]

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 25

**Time Warping**

‘stutters’:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 26

**Time Warping DETAILS**

Q: how to compute it?  
 A: dynamic programming

$$D(i, j) = \text{cost to match prefix of length } i \text{ of first sequence } x \text{ with prefix of length } j \text{ of second sequence } y$$

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 27

**Time Warping DETAILS**

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 28

**Time Warping**

• Time warping matrix & optimal path:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 29

**Time Warping**

• Time warping matrix & optimal path:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 30

**Time Warping - variations**

- Time warping matrix & optimal path:

At most  $k$  stutters:  
Sakoe-Chiba band

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 31

**Time Warping - variations**

- Time warping matrix & optimal path:

At most  $x\%$  stutters:  
Itakura parallelogram

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 32

**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]

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 33

**A variation: Uniform axis scaling**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 34

**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]

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 35

**More distance functions.**

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

**On The Marriage of  $L_p$ -norms and Edit Distance, [Lei Chen, Raymond T. Ng](#), VLDB' 04**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 36

Kumamoto U CMU CS

## Conclusions

- Prevailing distances:
  - Euclidean and
  - time-warping

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 37

Kumamoto U CMU CS

## Roadmap

- Motivation
- Similarity Search and Indexing
  - distance functions: Euclidean; Time-warping
  - ➔ – indexing
- Feature extraction
- Linear forecasting
- Streaming pattern discovery
- Automatic mining

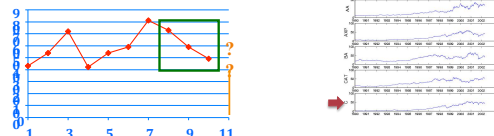
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 38

Kumamoto U CMU CS

## Indexing

Problem 2':

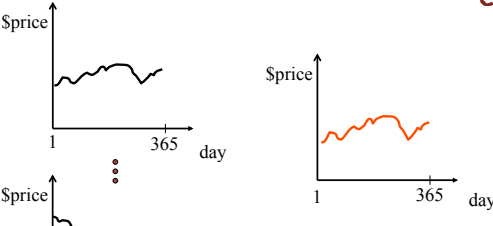
- given a set of time sequences,
- find the ones similar to a desirable query sequence



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 39

Kumamoto U CMU CS

## Indexing



distance function: by expert  
(Euclidean; DTW; ...)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 40

Kumamoto U CMU CS

## Idea: 'GEMINI'

Eg., 'find stocks similar to MSFT'

Seq. scanning: too slow

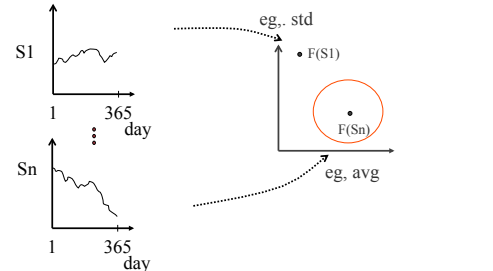
How to accelerate the search?

[Faloutsos96]

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 41

Kumamoto U CMU CS

## 'GEMINI' - Pictorially



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 42

**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' – R-tree, etc)
- discard false alarms

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 43

**Examples of GEMINI**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 44

**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:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 45

**R-trees**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 46

**R-trees**

- eg., w/ fanout 4:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 47

**R-trees**

- eg., w/ fanout 4:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 48



**R-trees - range search?**


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 49

**R-trees - range search?**

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 50

**Recent work**

- Rakthanmanon+ [kdd' 12]: EXCELLENT Software, the UCR Suite for ultrafast subsequence search
- Zoumpatianos+ [sigmod' 14]: ADS+, exploratory analysis
- Camerra+ [KAIS' 14]: iSAX2+, indexing for bulk loading



Eamonn Keogh (UCR)

*Searching and Mining Trillions of Time Series Subsequences under Dynamic Time Warping, T. Rakthanmanon, B. Campana, A. Mueen, G. Batista, B. Westover, Q. Zhu, J. Zakaria, E. Keogh, KDD' 12*

*Indexing for Interactive Exploration of Big Data Series, K. Zoumpatianos, S. Idreos, T. Palpanas, SIGMOD' 14*

*Beyond One Billion Time Series: Indexing and Mining Very Large Time Series Collections with iSAX2+, A. Camerra, J. Shieh, T. Palpanas, T. Rakthanmanon, E. Keogh, KAIS 2014*

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 52

**Conclusions**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 52

**Roadmap**

- Motivation
- Similarity Search and Indexing
- ➔ Feature extraction
- Linear forecasting
- Streaming pattern discovery
- Automatic mining

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 53

**Roadmap**

- Motivation
- Similarity Search and Indexing
- Feature extraction
  - ➔ – DFT, DWT, DCT (data independent)
  - SVD, ICA (data independent)
  - MDS, FastMap
- Linear forecasting
- Streaming pattern discovery
- Automatic mining

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 54

Kumamoto U CMU CS

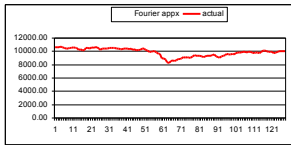
## DFT and cousins

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 55

Kumamoto U CMU CS

## DFT and stocks

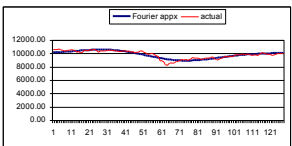


- Dow Jones Industrial index, 6/18/2001-12/21/2001

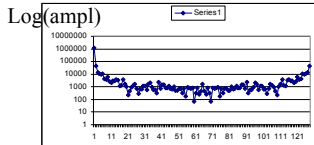
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 56

Kumamoto U CMU CS

## DFT and stocks



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

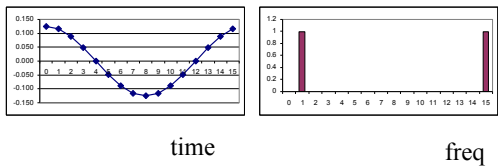


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 57

Kumamoto U CMU CS

## DFT (and DWT)

- Many more details, soon



time freq

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 58

Kumamoto U CMU CS

## Roadmap

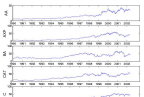
- Motivation
- Similarity Search and Indexing
- Feature extraction
  - DFT, DWT, DCT (data independent)
  - ➡ – SVD, ICA (data independent)
  - MDS, FastMap
- Linear forecasting
- Streaming pattern discovery
- Automatic mining

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 59

Kumamoto U CMU CS

## SVD

- THE optimal method for dimensionality reduction
  - (under the Euclidean metric)
- Given: many time sequences
- Find: the latent ('hidden') variables



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos

**SVD**

Two (equivalent) interpretations:

- Geometric (each sequence  $\rightarrow$  point in T-d space)
- Matrix algebra ( $N \times T$  matrix)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 61

**Singular Value Decomposition (SVD)**

- SVD ( $\sim$ LSI  $\sim$ KL  $\sim$ PCA  $\sim$  spectral analysis...) – Geometric interpretation

day2

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 62

**Reminder:**

- SVD  $\rightarrow$  matrix factorization: finds blocks

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 63

**SVD – matrix interpretation**

- SVD  $\rightarrow$  matrix factorization: finds blocks

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 64

**SVD**

- **Extremely** useful tool
  - (also behind PageRank/google and Kleinberg's algorithm for hubs and authorities)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 65

**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?

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 66

**SVD shortcuts**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 67

**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])

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 68

**SVD & improvement**

- Q: Can we do even better?
- A: sometimes, yes – by shooting for sparsity

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 69

**Independent Component Analysis (ICA)**

- PCA sometimes misses essential features – PCA vs. ICA

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 70

**A.k.a.: BSS = cocktail party problem**

**Find hidden variables**

- Untangle two sound sources

= "blind source separation"  
 • unknown sources,  
 • unknown mixing

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos #71

**ICA**

- Why not PCA

Source

Mix

Sequence #1 (Sources #1 & #3) Sequence #2 (Sources #2 & #3) Sequence #3 (Mix of all 3 sources)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 72

**ICA**

- Why not PCA

PCA: PC1, PC2, PC3

ICA: IC1, IC2, IC3

ICA recognizes the components successfully and separately

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 73

**Hidden variables**

Alcoa, American Express, Boeing, Caterpillar, Citi Group, Dow Jones Industrial Average

Find common hidden variables, and weights.

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 74

**Hidden variables**

- Find hidden variables [Pan+04]

Caterpillar, Intel

B1,CAT, B1,INTC, B2,CAT, B2,INTC

Hidden variable 1, Hidden variable 2

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 75

**Hidden variables**

- Find hidden variables [Pan+04]

Caterpillar, Intel

0.94, 0.63, 0.03, 0.64

"Hidden variable 1", "Hidden variable 2"

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 76

**Hidden variables**

- Find hidden variables [Pan+04]

Caterpillar, Intel

0.94, 0.63, 0.03, 0.64

"General trend", "Internet bubble"

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 77

**Hidden variables**

- Local component analysis [Sakurai+11]

Original sequence, Anomaly spikes, Weekly pattern, Daily pattern

(b) Weekly pattern (WindMine), (c) Daily pattern (WindMine)

(d) Weekly pattern (PCA), (e) Daily pattern (PCA)

PCA: failed

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 78

**Motivation: Find hidden variables**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 79

**Motivation: Find patterns in data**

- Motion capture data (broad jumps)

Left Knee  
Right Knee  
Energy exerted

Take-off  
Landing

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 80

**Motivation: Find patterns in data**

- Best SVD axis: not always meaningful!

Right Knee  
Left Knee

Auto-Split Bases  
PCA Bases

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 81

**Motivation: Find patterns in data**

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

Right Knee  
Left Knee

60:1  
1:1

Auto-Split Bases  
PCA Bases

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 82

**Problem formulation**

- Given  $n$  data items, each has  $m$  attributes
- Find the  $m$  hidden variables and the  $m$  bases

$$\begin{bmatrix} X_{11} & X_{12} & \dots & X_{1m} \\ \dots & \dots & \dots & \dots \\ X_{n1} & X_{n2} & \dots & X_{nm} \end{bmatrix} = \begin{bmatrix} H_{11} & H_{12} & \dots & H_{1m} \\ \dots & \dots & \dots & \dots \\ H_{n1} & H_{n2} & \dots & H_{nm} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} & \dots & B_{1m} \\ \dots & \dots & \dots & \dots \\ B_{m1} & B_{m2} & \dots & B_{mm} \end{bmatrix}$$

$X=HB$

Samples of the  $m$ -th hidden variable

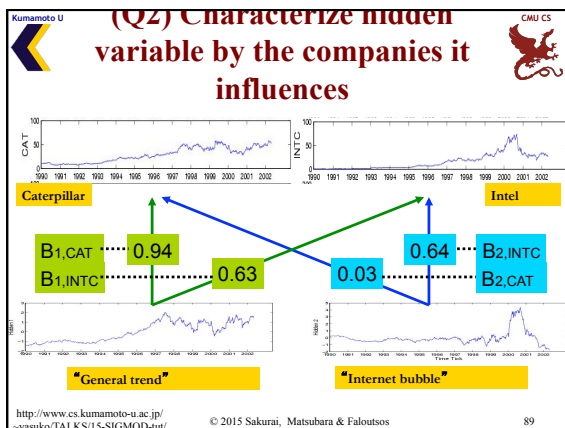
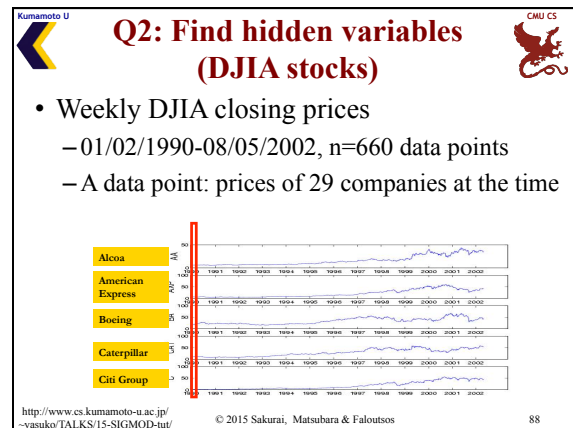
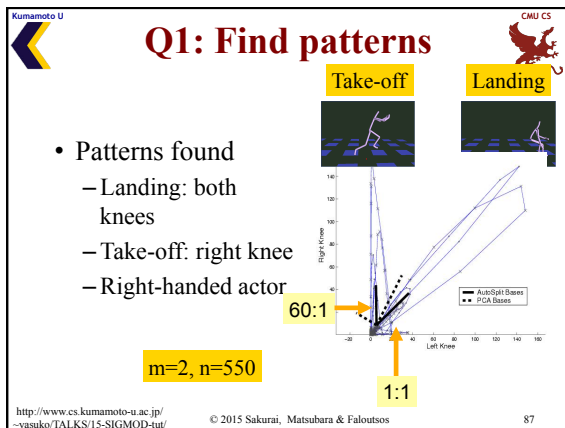
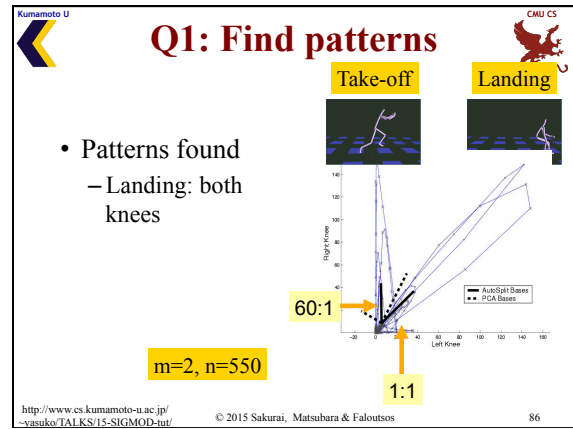
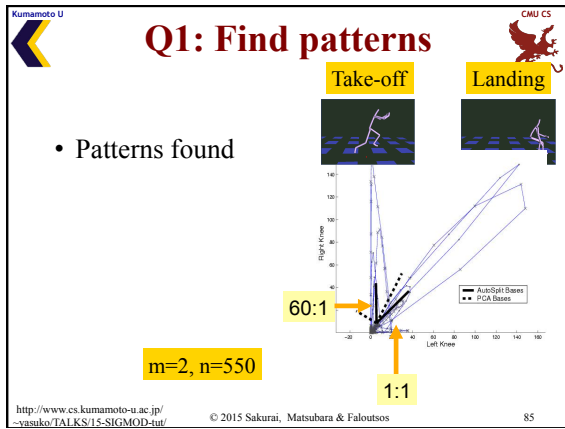
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 83

**Formulation: (Q1) Find patterns in data**

Left Knee  
Right Knee

Basis 1

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 84



**Companies related to hidden variable 1**

	B <sub>1j</sub>	
	Highest	Lowest
Caterpillar	0.938512	AT&T 0.021885
Boeing	0.911120	WalMart 0.624570
MMM	0.906542	Intel 0.638010
Coca Cola	0.903858	Home Depot 0.647774
Du Pont	0.900317	Hewlett-Packard 0.658768


Hidden 1

"General trend"

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 90

**Citation**

- *AutoSplit: Fast and Scalable Discovery of Hidden Variables in Stream and Multimedia Databases*, **Jia-Yu Pan**, Hiroyuki Kitagawa, Christos Faloutsos and Masafumi Hamamoto, PAKDD 2004, Sydney, Australia.
- *WindMine: Fast and Effective Mining of Web-click Sequences*, Yasushi Sakurai, Lei Li, Yasuko Matsubara, Christos Faloutsos, SDM 2011, Mesa, Arizona.



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos #91

**Roadmap**

- Motivation
- Similarity Search and Indexing
- Feature extraction
  - DFT, DWT, DCT (data independent)
  - SVD, ICA (data independent)
  - ➔ – MDS, FastMap
- Linear forecasting
- Streaming pattern discovery
- Automatic mining

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 92

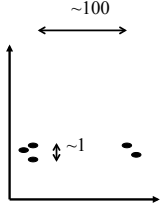
**MDS / FastMap**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 93

**MDS/FastMap**

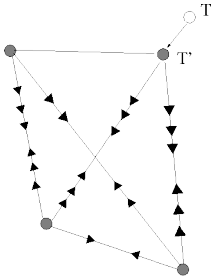
	01	02	03	04	05
01	0	1	1	100	100
02	1	0	1	100	100
03	1	1	0	100	100
04	100	100	100	0	1
05	100	100	100	1	0



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 94

**MDS**

Multi Dimensional Scaling



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 95

**FastMap**

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

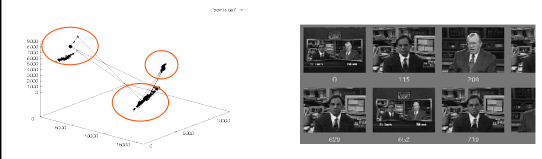
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 96



Kumamoto U CMU CS

## FastMap: Application

VideoTrails [Kobla+97]




scene-cut detection (about 10% errors)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 97

Kumamoto U CMU CS

## Variations


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



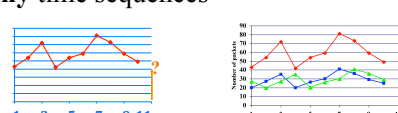
<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 98

Kumamoto U CMU CS

## Wish list



- Problem 1: find patterns/rules
- Problem 2: **forecast**
- ✓ • Problem 2': **similarity search**
- Problem 3: find patterns/rules/forecast, with **many** time sequences



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 99

Kumamoto U CMU CS

## Conclusions - Practitioner's guide

Similarity search in time sequences

- 1) establish/choose distance (Euclidean, time-warping,...)
- 2) extract features (SVD, ICA, 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...)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 100

Kumamoto U CMU CS

## 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)


<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 101

Kumamoto U CMU CS

## 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/>


<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 102

Kumamoto U 

## 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.


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 103

Kumamoto U 

## 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.


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 104

Kumamoto U 

## References

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


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos Part2.1 #105

Kumamoto U 

## References

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


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 106

Kumamoto U 

## 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.

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 107

Kumamoto U 

## 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.

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 108

Kumamoto U CMU CS

## References

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos #109

Kumamoto U CMU CS

## 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.

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos #110

Kumamoto U CMU CS

## References

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos #111

Kumamoto U CMU CS

## References

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos #112

Kumamoto U CMU CS


## Roadmap

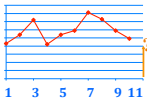
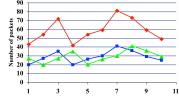
- Motivation
- Similarity Search and Indexing
- Feature extraction
- ➔ Linear forecasting
- Streaming pattern discovery
- Automatic mining

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos #113

Kumamoto U CMU CS

## Wish list

- Problem 1: find patterns/rules 
- ➔ Problem 2: **forecast**
  - Problem 2': **similarity search**
- Problem 3: find patterns/rules/forecast, with **many** time sequences

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos #114

## Forecasting

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

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



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 115

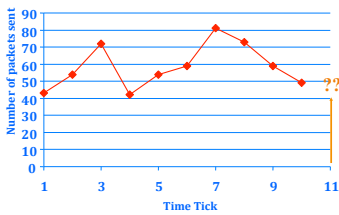
## Roadmap

- Motivation
- Similarity Search and Indexing
- Feature extraction
- Linear forecasting
  - ➔ – Auto-regression: Least Squares; RLS
  - Co-evolving time sequences
- Streaming pattern discovery
- Automatic mining

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 116

## Problem: Forecasting

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

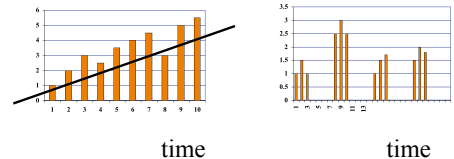


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 117

## Forecasting: Preprocessing

MANUALLY:  
remove trends  
periodicities

spot  
7 days



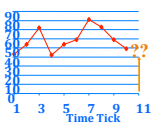
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 118

## Problem: 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:

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

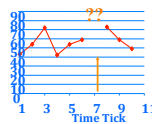


(if we **know** it is a non-linear model, see Part 2)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 119

## (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



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 120

**Linear Regression: idea**

patient	weight	height
1	27	43
2	43	54
3	54	72
...	...	...
N	25	??

Body height

Body weight

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 121

**Linear Auto Regression:**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 122

**Linear Auto Regression:**

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

'lag-plot'

Number of packets sent (t)

Number of packets sent (t-1)

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 123

**More details:**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 124

**More details:**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 125

**More details:**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 126

**More details: 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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 127

**More details: DETAILS**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 128

**More details: DETAILS**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 129

**More details: 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  $\mathbf{y}$

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 130

**Even more details: 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:

**[Yi+00] Byoung-Kee Yi et al.: Online Data Mining for Co-Evolving Time Sequences, ICDE 2000.**

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 132

**Even more details: DETAILS**

- Given:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 132

**Even more details**

- Given:

Dependent Variable

Independent Variable

new point

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 133

**Even more details**

Recursive Least Squares (RLS):  
quickly compute new best fit

Dependent Variable

Independent Variable

new point

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 134

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

49,000,000 ←→ 49

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 135

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

49,000,000 ←→ 49

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

**RLS: GREAT for streams**

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 136

**Even more detail DETAILS**

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

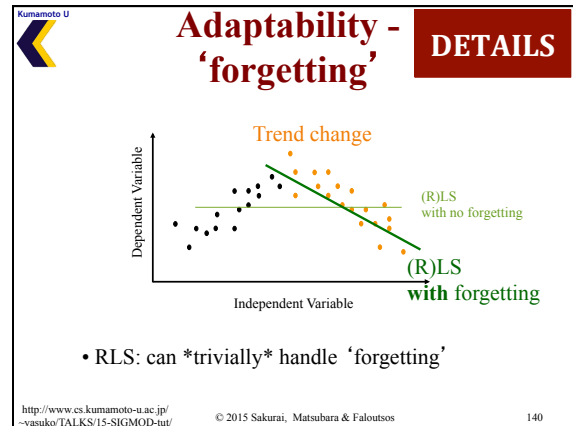
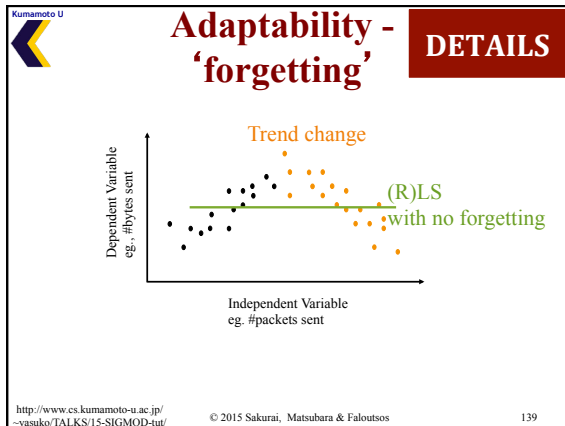
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 137

**Adaptability - 'forgetting' DETAILS**

Dependent Variable  
eg., #bytes sent

Independent Variable  
eg., #packets sent

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 138



**How to choose 'w'?**

- Quick & dirty answer:  $w=1$  or  $w=2$
- Better answer: Model selection (say, with BIC or MDL – see later)
- Even better answer: **multi-scale windows** [Papadimitriou+, vldb2003]

Spiros Papadimitriou, Anthony Brockwell and Christos Faloutsos *Adaptive, Hands-Off Stream Mining* VLDB 2003, Berlin, Germany, Sept. 2003

**How to choose 'w'?**

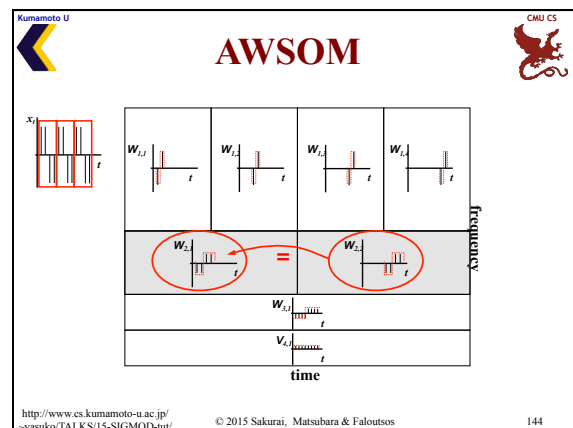
- goal: capture arbitrary periodicities
- with NO human intervention
- on a semi-infinite stream

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 142

**Answer:**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 143





**AWSOM**

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 145

**AWSOM - idea**

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 146

**More details...**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 147

**Results - Synthetic data**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 148

**Results - Real data**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 149

**Results - real data**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 150

Kumamoto U CMU CS

## Complexity

Skip

- 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;  $O(\lg N)$

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 151

Kumamoto U CMU CS

## Roadmap

- Motivation
- Similarity Search and Indexing
- Feature extraction
- Streaming pattern discovery
- Linear forecasting
  - Auto-regression: Least Squares; RLS
- ➔ Co-evolving time sequences
- Automatic mining

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 152

Kumamoto U CMU CS

## Co-Evolving Time Sequences

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 153

Kumamoto U CMU CS

## Solution:

Q: what should we do?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 154

Kumamoto U CMU CS

## 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])



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 155

Kumamoto U CMU CS

## Practitioner's guide



- AR(IMA) methodology: prevailing method for linear forecasting
- Brilliant method of Recursive Least Squares for fast, incremental estimation.
- See [Box-Jenkins]

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 156

Kumamoto U  **Resources: software and urls** 



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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 157

Kumamoto U  **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.

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 158

Kumamoto U  **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)



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 159

**Part 1**  

**Similarity search,  
pattern discovery and  
summarization**



Yasushi Sakurai (Kumamoto University)  
Yasuko Matsubara (Kumamoto University)  
Christos Faloutsos (Carnegie Mellon University)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 160

Kumamoto U  **Roadmap** 

- Motivation
- Similarity Search and Indexing
- Feature extraction
  - ➡ – DFT, DWT, DCT (data independent)
  - SVD, ICA (data independent)
  - (MDS, FastMap)
- Streaming pattern discovery

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 161

Kumamoto U  **Roadmap** 

- DFT
  - ➡ – Definition of DFT and properties
  - how to read the DFT spectrum
- DWT
  - Definition of DWT and properties
  - how to read the DWT scalogram

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos 162

**Wish list**

- Problem 1: find patterns/rules
- Problem 2: **forecast**
- Problem 2': **similarity** search
- Problem 3: find patterns/rules/forecast, with **many** time sequences

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 163

**Introduction - Problem#1**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 164

**DFT: definition DETAILS**

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

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

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 165

**DFT: definition**

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

```
x = [1,2,1,2];
X = Fourier[x];
Plot[ Abs[X] ];
```

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 166

**DFT: Amplitude spectrum**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 167

**DFT: examples** Skip

flat

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 168

Kumamoto U CMU CS

## DFT: examples

Skip

Low frequency sinusoid

time freq

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 169

Kumamoto U CMU CS

## DFT: examples

Skip

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

time freq

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 170

Kumamoto U CMU CS

## DFT: examples

Skip

- Higher freq. sinusoid

time freq

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 171

Kumamoto U CMU CS

## DFT: examples

Skip

examples

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 172

Kumamoto U CMU CS

## DFT: examples

Skip

examples

Ampl. Freq.

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 173

Kumamoto U CMU CS

## Roadmap

- DFT
  - Definition of DFT and properties
  - ➡ – how to read the DFT spectrum
- DWT
  - Definition of DWT and properties
  - how to read the DWT scalogram

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 174

**DFT: Amplitude spectrum**

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

count

year

Ampl.

Freq.

freq=0

freq=12

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 175

**DFT: Amplitude spectrum**

count

year

Ampl.

Freq.

freq=0

freq=12

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 176

**DFT: Amplitude spectrum**

count

year

Ampl.

Freq.

freq=0

freq=12

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 177

**DFT: Amplitude spectrum**

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

count

year

Ampl.

Freq.

freq=0

freq=12

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 178

**DFT: Amplitude spectrum**

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

count

year

Ampl.

Freq.

freq=0

freq=12

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 179

**DFT: Amplitude spectrum**

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

count

year

Ampl.

Freq.

freq=0

freq=12

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 180

**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)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 181

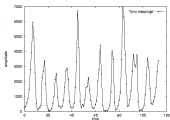
**Roadmap**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 182

**Problem #1:**

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

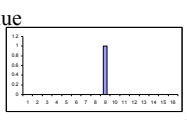


lynx caught per year  
(packets per day;  
virus infections per month)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 183

**Wavelets - DWT**

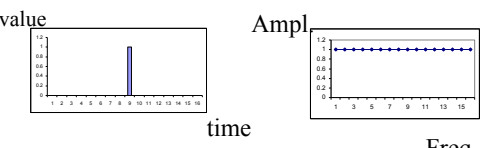
- DFT is great - but, how about compressing a spike?



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 184

**Wavelets - DWT**

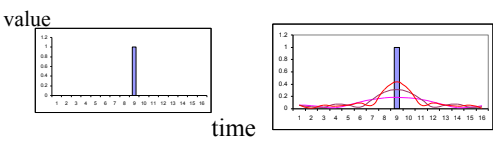
- DFT is great - but, how about compressing a spike?
- A: Terrible - all DFT coefficients needed!



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 185

**Wavelets - DWT**

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



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 186

**Wavelets - DWT**

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

value

time

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 187

**Wavelets - DWT**

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

freq

time

value

time

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 188

**Wavelets - DWT**

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

**'Multi-scale windows'**: brilliant idea that we'll see several times in this tutorial (BRAID, TriMine, etc)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 189

**Wavelets - DWT**

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

**Multi-scale windows**

Time domain

	DFT	SWFT	DWT
freq			
time			

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 190

**Haar Wavelets**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 191

**Wavelets - construc DETAILS**

x0 x1 x2 x3 x4 x5 x6 x7

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 192



**Wavelets - construc DETAILS**

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 193

**Wavelets - construc DETAILS**

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

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 194

**Wavelets - construc DETAILS**

etc ...

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

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 195

**Wavelets - construc DETAILS**

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

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

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 196

**Wavelets - construc DETAILS**

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

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

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 197

**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
# haarpl <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 "d", $vals[0], "n"; # the final, smooth component
    
```

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 198

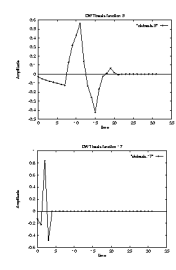
**Wavelets - construc DETAILS**

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, ...

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 199

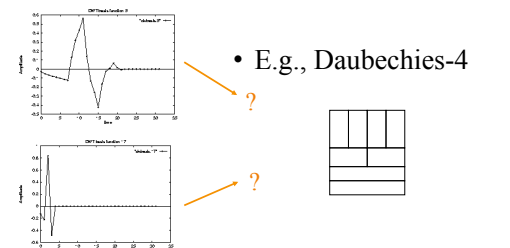
**Wavelets - how do they look like?**



- E.g., Daubechies-4

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 200

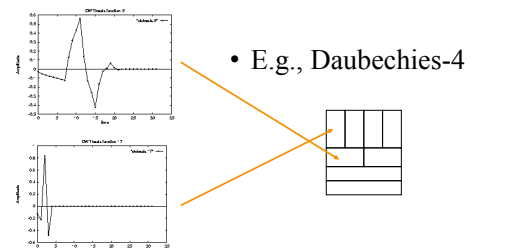
**Wavelets - how do they look like?**



- E.g., Daubechies-4

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 201

**Wavelets - how do they look like?**



- E.g., Daubechies-4

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 202

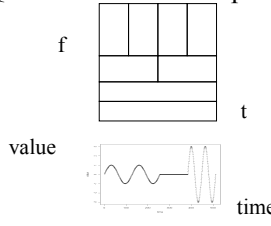
**Roadmap**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 203

**Wavelets - Drill#1:**

- Q: baritone/silence/soprano - DWT?



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 204

**Wavelets - Drill#1:**

- Q: baritone/silence/soprano - DWT?

Diagram showing frequency (f) vs time (t) with shaded cells and a corresponding waveform plot labeled 'value' vs 'time'.

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 205

**Wavelets - Drill#2:**

- Q: spike - DWT?

Diagram showing frequency (f) vs time (t) and a corresponding waveform plot with a single sharp spike.

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 206

**Wavelets - Drill#2:**

- Q: spike - DWT?

Diagram showing frequency (f) vs time (t) with numerical values and a corresponding waveform plot with a spike.

0.00	0.00	0.71	0.00
0.00	0.50		
	-0.35		
	0.35		

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 207

**Wavelets - Drill#3:**

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

Diagram showing frequency (f) vs time (t) and a corresponding waveform plot with periodic oscillations and a spike.

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 208

**Wavelets - Drill#3:**

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

Diagram showing frequency (f) vs time (t) with a shaded band and a corresponding waveform plot with periodic oscillations and a spike.

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 209

**Wavelets - Drill#3:**

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

Diagram showing frequency (f) vs time (t) with a shaded band and a corresponding waveform plot with periodic oscillations and a spike.

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 210

**Wavelets - Drill#3:**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 211

**Wavelets - Drill#3:**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 212

**Wavelets - Drill#3:**

- Q: DFT?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 213

**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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 214

**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!)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 215

**Conclusions**

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 216

Kumamoto U 

## Part 1 - Roadmap

- Motivation
- Sim. Search and Indexing { Euclidean/DTW + Feature extraction + R-trees
- Feature extraction { DFT, DWT (SVD, ICA)
- Linear forecasting { AR, RLS

---

- Streaming pattern discovery
- Automatic mining


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 217

Kumamoto U 

## 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)


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 218

Kumamoto U 

## 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


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 219

Kumamoto U 

## 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)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 220

Kumamoto U 

## 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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 221