**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 1

---

Kumamoto U  CMU CS

# Outline

- 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 2
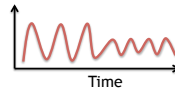
---

Kumamoto U  CMU CS

# 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/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 3

---

Kumamoto U  CMU CS

# Stream mining

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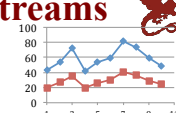    and feature extraction

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

---

Kumamoto U  CMU CS

# Monitoring data streams

- Correlation coefficient

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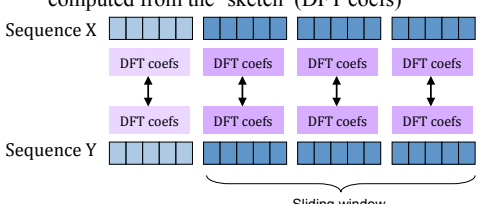
- Correlation coefficient and the (Euclidean) distance

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

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

---

Kumamoto U  CMU CS

# Monitoring data streams

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

Dennis Shasha

Sequence X

DFT coefs  DFT coefs  DFT coefs  DFT coefs

DFT coefs  DFT coefs  DFT coefs  DFT coefs

Sequence Y

Sliding window

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

## Slide 7

### 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

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

## Slide 8

### Monitoring data streams

- Lag correlation [Sakurai+, sigmod05]



CCF (Cross-Correlation Function)

## Slide 9

### Monitoring data streams

- Lag correlation [Sakurai+, sigmod05]



l=1300

correlated with lag $l$=1300

CCF (Cross-Correlation Function)

## Slide 10

### Lag correlation

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

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

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

## Slide 11

### Lag correlation

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

## Slide 12

### Lag correlation

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

Multi-scale windows



$2^3$   $2^2$ $2^1$ $2^0$

$2^2$

$2^3$

# Lag correlation

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

# Lag correlation

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

# Lag correlation

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

# Lag correlation

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

# Lag correlation

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

# Lag correlation

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

---

**Kumamoto U** | **CMU CS**

## Lag correlation

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

**BRAID**
- $O(\log n)$ space
- $O(1)$ time

*Time*   $t=n$

**Multi-scale windows**

---

**Kumamoto U** | **CMU CS**

## BRAID in the real world

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

---

**Kumamoto U** | **CMU CS**

## 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

  - Uses BRAID
  - Metropolitan Expressway (Tokyo, Japan)

Structural health monitoring      Vibration/shock sensor

---

**Kumamoto U** | **CMU CS**

## BRAID in the real world

- Bridge structural health monitoring with BRAID

Metropolitan Expressway (Tokyo, Japan)

Tokyo Gate Bridge (Tokyo, Japan)

Can Tho Bridge (Vietnam)

---

**Kumamoto U** | **CMU CS**

## Feature extraction from streams

major leak

normal operation

water distribution network

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

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

---

**Kumamoto U** | **CMU CS**

## Feature extraction from streams

hidden variables

Phase 1   Phase 2   Phase 3

sensors near leak

sensors away from leak

chlorine concentrations

water distribution network

normal operation      major leak

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

---

**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)

**How to capture correlations?**

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

**Incremental update**

For each new point
- Project onto current line
- Estimate error

**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

**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

**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]
- …
- …

## Related work

Tutorial@PODS'15

- 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]…
- …

Graham Cormode

Tutorial@SIGMOD'15

Divesh Srivastava

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

## Outline

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

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

## Motivation

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

"Chicken dance"

left/right legs 1
0.5
left/right arms 0

Time: 200 400 600 800 1000 1200

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

## Motivation

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

"Chicken dance"

Q. Any distinct patterns?

Q. If yes, how many?

Q. What kind?

left/right legs 1
0.5
left/right arms 0

Time: 600 800 1200

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

## Motivation

Challenges: co-evolving sequences
- Unknown # of patterns (e.g., beaks)
- Different durations

beaks   wings   tail feathers   claps   ...

left/right legs 1
0.5
left/right arms 0

Time: 200 400 600 800 1000 1200

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

## Motivation

Goal: find patterns that agree with human intuition

Input

left/right legs 1
0.5
left/right arms 0

200 400 600 800 1000 1200

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

## Slide 43

**Kumamoto U** · **CMU CS**

# Motivation

Goal: find patterns that agree with human intuition

**Input**

left/right legs
left/right arms



**Output**

Tail feathers · Tail feathers
Beaks · Beaks
Claps · Claps
Wings · Wings

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

## Slide 44

**Kumamoto U** · **CMU CS**

# Motivation

Goal: find patterns that agree with human intuition

**Input**

left/right legs
left/right arms

**NO magic numbers !**

**Automatic!**

**Output**

Tail feathers
Claps
Wings · Wings

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

## Slide 45

**Kumamoto U** · **CMU CS**

# 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/15-SIGMOD-tut/
© 2015 Sakurai, Matsubara & Faloutsos · 45

## Slide 46

**Kumamoto U** · **CMU CS**
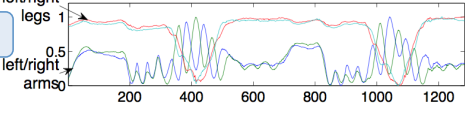
# How: Automatic mining

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/15-SIGMOD-tut/
© 2015 Sakurai, Matsubara & Faloutsos · 46

## Slide 47

**Kumamoto U** · **CMU CS**

# 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/15-SIGMOD-tut/
© 2015 Sakurai, Matsubara & Faloutsos · 47

## Slide 48

**Background** · **CMU CS**

# 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**

M=0 · M=1 · M=3 · M=9

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

## Slide 49

**Background**

### Solution: MDL (Minimum description length)

Solution: Minimize total encoding cost $ !

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

**Total cost**      Model cost      Coding cost (error)

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

$C_M=0$    $C_C=\$\$\$\$\$$

$C_M=1$    $C_C=\$\$\$$

$C_M=3$    $C_C=\$$

$C_M=9$    $C_C=0$

[Bishop: PR&ML]

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

## Slide 50

[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/15-SIGMOD-tut/      © 2015 Sakurai, Matsubara & Faloutsos      50

## Slide 51

### Problem definition

Goal: find patterns that agree with human intuition

**Input** — left/right legs, left/right arms

**Output** — Beaks, Tail feathers, Claps, Wings

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/      © 2015 Sakurai, Matsubara & Faloutsos      51

## Slide 52

### Problem definition

- Bundle : set of d co-evolving sequences

given

$$X = \{x_1,...,x_n\} \atop d \times n$$

Bundle X (d=4)

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

## Slide 53

### Problem definition

- Segment: convert X -> m segments, S

hidden

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

Segment (m=8)   s1 s2 s3 s4 s5 s6 s7 s8

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

## Slide 54

### Problem definition

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

hidden

Regimes (r=4)

$\theta_r$ : model of regime r

beaks — $\theta_1$, $\theta_2$
wings — $\theta_3$, $\theta_4$

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

## Problem definition

- Segment-membership: assignment

  hidden

$$F = \{f_1, \ldots, f_m\}$$

$F = \{ \; 2, \; 4, \; 1, \; 3, \; 2, \; 4, \; 1, \; 3 \; \}$

Segment-membership (m=8)

## Problem definition

- Given: bundle X

  left/right legs

  left/right arms

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

- Find: compact description C of X

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

## Problem definition

- Given: bundle X

  left/right legs

  left/right arms

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

- Find: compact description C of X

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

m segments

r regimes

Segment-membership

## 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 ?

## 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

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

Q1. How to generate 'informative' regimes ?

Model

beaks    claps

wings

Sequences          Regimes

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

Q1. How to generate 'informative' regimes?

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

Idea (1): Multi-level chain model
– HMM-based probabilistic model
– with "across-regime" transitions

---

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

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

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



Regime switch

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

Regime1 "beaks"   Regime2 "wings"

---

**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

---

**Idea (2): model description cost**

Idea: Minimize encoding cost!

$$\min \left( \boxed{Cost_M(M)} + \boxed{Cost_C(X|M)} \right)$$

Model cost   Coding cost

CostM, CostC, CostT

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

Good compression ⟳ Good description

---

**Idea (2): model description cost**

Total cost of bundle X, given C

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

$$Cost_T(\boldsymbol{X}; \mathcal{C}) = Cost_T(\boldsymbol{X}; m, r, \mathcal{S}, \boldsymbol{\Theta}, \mathcal{F})$$
$$= \log^*(n) + \log^*(d) + \log^*(m) + \log^*(r) + m\log(r)$$
$$+ \sum_{i=1}^{m-1} \log^* |s_i| + Cost_M(\boldsymbol{\Theta}) + Cost_C(\boldsymbol{X}|\boldsymbol{\Theta}) \qquad (6)$$
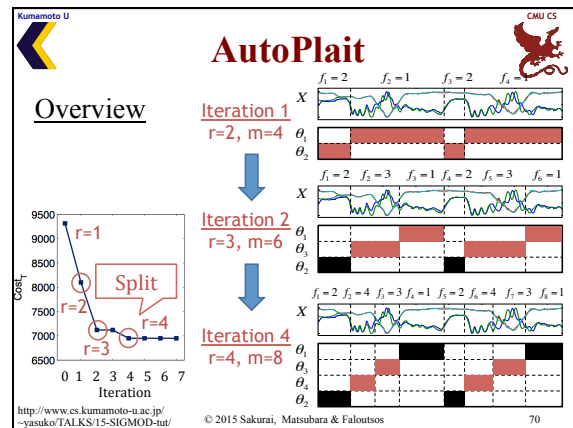
---

**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

$$Cost_T(\boldsymbol{X}; \mathcal{C}) = Cost_T(\boldsymbol{X}; m, r, \mathcal{S}, \boldsymbol{\Theta}, \mathcal{F})$$
$$= \boxed{\log^*(n) + \log^*(d)} + \boxed{\log^*(m) + \log^*(r)} + \boxed{m\log(r)}$$
$$+ \sum_{i=1}^{m-1} \boxed{\log^* |s_i|} + \boxed{Cost_M(\boldsymbol{\Theta})} + \boxed{Cost_C(\boldsymbol{X}|\boldsymbol{\Theta})} \qquad (6)$$

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

---

## AutoPlait

### Algorithms

1. CutPointSearch — Inner-most loop

   Find good cut-points/segments

2. RegimeSplit — Inner loop

   Estimate good regime parameters $\Theta$

3. AutoPlait — Outer loop

   Search for the best number of regimes (r=2,3,4…)

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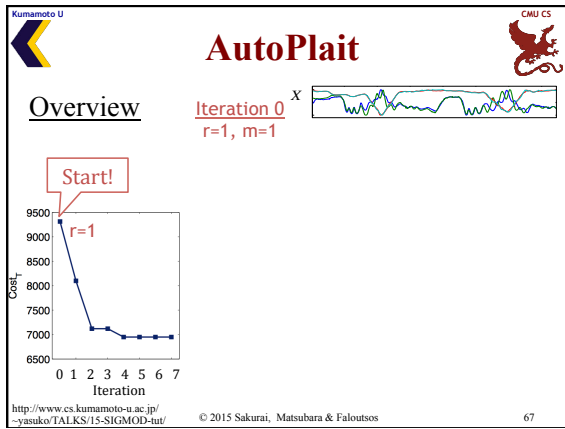~yasuko/TALKS/15-SIGMOD-tut/          © 2015 Sakurai, Matsubara & Faloutsos          71

## 1. CutPointSearch

Inner-most loop

Given:

- bundle $X$

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

Find: cut-points of segment sets $S_1$, $S_2$,

$$\{S_1, S_2\} = \arg\max P(X \mid S_1, S_2, \Theta)$$

$\{\theta_1, \theta_2, \Delta\}$          $S_1 = \{s_2, s_4\}$

$S_2 = \{s_1, s_3\}$

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

Slide 73 — **1. CutPointSearch**

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

$\theta_1$, $\theta_2$

states with $L_{1:1}(1) = \phi$, $L_{1:1}(3) = \phi$, $L_{1:3}(2) = \phi$, switch?? $\delta_{12}$ switch??

$L_{2:1}(3) = \{3\}$, $L_{2:1}(4) = \{3\}$, $L_{2:3}(5) = \{4\}$, $L_{2:2}(4) = \{4\}$, $L_{2:2}(5) = \{3\}$, $L_{2:2}(6) = \{4\}$

$X$ : $t=1$, $t=2$, $t=3$, $t=4$, $t=5$, $t=6$ …

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

Slide 74 — **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/15-SIGMOD-tut/ © 2015 Sakurai, Matsubara & Faloutsos 74

Slide 75 — **2. RegimeSplit**

Inner loop

Given:
- bundle $X$

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

$X$, $\theta_1$, $\theta_2$

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

Slide 76 — **2. RegimeSplit**

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

$X$, $\theta_1$, $\theta_2$

$S_1 = \{s_2, s_4\}$ Phase 1
$S_2 = \{s_1, s_3\}$

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

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

Slide 77 — **3. AutoPlait**

Outer loop

Given:
- bundle $X$

Find: r regimes (r=2, 3, 4, … )

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

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

Slide 78 — **3. AutoPlait**

Split regimes r=2,3,…, as long as cost keeps decreasing
- Find appropriate # of regimes

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

r=2, m=4
$f_1 = 2$ $f_2 = 1$ $f_3 = 2$ $f_4 = 1$
$X$, $\theta_1$, $\theta_2$

r=4, m=8
$f_1 = 2$ $f_2 = 4$ $f_3 = 3$ $f_4 = 1$ $f_5 = 2$ $f_6 = 4$ $f_7 = 3$ $f_8 = 1$
$X$, $\theta_1$, $\theta_3$, $\theta_4$, $\theta_2$

r=3, m=6
$f_1 = 2$ $f_2 = 3$ $f_3 = 1$ $f_4 = 2$ $f_5 = 3$ $f_6 = 1$
$X$, $\theta_1$, $\theta_3$, $\theta_2$

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

**Results**

- Mocap data
- WebClick data
- Google Trends



**Q1. Sense-making**

MoCap data



**Q1. Sense-making**

MoCap data



**Q2. Accuracy**



**Q3. Scalability**



**App. Event discovery (GoogleTrend)**

## Slide 85

**App. Event discovery (GoogleTrend)**

Anomaly detection (flu-related topics, 10 years)



(a) Flu-related topics (regimes $r = 2$)

AutoPlait detects 1 unusual spike in 2009 (i.e., swine flu)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/     © 2015 Sakurai, Matsubara & Faloutsos     85

## Slide 86

**App. Event discovery (GoogleTrend)**

Turning point detection (seasonal sweets topics)



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/     © 2015 Sakurai, Matsubara & Faloutsos     86

## Slide 87

**App. Event discovery (GoogleTrend)**

Turning point detection (seasonal sweets topics)



(b) Seasonal sweets topics (regimes $r = 2$)

Trend suddenly changed in 2010 (release of android OS "Ginger bread", "Ice Cream Sandwich")

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/     © 2015 Sakurai, Matsubara & Faloutsos     87

## Slide 88

**App. Event discovery (GoogleTrend)**

Trend discovery (game-related topics)



http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/     © 2015 Sakurai, Matsubara & Faloutsos     88

## Slide 89

**App. Event discovery (GoogleTrend)**

Trend discovery (game-related topics)



(c) Game-related topics (regimes $r = 3$)

It discovers 3 phases of "game console war" (Xbox&PlayStation/Wii/Mobile social games)

http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/     © 2015 Sakurai, Matsubara & Faloutsos     89

## Slide 90

**Industrial contribution**

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

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

---

**Kumamoto U**  **CMU CS**

# Code at

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

---

**Kumamoto U**  **CMU CS**

# Part 1 – Conclusions

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

---

**Kumamoto U**  **CMU CS**

# 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)

---

**Kumamoto U**  **CMU CS**

# Part 1 – Conclusions

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

---

**Kumamoto U**  **CMU CS**

# 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.

---

**Kumamoto U**  **CMU CS**

# 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.

---

**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/15-SIGMOD-tut/      © 2015 Sakurai, Matsubara & Faloutsos      97

**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/15-SIGMOD-tut/      © 2015 Sakurai, Matsubara & Faloutsos      98

**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      99