WindMine: Fast and Effective Mining of Web-click Sequences

Yasushi Sakurai (NTT) Lei Li (Carnegie Mellon Univ.) Yasuko Matsubara (Kyoto Univ.) Christos Faloutsos (Carnegie Mellon Univ.)

Introduction

Web-click sequence applications

Web masters and web-site owners

- Capacity planning
- Intrusion detection
- Advertisement design

<u>Goal</u>

- Find meaningful patterns for web-click data (e.g., the lunch-break trend, huge spike, anomalies)
- Find periodicity (daily and/or weekly, etc)
- Determine suitable window sizes automatically

Introduction

Examples

access count from a business news site





0

Problem definition

Web-click sequences of *m* URLs:

$$(X_1, \ldots, X_m)$$

Web-click sequence *X* of duration *n* :

 $X = (x_1, ..., x_t, ..., x_n)$

Local Component Analysis: Given *m* sequences of duration *n*, $(X_1, ..., X_m)$

- Find patterns, main components of the sequences
- Find the 'best window' size *w* for the analysis

Final challenge: scalable algorithm for the local component analysis

Background

Independent component analysis (ICA)

- PCA vs. ICA



Why not 'PCA'?

Example of component analysis



Why not 'PCA'?

Example of component analysis



Main idea (1)

Multi-scale local component analysis



Main idea (2)

Best window size selection

Q: How to estimate a 'good window size' automatically when we have multiple sequences?

Proposed criterion:

- CEM (Component Entropy Maximization)
- Estimate the optimal number of *w* for the sequence set
- Compute the entropy of the weight values of the mixing matrix A
- 'popular' (widely-used) components show high CEM scores

Main idea (2)

CEM criterion:

- CEM score of the *j*-th component for the window size *w*

$$C_{w,j} = -\frac{1}{\sqrt{w}} \sum_{i} p_{i,j} \log p_{i,j} \quad \begin{array}{c} \text{$k: \# $ of $ components$} \\ M: \# $ of $ subsequences$} \\ (i = 1, \dots, M; j = 1, \dots, k) \end{array}$$

- Probability for the *j*-th component (size of the *j*-th component's contribution to each subsequence)

$$p_{i,j} = \|a'_{i,j}\| / \sum_{i} \|a'_{i,j}\|$$

- Normalized weight values for each subsequences

$$a_{i,j}' = a_{i,j} / \sum_{j} a_{i,j}^2$$

- Mixing matrix $A_w = [a_{i,j}]$

WindMine-part

Efficient solution

Q: How do we efficiently extract the best local component from large sequence sets?

Hierarchical partitioning approach: WindMine-part

- Partition the original window matrix into sub-matrices
- Extract local components each from the sub-matrices
- Reuse the local components for the component analysis on the higher level

WindMine-part

X : original sequence



Experimental Results

Experiments with real and datasets

Ondemand TV, WebClick,

Automobile, Temperature, Sunspots

Evaluation

Accuracy for pattern discovery Accuracy for the best window size Computation time





SDM 2011



SDM 2011

WebClick

other websites





SDM 2011

Generalization of WindMine



SDM 2011

Y. Sakurai et al.

Choice of best window size

CEM score for various window sizes



(a) Ondemand TV









(c) Temperature

Y. Sakurai et al.



Wall clock time vs. # of subsequences

- Up to 70 times faster

Computation time



Ondemand TV

Wall clock time vs. duration

SDM 2011

Conclusions

Scalable pattern extraction and anomaly detection in large web-click sequences

- 1. Scalable, parallelizable method for breaking sequences into a few, fundamental ingredients
- 2. Linearly over the sequence duration, and near-linearly on the number of sequence