



Mining Big Time-series Data on the Web


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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 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/16-WWW-tut/> © 2016 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/16-WWW-tut/> © 2016 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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 4



Motivation - Applications

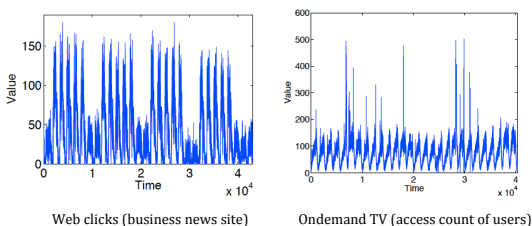
- Web online activities
 - Web access logs
 - Search volume
 - Online reviews
- IoT device data
- Medical, healthcare data

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 5



Motivation - Applications

- Web access logs



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 6

Motivation - Applications

- Web search volume from Google trends

Compare Search terms: Internet of Things

Interest over time

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 7

Motivation - Applications

- IoT (Internet of Things) device data
 - Civil/automobile infrastructure
 - Bridge vibrations [Oppenheim+02]
 - Road conditions / traffic monitoring
 - Environmental data (air/water pollutant monitoring)

Automobile traffic

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 8

Motivation - Applications

- Medical (epidemic) time-series data e.g., measles cases in the U.S.

Count x 10⁷

Year (Weekly)

Shocks, e.g., 1941

Yearly periodicity

Vaccine effect

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 9

Wish list

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 10

Problem #1

Given: time-series data (e.g., #clicks over time)

Find: patterns, periodicities, and/or compress

Original web-click sequence

Weekday component

Weekend component

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 11

Problem #2

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

Number of packets sent

Time Tick

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 12

Problem #3

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

Time Tick	sent	lost	repeated
1	40	20	20
2	50	25	25
3	70	30	30
4	40	25	25
5	55	20	20
6	60	30	30
7	80	35	35
8	70	30	30
9	55	25	25
10	45	20	20
11	40	20	20

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 13

Important observations

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

- For forecasting, we need
 - patterns/rules/models
 - 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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 14

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 15

Roadmap

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 16

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 17

Importance of distance functions

Subtle, but **absolutely necessary**:

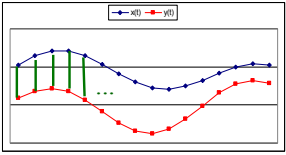
- A 'must' for similarity search, indexing and clustering

Two major families

- Euclidean and Lp norms
- Time warping and variations

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 18

Euclidean and Lp



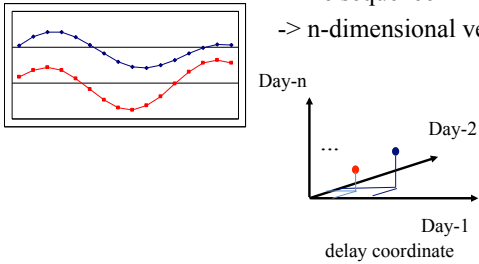
$$D(\vec{x}, \vec{y}) = \sum_{i=1}^n (x_i - y_i)^2$$

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

- L_1 : city-block = Manhattan
- L_2 = Euclidean
- L_∞

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 19

Observation #1

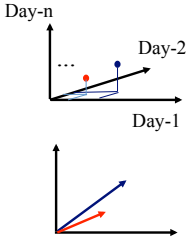


- Time sequence
-> n-dimensional vector

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 20

Observation #2

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



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 21

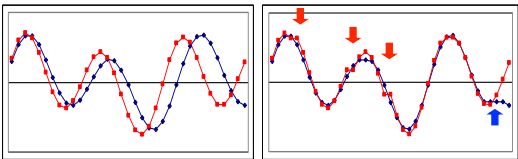
Time Warping

- 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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 22

Time Warping

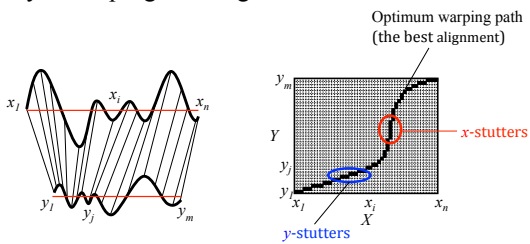
- Allow sequences to be stretched along the time axis
 1. minimize the distance of sequences
 2. insert 'stutters' into a sequence
 3. THEN compute the (Euclidean) distance



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 23

Time Warping

Q: how to compute it?
A: dynamic programming



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 24

Time Warping DETAILS

Q: how to compute it?
 A: dynamic programming

$$X = \{x_1, x_2, \dots, x_i\}, Y = \{y_1, y_2, \dots, y_j\}$$

$$D_{dw}(X, Y) = f(n, m)$$

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 25

Time Warping

- Time warping matrix & optimal path:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 26

Time Warping

- Time warping matrix & optimal path:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 27

Time Warping - variations

- Time warping matrix & optimal path:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 28

Time Warping - variations

- Time warping matrix & optimal path:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 29

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 30

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 31

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 32

Related work

- 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
- Rakthanmanon+ [kdd’12]: EXCELLENT Software, the UCR Suite for ultrafast subsequence search

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 33

Conclusions

- Prevailing distances:
 - Euclidean and
 - time-warping

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 34

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 35

Indexing

- Given a set of time sequences,
- Find the ones similar to a desirable query sequence

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 36

Indexing

Price

1 365 day

Price

1 365 day

Price

1 365 day

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 37

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 38

'GEMINI' - Pictorially

S1

1 365 day

Sn

1 365 day

eg., std

eg., avg

feature vectors

$F(S1)$

$F(Sn)$

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 39

GEMINI

Solution: Quick-and-dirty' filter:

- extract d features (numbers, eg., avg., etc.)
- map into a point in the d -dimensional 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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 40

Examples of GEMINI

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 41

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 42

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 43

R-trees

- eg., w/ fanout 4:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 44

R-trees

- eg., w/ fanout 4:

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 45

R-trees - range search?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 46

R-trees - range search?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 47

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 48

Roadmap

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 49

Roadmap

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 50

DFT: definition

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

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

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 51

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 52

DFT: examples

- Flat

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 53

DFT: examples

- Low frequency sinusoid

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 54

DFT: examples

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

time freq

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 55

DFT: examples

- Higher freq. sinusoid

time freq

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 56

DFT: examples

- Examples

time

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 57

DFT: examples

- Examples

time

Ampl.

Freq.

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 58

DFT: Amplitude spectrum

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

count

year

Freq.

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 59

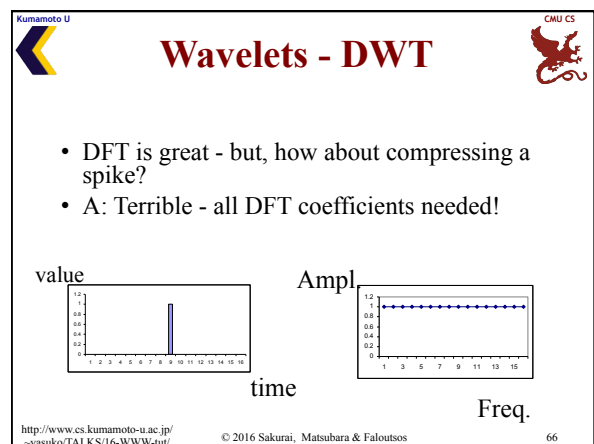
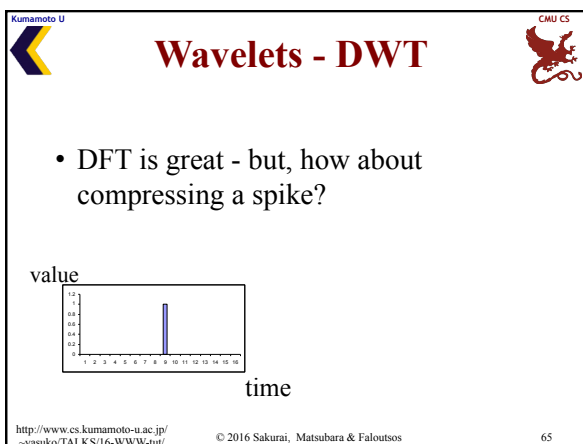
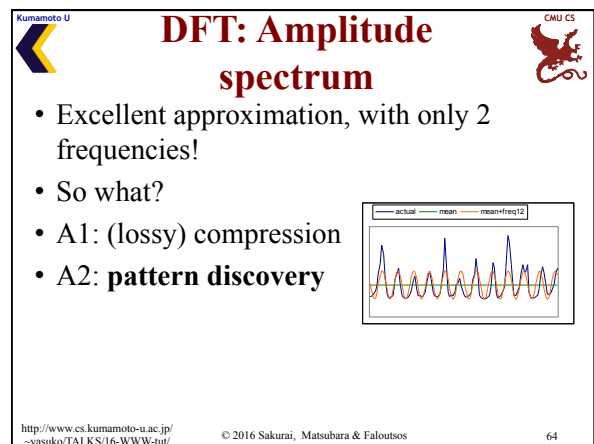
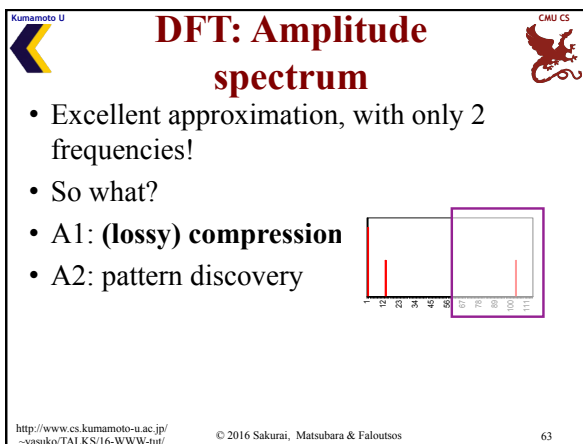
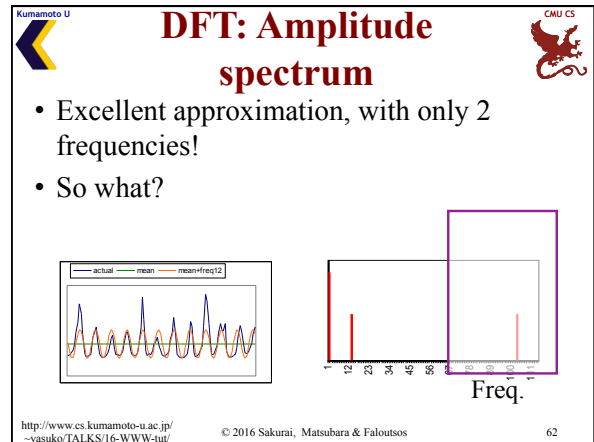
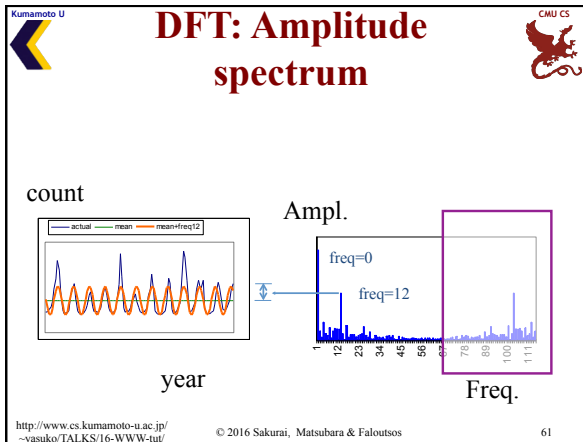
DFT: Amplitude spectrum

count

year

Freq.

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 60



Wavelets - DWT

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

value

time

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 67

Wavelets - DWT

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

value

time

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 68

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 69

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 70

Wavelets - DWT

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

Time domain

Multi-scale windows

freq

DFT

SWFT

DWT

time

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 71

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 72

Wavelets - construc DETAILS

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 73

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 74

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 75

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 76

Wavelets - construc DETAILS

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

f

t

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

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 77

Wavelets - construc DETAILS

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

f

t

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

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

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 78

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 79

Wavelets - how do they look like?

- E.g., Daubechies-4

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 80

Wavelets - how do they look like?

- E.g., Daubechies-4

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 81

Wavelets - Drill#1:

- Q: baritone/silence/soprano - DWT?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 82

Wavelets - Drill#1:

- Q: baritone/silence/soprano - DWT?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 83

Wavelets - Drill#2:

- Q: spike - DWT?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 84

Wavelets - Drill#2:

- Q: spike - DWT?

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 85

Wavelets - Drill#3:

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 86

Wavelets - Drill#3:

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 87

Wavelets - Drill#3:

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 88

Wavelets - Drill#3:

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 89

Wavelets - Drill#3:

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 90

Wavelets - Drill#3:

- Q: DFT?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 91

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'

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 92

DFT & DWT: conclusions

- **DFT** 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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 93

DFT & DWT: conclusions

- **DWT**: multi-resolution
 - very suitable for self-similar traffic
 - used for summarization of streams [Gilbert+01], db histograms, etc
- **DFT&DWT**: powerful tools for **compression, pattern detection** in real signals
 - included in math packages (matlab, 'R', mathematica, ... - often in spreadsheets!)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 94

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 95

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 96

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 97

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 98

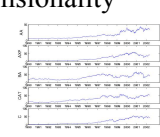
Roadmap

- 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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 99

SVD

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



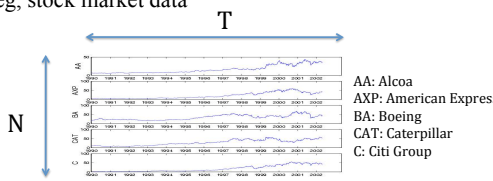
<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 100

SVD

Two (equivalent) interpretations:

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

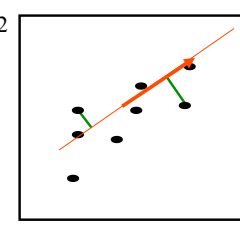
eg, stock market data



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 101

Singular Value Decomposition (SVD)

- SVD (~LSI ~ KL ~ PCA ~ spectral analysis...) – Geometric interpretation



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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 102

SVD – matrix interpretation

- SVD -> matrix factorization: finds blocks

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 103

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 104

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 105

SVD shortcuts

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 106

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 107

SVD & improvement

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 108

Independent Component Analysis (ICA)

- PCA (or SVD) sometimes misses essential features
- PCA vs. ICA

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 109

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 110

ICA

- Why not PCA

Source #1 Source #2 Source #3

Mix

Sequence #1 Sequence #2 Sequence #3

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 111

ICA

- Why not PCA

Source #1 Source #2 Source #3

Mix

Sequence #1 Sequence #2 Sequence #3

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 112

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 113

Hidden variables

- Local component analysis [Sakurai+11]

Original sequence

Anomaly spikes

Weekly pattern

Daily pattern

(b) Weekly pattern (WindMine)

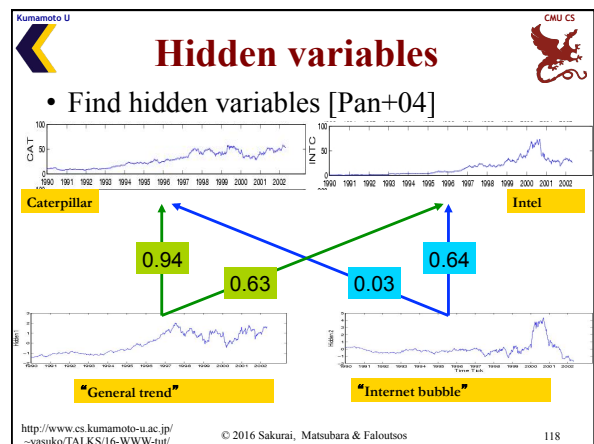
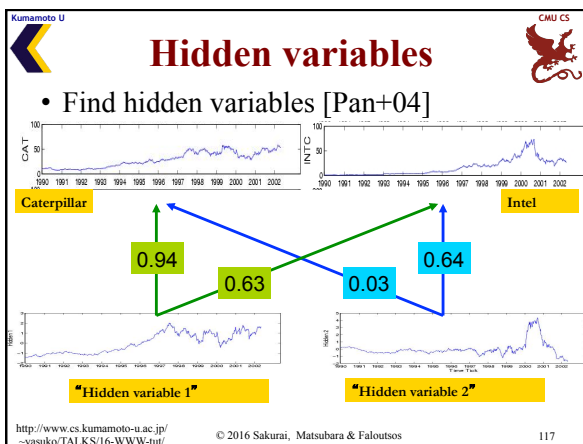
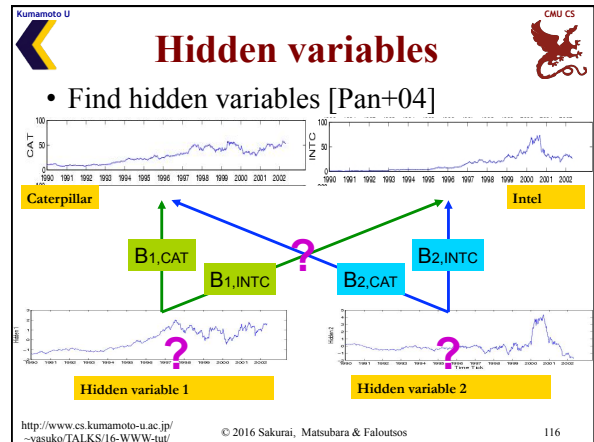
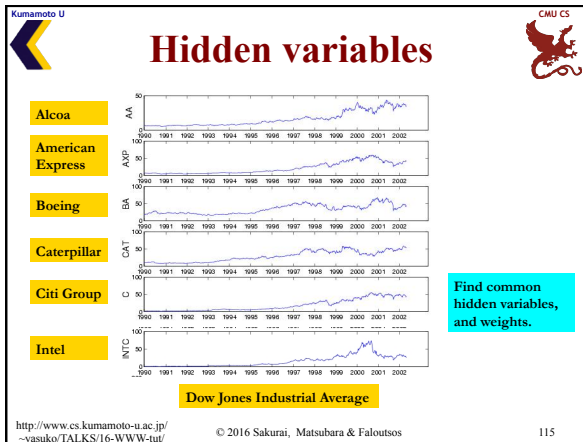
(c) Daily pattern (WindMine)

(d) Weekly pattern (PCA)

(e) Daily pattern (PCA)

PCA: failed

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 114



Motivation: Find hidden variables

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 119

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 120

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 121

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 122

MDS/FastMap

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 123

MDS

Multi Dimensional Scaling

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 124

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 125

FastMap: Application

VideoTrails [Kobla+97]


scene-cut detection (about 10% errors)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 126

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/16-WWW-tut/>
© 2016 Sakurai, Matsubara & Faloutsos
127

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/16-WWW-tut/>
© 2016 Sakurai, Matsubara & Faloutsos
128

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/16-WWW-tut/>
© 2016 Sakurai, Matsubara & Faloutsos
129

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/16-WWW-tut/>
© 2016 Sakurai, Matsubara & Faloutsos
130

Kumamoto U CMU CS

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/16-WWW-tut/>
© 2016 Sakurai, Matsubara & Faloutsos
131



Kumamoto U CMU CS

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/16-WWW-tut/>
© 2016 Sakurai, Matsubara & Faloutsos
132

Kumamoto U CMU CS






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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 133

Kumamoto U CMU CS






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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 134

Kumamoto U CMU CS






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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 135

Kumamoto U CMU CS






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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 136

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 137

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 138

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 139

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 140

Kumamoto U CMU CS


Roadmap

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

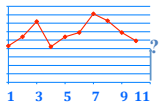
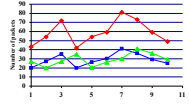
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 141

Kumamoto U CMU CS

Wish list



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


http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 142

Kumamoto U CMU CS

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 143

Kumamoto U CMU CS

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 144

Problem: Forecasting

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 145

Forecasting: Preprocessing

MANUALLY:
remove trends
periodicities

spot 7 days

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 146

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:

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

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 147

(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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 148

Background: Linear Regression

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

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 149

Linear Auto Regression:

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 150

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 ($S[t-1]$)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 151

More details:

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 152

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 153

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 154

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 155

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

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

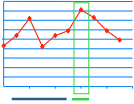
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 156

More details: DETAILS

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

Ind-var1 Ind-var-w

time



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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 157

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 158

Even more details: DETAILS

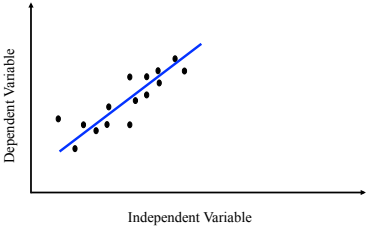
- Q3: Can we estimate \mathbf{a} incrementally?
- A3: Yes, with the brilliant, classic method of 'Recursive Least Squares' (RLS) (see, e.g., [Yi+00], for details) - pictorially:

[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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 160

Even more details

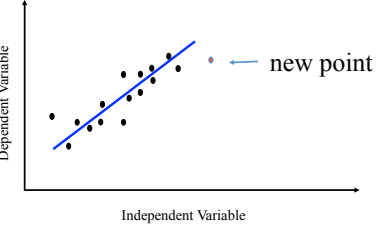
- Given:



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 160

Even more details

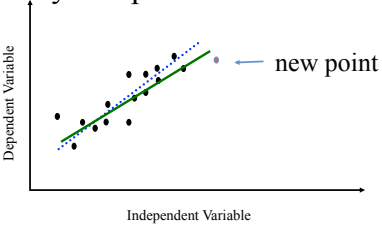
- Given:



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 161

Even more details

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



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 162

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 \longleftrightarrow 49

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 163

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 \longleftrightarrow 49

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

RLS: GREAT for streams

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 164

Even more detail DETAILS

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 165

Adaptability - 'forgetting' DETAILS

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 166

Adaptability - 'forgetting' DETAILS

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 167

Adaptability - 'forgetting' DETAILS

- RLS: can *trivially* handle 'forgetting'

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 168

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 170

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 171

AWSOM

The diagram shows a time series x_t on the left. Below it, a grid represents the decomposition of the signal into wavelet levels. The top row shows $W_{i,t}$ for $i=1, 2, 3, 4$. The bottom row shows $V_{i,t}$ for $i=1, 2$. A shaded region highlights the relationship $W_{2,t} = W_{1,t} + V_{2,t}$, with an equals sign between the two terms.

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 172

AWSOM

The diagram shows a time series x_t on the left. Below it, a grid represents the decomposition of the signal into wavelet levels. The top row shows $W_{i,t}$ for $i=1, 2, 3, 4$. The bottom row shows $V_{i,t}$ for $i=1, 2$. The relationship $W_{2,t} = W_{1,t} + V_{2,t}$ is highlighted with a shaded region and an equals sign.

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 173

AWSOM - idea

The diagram shows a time series x_t on the left. Below it, a grid represents the decomposition of the signal into wavelet levels. The top row shows $W_{i,t}$ for $i=1, 2, 3, 4$. The bottom row shows $V_{i,t}$ for $i=1, 2$. The relationship $W_{2,t} = W_{1,t} + V_{2,t}$ is highlighted with a shaded region and an equals sign.

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

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

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 174

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 175

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 176

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 177

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 178

Complexity

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 179

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 180

Kumamoto U CMU CS

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 181

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 182

Kumamoto U CMU CS

Solution:

Q: what should we do?
 A: 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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 183

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 184

Kumamoto U CMU CS

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/>

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 185

Kumamoto U CMU CS

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 186

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/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 187

Outline

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 188

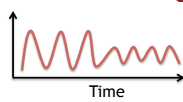
Stream mining

- Applications
 - Sensor monitoring
 - Network analysis
 - Financial and/or business transaction data
 - Web access and media service logs
 - Moving object tracking
 - Industrial manufacturing

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 189

Stream mining

- Requirements
 - **Fast**
high performance and quick response
 - **Nimble**
low memory consumption, single scan
 - **Accurate**
good approximation for pattern discovery and feature extraction



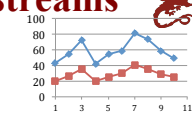
http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 190

Monitoring data streams

- Correlation coefficient

$$\rho = \frac{\sum_{t=1}^n (x_t - \bar{x}) \cdot (y_t - \bar{y})}{\sigma(x) \cdot \sigma(y)} \quad \sigma(x) = \sqrt{\sum_{t=1}^n (x_t - \bar{x})^2}$$
- Correlation coefficient and the (Euclidean) distance

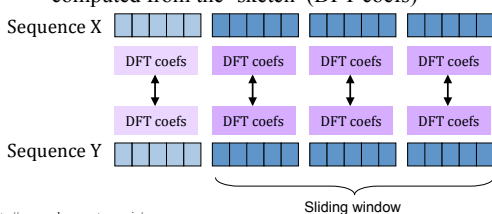
$$\rho = 1 - \frac{1}{2} \sum_{t=1}^n (\hat{x}_t - \hat{y}_t)^2 \quad \hat{x}_t = (x_t - \bar{x}) / \sigma(x)$$



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 191

Monitoring data streams

- Correlation monitoring [Zhu+, vldb02]
 - DFT coefficients for each basic window
 - Correlation coefficient of each sliding window computed from the 'sketch' (DFT coeffs)



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 192

Monitoring data streams

- Grid structure (to avoid checking all pairs)
 - DFT coefficients yields a vector
 - High correlation -> closeness in the vector space

Vector V_X of sequence X
 Vector V_Y of sequence Y

Correlation coefficients and the Euclidean distance

$$\rho = 1 - \frac{1}{2} \sum_{t=1}^n (\hat{x}_t - \hat{y}_t)^2$$

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 193

Monitoring data streams

- Lag correlation [Sakurai+, sigmod05]

CCF (Cross-Correlation Function)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 194

Monitoring data streams

- Lag correlation [Sakurai+, sigmod05]

correlated with lag $l=1300$

CCF (Cross-Correlation Function)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 195

Lag correlation

- Definition of 'score', absolute value of $R(l)$

$$score(l) = |R(l)| \quad R(l) = \frac{\sum_{t=l+1}^n (x_t - \bar{x})(y_{t-l} - \bar{y})}{\sqrt{\sum_{t=l+1}^n (x_t - \bar{x})^2} \sqrt{\sum_{t=1}^{n-l} (y_t - \bar{y})^2}}$$

- Lag correlation
 - Given a threshold γ , $score(l) > \gamma$
 - A local maximum
 - The earliest such maximum, if more maxima exist

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 196

Lag correlation

- Why not naïve?
 - Compute correlation coefficient for each lag
 - $l = \{0, 1, 2, 3, \dots, n/2\}$
- But
 - $O(n)$ space
 - $O(n^2)$ time
 - or $O(n \log n)$ time w/ FFT

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 197

Lag correlation

- BRAID
 - Geometric lag probing + smoothing
 - Use colored windows
 - Keep track of only a geometric progression of the lag values: $l = \{0, 1, 2, 4, 8, \dots, 2^h, \dots\}$

Multi-scale windows

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 198

Lag correlation

- BRAID
 - Geometric lag probing + smoothing
 - Keep track of only a geometric progression of the lag values: $l = \{0, 1, 2, 4, 8, \dots, 2^h, \dots\}$

Y $l=0$ $h=0$
X $l=n$ $h=0$

Correlation vs Lag

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 199

Lag correlation

- BRAID
 - Geometric lag probing + smoothing
 - Keep track of only a geometric progression of the lag values: $l = \{0, 1, 2, 4, 8, \dots, 2^h, \dots\}$

Y $l=1$ $h=0$
X $l=n$ $h=0$

Correlation vs Lag

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 200

Lag correlation

- BRAID
 - Geometric lag probing + smoothing
 - Keep track of only a geometric progression of the lag values: $l = \{0, 1, 2, 4, 8, \dots, 2^h, \dots\}$

Y $l=2$ $h=1$
X $l_n=n/2$ $h=1$

Correlation vs Lag

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 201

Lag correlation

- BRAID
 - Geometric lag probing + smoothing
 - Keep track of only a geometric progression of the lag values: $l = \{0, 1, 2, 4, 8, \dots, 2^h, \dots\}$

Y $l=4$ $h=2$
X $l_n=n/4$ $h=2$

Correlation vs Lag

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 202

Lag correlation

- BRAID
 - Geometric lag probing + smoothing
 - Keep track of only a geometric progression of the lag values: $l = \{0, 1, 2, 4, 8, \dots, 2^h, \dots\}$

Y $l=8$ $h=3$
X $l_n=n/8$ $h=3$

Correlation vs Lag

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 203

Lag correlation

- BRAID
 - Geometric lag probing + smoothing
 - Keep track of only a geometric progression of the lag values: $l = \{0, 1, 2, 4, 8, \dots, 2^h, \dots\}$
 - Use a cubic spline to interpolate

Y $l=n$ $h=0$

Correlation vs Lag

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 204

Lag correlation

- Why not naïve?
 - Compute correlation coefficient for each lag $l = \{0, 1, 2, 3, \dots, n/2\}$
- But
 - $O(n)$ space
 - $O(n^2)$ time
 - or $O(n \log n)$ time w/

BRAID

- $O(\log n)$ space
- $O(l)$ time

Multi-scale windows

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 205

BRAID in the real world

- Bridge structural health monitoring
 - Structural monitoring using vibration/shock sensors
 - Keep track of lag correlations for sensor data streams

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 206

BRAID in the real world

- Bridge structural health monitoring
 - Goal: real-time anomaly detection for disaster prevention
 - Several thousands readings (per sec) from several hundreds sensor nodes

Structural health monitoring

- Uses BRAID
- Metropolitan Expressway (Tokyo, Japan)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 207

BRAID in the real world

- Bridge structural health monitoring with BRAID

Metropolitan Expressway (Tokyo, Japan)

Tokyo Gate Bridge (Tokyo, Japan)

Can Tho Bridge (Vietnam)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 208

Feature extraction from streams

- Find hidden variables from streams [Papadimitriou+, vldb2005]

water distribution network

May have hundreds of measurements, but it is **unlikely they are completely unrelated!**

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 209

Feature extraction from streams

	Phase 1	Phase 2	Phase 3
sensors near leak
sensors away from leak

water distribution network

chlorine concentrations

normal operation

major leak

May have hundreds of measurements, but it is **unlikely they are completely unrelated!**

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 210

Motivation

chlorine concentrations

actual measurements (n streams)

$k = 1$

k hidden variable(s)

We would like to discover a few "hidden (latent) variables" that summarize the key trends

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 211

Motivation

chlorine concentrations

actual measurements (n streams)

$k = 2$

k hidden variable(s)

We would like to discover a few "hidden (latent) variables" that summarize the key trends

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 212

Motivation

chlorine concentrations

actual measurements (n streams)

$k = 1$

k hidden variable(s)

We would like to discover a few "hidden (latent) variables" that summarize the key trends

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 213

How to capture correlations?

First sensor

Temperature T_1

time

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 214

How to capture correlations?

First sensor

Second sensor

Temperature T_2

Temperature T_1

time

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 215

How to capture correlations?

Correlations:

Let's take a closer look at the first three value-pairs...

Temperature T_2

Temperature T_1

20°C 30°C

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 216

How to capture correlations?

First three lie (almost) on a line in the space of value-pairs...

- $O(n)$ numbers for the slope, and
- *One* number for each value-pair (offset on line)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 217

How to capture correlations?

Other pairs also follow the same pattern: they lie (approximately) on this line

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 218

Incremental update

For each new point

- Project onto current line
- Estimate error

• New value

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 219

Incremental update

For each new point

- Project onto current line
- Estimate error
- Rotate line in the direction of the error and in proportion to its magnitude

→ $O(n)$ time

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 220

Incremental update

For each new point

- Project onto current line
- Estimate error
- Rotate line in the direction of the error and in proportion to its magnitude

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 221

Related work

- Wavelet over streams [Gilbert+, vldb01] [Guha+, vldb04]
- Fourier representations [Gilbert+, stoc02]
- KNN [Koudas+, 04] [Korn+, vldb02]
- Histograms [Guha+, stoc01]
- Clustering [Guha+, focs00] [Aggarwal+, vldb03]
- Sketches [Indyk+, vldb00] [Cormode+, J. Algorithms 05]
- ...
- ...

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 222

Related work

- Heavy hitters [Cormode+, vldb03]
- Data embedding [Indyk+, focs00]
- Burst detection [Zhu+, kdd03]
- Segmentation [Keogh+, icdm01]
- Multiple scale analysis [Papadimitriou+, sigmod06]
- Fractal [Korn+, sigmod06]
- Time warping [Sakurai+, icde07]...
- ...

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 223

Outline

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

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 224

Motivation

Given: co-evolving time-series
 – e.g., MoCap (leg/arm sensors)

“Chicken dance”

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 225

Motivation

Given: co-evolving time-series
 – e.g., MoCap (leg/arm sensors)

“Chicken dance”

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 226

Motivation

Challenges: co-evolving sequences

- Unknown # of patterns (e.g., beaks)
- Different durations

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 227

Motivation

Goal: find patterns that agree with human intuition

Input

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 228

Motivation

Goal: find patterns that agree with human intuition

Input

Output

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 229

Motivation

Goal: find patterns that agree with human intuition

Input

NO magic numbers!

Automatic!

Output

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 230

Why: Automatic mining

No magic numbers! ... because,

Manual (use magic)

- sensitive to the parameter tuning
- long tuning steps (hours, days, ...)

Automatic (no magic numbers)

- no expert tuning required

Big data mining:
-> we cannot afford human intervention!!

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 231

How: Automatic mining

Goal: fully-automatic modeling

- Given: **data X**
- Find: a compact description (**model M**) of X

Data (X) → **Ideal model (M)**

Q. How can we find the best model M?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 232

How: Automatic mining

Goal: fully-automatic modeling

- Given: **data X**
- Find: a compact description (**model M**) of X

Answer: MDL!

Data (X) → **Ideal model (M)**

Q. How can we find the best model M?

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 233

Solution: MDL (Minimum description length)

Solution: Minimize total encoding cost \$!

- Occam's razor (i.e., law of parsimony)
- **Fully automatic** parameter optimization
- No over-fitting

Ideal model

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos [Bishop: PR&ML] 234

Solution: MDL (Minimum description length)

Solution: Minimize total encoding cost \$!

$$\text{Cost}_T(X;M) = \min (\text{Cost}_M(M) + \text{Cost}_c(X|M))$$

Total cost Model cost Coding cost (error)

\$\$\$ \$\$ \$ (Ideal!) \$\$\$\$

$C_M=0$ $C_M=1$ $C_M=3$ $C_M=9$

$C_c=\$$ $C_c=\$$ $C_c=\$$ $C_c=0$

[Bishop: PR&ML]

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 235

[Matsubara+ SIGMOD'14]

AutoPlait: Automatic Mining of Co-evolving Time Sequences

Yasuko Matsubara (Kumamoto University)
Yasushi Sakurai (Kumamoto University),
Christos Faloutsos (CMU)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 236

Problem definition

Goal: find patterns that agree with human intuition

Input

Output

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 237

Problem definition

- Bundle : set of d co-evolving sequences

given

$$X = \{x_1, \dots, x_n\}$$

$d \times n$

Bundle X ($d=4$)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 238

Problem definition

- Segment: convert $X \rightarrow m$ segments, S

hidden

$$S = \{s_1, \dots, s_m\}$$

Segment ($m=8$)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 239

Problem definition

- Regime: segment groups: $\Theta = \{\theta_1, \theta_2, \dots, \theta_r, \Delta_{r \times r}\}$

hidden

Regimes ($r=4$)

θ_r : model of regime r

beaks θ_1

wings θ_2

θ_3

θ_4

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 240

Problem definition

- Segment-membership: assignment

hidden $F = \{f_1, \dots, f_m\}$

$F = \{ \begin{matrix} 1 & 2 & 4 & 1 & 3 & 2 & 4 & 1 & 3 \\ 0.5 & & & & & & & & \\ 0 & & & & & & & & \end{matrix} \}$

Segment-membership (m=8)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 241

Problem definition

- Given: bundle X

$X = \{x_1, \dots, x_n\}$

- Find: compact description C of X

$C = \{m, r, S, \Theta, F\}$

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 242

Problem definition

- Given: bundle X

$X = \{x_1, \dots, x_n\}$

- Find: compact description C of X

$C = \{m, r, S, \Theta, F\}$

m segments
r regimes
Segment-membership

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 243

Main ideas

Goal: compact description of X

$C = \{m, r, S, \Theta, F\}$

without user intervention!!

Challenges:

Q1. How to generate 'informative' regimes ?

Q2. How to decide # of regimes/segments ?

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 244

Main ideas

Goal: compact description of X

$C = \{m, r, S, \Theta, F\}$

without user intervention!!

Challenges:

Q1. How to generate 'informative' regimes ?

Idea (1): Multi-level chain model

Q2. How to decide # of regimes/segments ?

Idea (2): Model description cost

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 245

Idea (1): MLCM: multi-level chain model

Q1. How to generate 'informative' regimes ?

Sequences → Model (beaks, claps, wings) → Regimes

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 246

Idea (1): MLCM: multi-level chain model

Q1. How to generate 'informative' regimes ?

Sequences → Model → Regimes

Idea (1): Multi-level chain model

- HMM-based probabilistic model
- with "across-regime" transitions

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 247

Idea (1): MLCM: multi-level chain model

$\Theta = \{\theta_1, \theta_2, \dots, \theta_r, \Delta_{r \times r}\}$ ($\theta_i = \{\pi, A, B\}$)

r regimes (HMMs) across-regime transition prob. Single HMM parameters

Regimes $r=2$
Regime 1 ($k=3$)
Regime 2 ($k=2$)

Regime1 "beaks" Regime2 "wings"

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 248

Idea (2): model description cost

Q2. How to decide # of regimes/segments ?

Idea (2): Model description cost

- Minimize encoding cost
- find "optimal" # of segments/regimes

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 249

Idea (2): model description cost

Idea: Minimize encoding cost!

$\min (\text{Cost}_M(M) + \text{Cost}_C(X|M))$

Model cost Coding cost

1 2 3 4 5 6 7 8 9 10 (# of r, m)

Good compression ↔ Good description

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 250

Idea (2): model description cost

Total cost of bundle X, given C
 $C = \{m, r, S, \Theta, F\}$

$$\begin{aligned} \text{Cost}_T(\mathbf{X}; C) &= \text{Cost}_T(\mathbf{X}; m, r, S, \Theta, F) \\ &= \log^*(n) + \log^*(d) + \log^*(m) + \log^*(r) + m \log(r) \\ &\quad + \sum_{i=1}^{m-1} \log^* |s_i| + \text{Cost}_M(\Theta) + \text{Cost}_C(\mathbf{X}|\Theta) \end{aligned} \quad (6)$$

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 251

Idea (2): model description cost

Total cost of bundle X, given C
 $C = \{m, r, S, \Theta, F\}$

duration/dimensions # of segments/regimes segment-membership F

$$\begin{aligned} \text{Cost}_T(\mathbf{X}; C) &= \text{Cost}_T(\mathbf{X}; m, r, S, \Theta, F) \\ &= \log^*(n) + \log^*(d) + \log^*(m) + \log^*(r) + m \log(r) \\ &\quad + \sum_{i=1}^{m-1} \log^* |s_i| + \text{Cost}_M(\Theta) + \text{Cost}_C(\mathbf{X}|\Theta) \end{aligned} \quad (6)$$

segment lengths Model description cost of Θ Coding cost of X given Θ

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/ © 2016 Sakurai, Matsubara & Faloutsos 252

AutoPlait

Overview

Iteration 0
 $r=1, m=1$

Start!

Cost

Iteration

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 253

AutoPlait

Overview

Iteration 1
 $r=2, m=4$

Split

Cost

Iteration

$f_1=2, f_2=1, f_3=2, f_4=1$

X

θ_1

θ_2

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 254

AutoPlait

Overview

Iteration 1
 $r=2, m=4$

Iteration 2
 $r=3, m=6$

Split

Cost

Iteration

$f_1=2, f_2=3, f_3=1, f_4=2, f_5=3, f_6=1$

X

θ_1

θ_2

θ_3

θ_4

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 255

AutoPlait

Overview

Iteration 1
 $r=2, m=4$

Iteration 2
 $r=3, m=6$

Iteration 4
 $r=4, m=8$

Split

Cost

Iteration

$f_1=2, f_2=4, f_3=3, f_4=1, f_5=2, f_6=4, f_7=3, f_8=1$

X

θ_1

θ_2

θ_3

θ_4

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 256

AutoPlait

Algorithms

- CutPointSearch** Inner-most loop
 Find good cut-points/segments
- RegimeSplit** Inner loop
 Estimate good regime parameters Θ
- AutoPlait** Outer loop
 Search for the best number of regimes ($r=2,3,4\dots$)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 257

1. CutPointSearch

Inner-most loop

Given:

- bundle X
- parameters of two regimes $\Theta = \{\theta_1, \theta_2, \Delta\}$

Find: cut-points of segment sets S_1, S_2

$$\{S_1, S_2\} = \operatorname{argmax}_{S_1, S_2} P(X | S_1, S_2, \Theta)$$

X

θ_1

θ_2

$S_1 = \{s_2, s_4\}$

$S_2 = \{s_1, s_3\}$

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 258

1. CutPointSearch

DP algorithm to compute likelihood: $P(X|\Theta)$

state 1 $L_{11}(1)=\phi$ $L_{11}(3)=\phi$...

state 2 $L_{21}(2)=\phi$ $L_{21}(5)=\phi$...

state 3 $L_{31}(4)=\phi$ $L_{31}(6)=\phi$...

switch?? $\delta_{2,3}$ $\delta_{2,4}$ $\delta_{2,5}$ $\delta_{2,6}$

state 1 $L_{22}(3)=\{3\}$ $L_{22}(4)=\{3\}$ $L_{22}(5)=\{4\}$...

state 2 $L_{22}(4)=\{4\}$ $L_{22}(5)=\{3\}$ $L_{22}(6)=\{4\}$...

X $t=1$ $t=2$ $t=3$ $t=4$ $t=5$ $t=6$...

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 259

1. CutPointSearch

Theoretical analysis

Scalability

- It takes $O(ndk^2)$ time (only single scan)
- n: length of X
- d: dimension of X
- k: # of hidden states in regime

Accuracy

It guarantees the optimal cut points

- (Details in paper)

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 260

2. RegimeSplit

Given: bundle X

Inner loop

Find: two regimes

1. find cut-points of segment sets: S_1, S_2
2. estimate parameters of two regimes: $\Theta = \{\theta_1, \theta_2, \Delta\}$

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 261

2. RegimeSplit

Two-phase iterative approach

- Phase 1: (CutPointSearch)
- Split segments into two groups: S_1, S_2
- Phase 2: (BaumWelch)
- Update model parameters: $\Theta = \{\theta_1, \theta_2, \Delta\}$

$S_1 = \{s_2, s_4\}$ Phase 1

$S_2 = \{s_1, s_3\}$ Phase 1

$\{\theta_1, \theta_2, \Delta\}$ Phase 2

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 262

3. AutoPlait

Given: bundle X

Outer loop

Find: r regimes ($r=2, 3, 4, \dots$)

- i.e., find full parameter set $C = \{m, r, S, \Theta, F\}$

$r = \min_r Cost_r(X; m, r, S, \Theta, F)$

$r=2, m=4$
 $f_1=2, f_2=1, f_3=2, f_4=1$

$r=4, m=8$
 $f_1=2, f_2=4, f_3=3, f_4=1, f_5=2, f_6=4, f_7=3, f_8=1$

$r=3, m=6$
 $f_1=2, f_2=3, f_3=1, f_4=2, f_5=3, f_6=1$

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 263

3. AutoPlait

Split regimes $r=2,3,\dots$, as long as cost keeps decreasing

- Find appropriate # of regimes

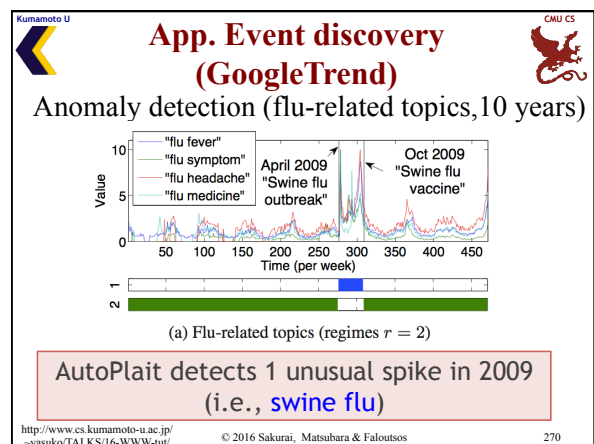
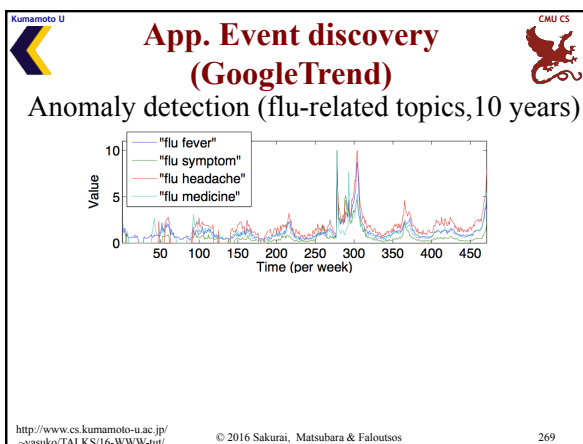
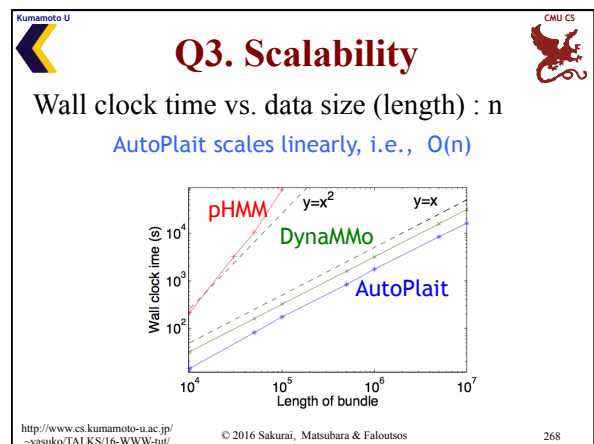
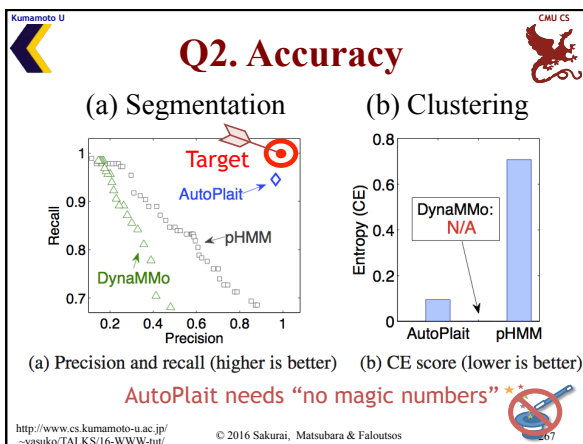
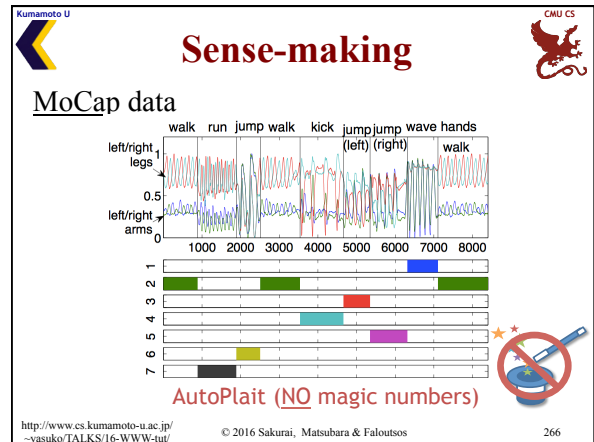
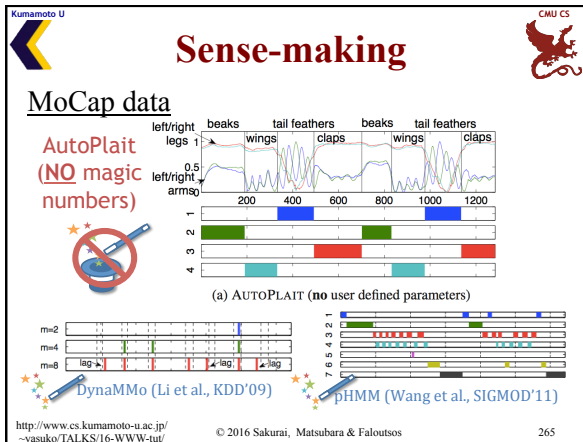
$r = \min_r Cost_r(X; m, r, S, \Theta, F)$

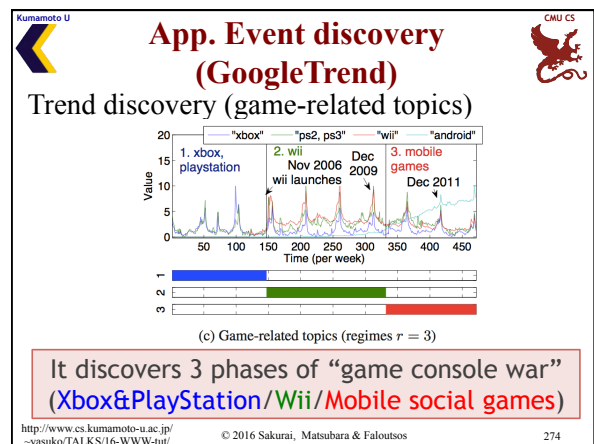
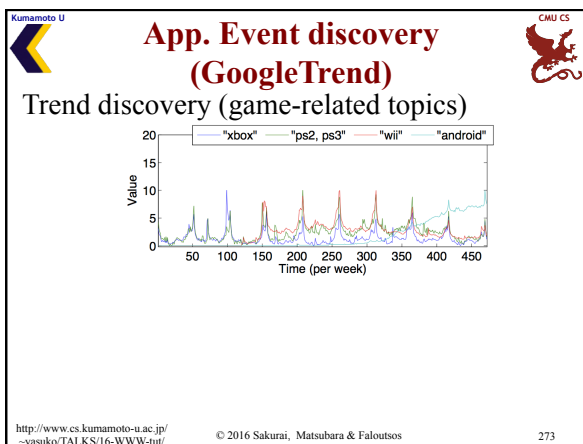
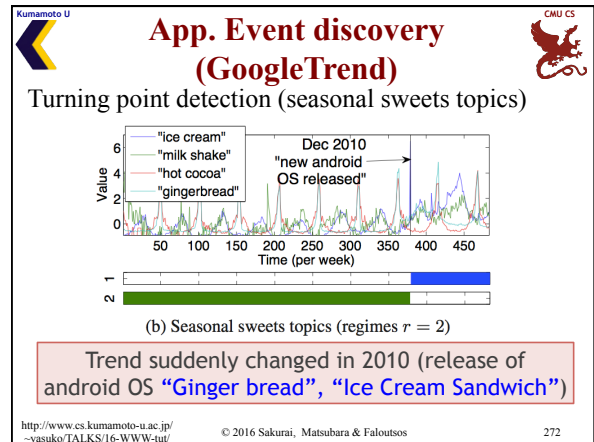
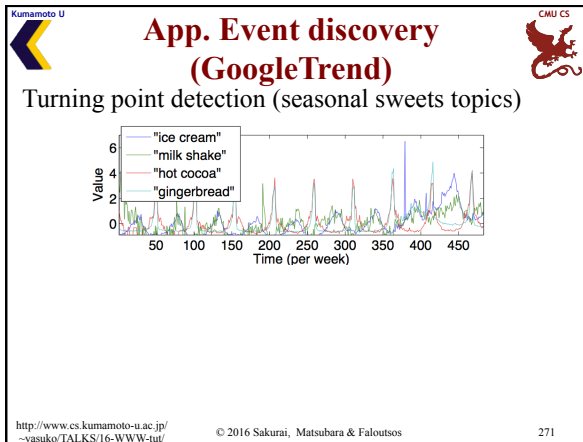
$r=2, m=4$
 $f_1=2, f_2=1, f_3=2, f_4=1$

$r=4, m=8$
 $f_1=2, f_2=4, f_3=3, f_4=1, f_5=2, f_6=4, f_7=3, f_8=1$

$r=3, m=6$
 $f_1=2, f_2=3, f_3=1, f_4=2, f_5=3, f_6=1$

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 264





Industrial contribution



- Automobile sensor data
 - location, velocity, longitudinal/lateral acceleration

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 275

Code at

- <http://www.cs.kumamoto-u.ac.jp/~yasuko/software.html>



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 276

Kumamoto U  

Part 1 – Conclusions

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



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 277

Kumamoto U  

Part 1 – Conclusions

- Motivation
- Similarity Search and Indexing
 - Euclidean/time-warping
 - extract features
 - index (SAM, R-tree)
- Feature extraction
 - SVD, ICA, DFT, DWT (multi-scale windows)



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 278

Kumamoto U  

Part 1 – Conclusions

- Linear forecasting
 - AR, RLS
- Streaming pattern discovery
 - RLS, “incremental” wavelet transform
 - Multi-scale windows
- Automatic mining
 - MDL



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 279

Kumamoto U  

References

- Yunyue Zhu, Dennis Shasha “StatStream: Statistical Monitoring of Thousands of Data Streams in Real Time” VLDB, August, 2002. pp. 358-369.
- Spiros Papadimitriou, Jimeng Sun, Christos Faloutsos. *Streaming Pattern Discovery in Multiple Time-Series*. VLDB 2005.
- Yasushi Sakurai, Spiros Papadimitriou, Christos Faloutsos. *BRAID: Stream Mining through Group Lag Correlations*. SIGMOD 2005.
- Anna C. Gilbert, Yannis Kotidis, S. Muthukrishnan, Martin Strauss. *Surfing Wavelets on Streams: One-Pass Summaries for Approximate Aggregate Queries*. VLDB 2001.
- Sudipto Guha, Nick Koudas, Amit Marathe, Divesh Srivastava. *Merging the Results of Approximate Match Operations*. VLDB 2004.
- Anna C. Gilbert, Sudipto Guha, Piotr Indyk, S. Muthukrishnan, Martin Strauss. *Near-optimal sparse fourier representations via sampling*. STOC 2002.



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 280

Kumamoto U  

References

- Nick Koudas, Beng Chin Ooi, Kian-Lee Tan, Rui Zhang. *Approximate NN queries on Streams with Guaranteed Error/performance Bounds*. VLDB 2004.
- Flip Korn, S. Muthukrishnan, Divesh Srivastava. *Reverse Nearest Neighbor Aggregates Over Data Streams*. VLDB 2002.
- Sudipto Guha, Nick Koudas, Kyuseok Shim. *Data-streams and histograms*. STOC 2001.
- Sudipto Guha, Nina Mishra, Rajeev Motwani, Liadan O’Callaghan. *Clustering Data Streams*. FOCS 2000.
- Charu C. Aggarwal, Jiawei Han, Jianyong Wang, Philip S. Yu. *A Framework for Clustering Evolving Data Streams*. VLDB 2003.



<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 281

Kumamoto U  

References



- Piotr Indyk, Nick Koudas, S. Muthukrishnan. *Identifying Representative Trends in Massive Time Series Data Sets Using Sketches*. VLDB 2000.
- Graham Cormode, S. Muthukrishnan. *An improved data stream summary: the count-min sketch and its applications*. J. Algorithms 55 (1), 2005.
- Graham Cormode, Flip Korn, S. Muthukrishnan, Divesh Srivastava. *Finding Hierarchical Heavy Hitters in Data Streams*. VLDB 2003.
- Piotr Indyk. *Stable Distributions, Pseudorandom Generators, Embeddings and Data Stream Computation*. FOCS 2000.
- Yunyue Zhu, Dennis Shasha. *Efficient elastic burst detection in data streams*. KDD 2003.

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 282

Kumamoto U  **References** 

- Eamonn J. Keogh, Selina Chu, David M. Hart, Michael J. Pazzani. *An Online Algorithm for Segmenting Time Series*. ICDM 2001.
- Spiros Papadimitriou, Philip S. Yu. *Optimal multi-scale patterns in time series streams*. SIGMOD 2006.
- Flip Korn, S. Muthukrishnan, Yihua Wu. *Modeling skew in data streams*. SIGMOD 2006.
- Yasushi Sakurai, Christos Faloutsos, Masashi Yamamuro. *Stream Monitoring under the Time Warping Distance*. ICDE 2007.

<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 283

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/16-WWW-tut/> © 2016 Sakurai, Matsubara & Faloutsos 284