

The 18TH European Conference on Machine Learning and the 11TH European Conference on Principles and Practice of Knowledge Discovery in Databases

STATE-OF-THE-ART IN DATA STREAM MINING

TUTORIAL NOTES

presented by Mohamed Gaber and Joao Gama

> September 17, 2007 Warsaw, Poland

Prepared and presented by:

Mohamed Gaber Tasmanian ICT Centre, CSIRO ICT Centre, Australia João Gama Laboratory of Artificial Intelligence and Decision Support, INESC-Porto, University of Porto, Portugal

Tutorial Summary

Data streams became ubiquitous as many sources produce data continuously and rapidly. Examples of streaming data include customer click streams, telephone records, web logs, multimedia data, and sets of retail chain transactions. Data streams have brought new challenges to the data mining research community. In consequence, new techniques are needed to process streaming data in reasonable time and space. The goal of this tutorial is to present and discuss the research problems, issues and challenges in learning from data streams. We will present the state-of-the-art techniques in change detection, clustering, classification, frequent patterns, and time series analysis from data streams. Applications of mining data streams in different domains are highlighted. Open issues and future directions will conclude this tutorial. The tutorial also points to data stream mining resources.

Specific goals and objectives

- Introducing the area of data stream mining
- Giving a detailed explanation of the major techniques in the area
- Emphasizing the research issues and challenges

Expected background of the audience

Basic knowledge of data mining concepts and techniques is required.

Outline

- 1. Introduction
- 2. Data Streams
- 3. Change Detection
- 4. Learning Descriptive Models from Data Streams
- 5. Learning Predictive Models from Data Streams
- 6. Frequent pattern mining
- 7. Time series analysis in data streams
- 8. Applications of mining data streams
- 9. Future Directions

Prepared and presented by:

Mohamed Gaber Tasmanian ICT Centre, CSIRO ICT Centre, Australia

Mohamed Medhat Gaber is a research scientist at Commonwealth Scientific and Industrial Research Organization (CSIRO), Australia. He has published more than 40 articles. Mohamed has served in the program committees of several international and local conferences and workshops in the area of data mining. He has also been serving as a reviewer for the special issues of international journals in the area of data stream mining. He was the co-chair of the International Workshop on Mining Evolving and Streaming Data held in conjunction with ICDM 2006. He is the co-chair of the International Workshop on Knowledge Discovery from Ubiquitous Data Streams to be held in conjunction with ECML/PKDD 2007 and the ACM Workshop on Knowledge Discovery from Sensor Data to be held in conjunction with ACM SIGKDD 2007.

João Gama

Laboratory of Artificial Intelligence and Decision Support, INESC-Porto, University of Porto.

Joao Gama is a researcher at LIACC, the Laboratory of Artificial Intelligence and Computer Science of the University of Porto, working at the Machine Learning group. His main research interest is in Learning from Data Streams. He has published several articles in change detection, learning decision trees from data streams, hierarchical Clustering from streams, etc. Editor of special issues on Data Streams in Intelligent Data Analysis, J. Universal Computer Science, and New Generation Computing Co-chair of a series of Workshops on Knowledge Discovery in Data Streams, ECML 2004, Pisa, Italy, ECML 2005, Porto, Portugal, ICML 2006, Pittsburg, US, ECML 2006 Berlin, Germany, SAC2007, Korea, and the ACM Workshop on Knowledge Discovery from Sensor Data to be held in conjunction with ACM SIGKDD 2007.

Joao and Mohamed are editing a book titled: *Learning from Data Streams-Processing Techniques in Sensor Networks* to be published by Springer.











Data Streams

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Stream

Predictive Models from Data Streams Decision Trees Neural Networks **Continuous flow** of data generated at **high-speed** in **Dynamic**, **Time-changing** environments.

The usual approaches for querying, clustering and prediction use **batch procedures** cannot cope with this streaming setting. We need to maintain **Decision models** in **real time**. Decision Models must be capable of:

incorporating new information at the speed data arrives;

- **forgetting** outdated information;
- detecting changes and adapting the decision models to the most recent information.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ⊙

Massive Data Sets

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Stream

Predictive Models from Data Streams Decision Trees Neural Networks

- Data analysis is complex, interactive, and exploratory over very large volumes of historic data, eventually stored in distributed environments.
- Traditional pattern discovery process requires online ad-hoc queries, not previously defined, that are successively refined.

Due to the exploratory nature of these queries, an exact answer may not be required. A user may prefer a fast approximate answer.

Approximate Answers

State-of-the-Art in Data Stream Mining (Part I)

Approximate answers:

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering

Predictive Models from Data Streams Decision Trees Neural Networks

- Find an answer that is within 10% of correct result
- More generally, a $(1 \pm \epsilon)$ factor approximation
- Randomization: allow a small probability of failure
 - Answer is correct, except with probability 1 in 10,000
 - More generally, success probability (1δ)

Actual answer is within 5 ± 1 with probability > 0.9.

• Approximation and Randomization: (ϵ, δ) -approximations

The constants ϵ and δ have great influence in the space used. Typically the space is $O(1/\epsilon^2 \log(1/\delta))$.

Tail Inequalities

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Stream

Predictive Models from Data Streams Decision Trees Neural Networks

Approximate answers:

Trade-off between accuracy of the answer and computational resource required to compute an answer.

Tail inequalities:

General bounds on the tail probability of random variables. The probability that a random variable deviates far from its expectation.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ⊙

Chebyshev Inequality

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outime

Wotivation

Data Streams Basic Methods

Change Detection Predictive

Clustering

Predictive Models from Data Streams Decision Trees Neural Networks if X is a random variable with standard deviation σ , the probability that the outcome of X is no less than $k\sigma$ away from its mean is no more than $1/k^2$:

$$P(|X-\mu| \le k\sigma) \le \frac{1}{k^2}$$

No more than 1/4 of the values are more than 2 standard deviations away from the mean, no more than 1/9 are more than 3 standard deviations away, no more than 1/25 are more than 5 standard deviations away, and so on.

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

<ロ> <四> <四> <三</td>

9 Q (?

Chernoff Bound

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Streams

Predictive Models from Data Streams Decision Trees Neural Networks Consider a biased coin. One side is more likely to come up than other, but we don't know which and would like to find it.

- Flip it many times and then choose the side that comes up the most.
- How many times do you have to flip it to be confident that you've chosen correctly?

Example: p=0.6; $\delta = 95\%$

$$n \geq rac{\ln(1/\sqrt{\delta})}{(p-1/2)^2}$$

Hoeffding Bound

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive

Learning Clustering

Predictive Models from Data Streams Decision Trees Neural Networks Characterize the deviation between the true probability of some event and its frequency over m independent trials.

 $P(|\overline{X} - \mu| \ge \epsilon) \le 2exp(-2m\epsilon^2/R^2),$ where R is the range of the random variables.

Example: After seeing 100 examples of a random variable X, $x_i \in [0, 1]$, the sample mean is $\overline{x} = 0.6$; The true mean is with confidence δ in $\overline{x} \pm \epsilon$, where

 $\epsilon = \frac{\sqrt{R^2 \ln(1/\delta)}}{2n}$

Sampling

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Streams

Predictive Models from Data Streams Decision Trees Neural Networks To otain an unbiased sampling of the data, we need to know the lenght of the stream. In Data Streams, we need to modify the approach!

Strategy

- Sample instances at periodic time intervals
- Useful to *slow down* data.
- Involves *loss* of information.

Known Problems

Not possible to detect:

- Changes
- Anomalies

The reservoir Sample Technique

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivatior

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Stream

Predictive Models from Data Streams Decision Trees Neural Networks Vitter, J.; Random Sampling with a Reservoir, ACM, 1985.

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 三臣 - のへで

<ロ> <四> <四> <三</td>

 $\mathfrak{I} \mathfrak{Q} \mathfrak{Q}$

- Creates uniform sample of fixed size k;
- Insert first k elements into sample
- Then insert *i*th element with prob. $p_i = k/i$
- Delete an instance at random.

Illustrative Problems

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive

Clustering Data Streams

Predictive Models from Data Streams Decision Trees Neural Networks

Illustrative Problems

Illustrative Problems:

- Count the number of distinct values in a stream;
- Count the number of 1's in a sliding window of a binary string;
- Count frequent items above a given support.

Illustrative Problems

Stream Mining (Part I)

João Gama

State-of-the-Art in Data

Outline

Motivatior

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Streams

Predictive Models from Data Streams Decision Trees Neural Networks

- Illustrative Problems:
 - Count the number of distinct values in a stream;
 - Count the number of 1's in a sliding window of a binary string;
 - Count frequent items above a given support.

Count the Number of Distinct Values in a Stream

Assume that the domain of the attribute is $\{0, 1, ..., M - 1\}$. The problem is trivial if we have space linear in M. Is there an approximate solution is space log(M)?

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

Exponential Histograms

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Basic Methods

Change Detection Predictive Learning

Clustering Data Stream

Predictive Models from Data Streams Decision Trees Neural Networks Maintaining Stream Statistics over Sliding Windows, M.Datar, A.Gionis, P.Indyk, R.Motwani; ACM-SIAM Symposium on Discrete Algorithms;2002

The basic idea:

- Use buckets of different sizes to hold the data
- Each bucket has a timestamp associated with it
- It is used to decide when the bucket is out of the window

Data Structures for Exponential Histograms:

- Buckets: counts and time stamp
- LAST: stores the size of the last bucket.
- TOTAL: keeps the total size of the buckets.

The estimate of the sum of data elements is proven to be bounded within a user-specified parameter.

Exponential Histograms

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Stream

Predictive Models from Data Streams Decision Trees Neural Networks Consider a simplified data stream environment where each element comes from the same data source and is either 0 or 1. When a new data element arrives:

- If the new data element is 0, ignore it
- Otherwise create a new bucket of size 1 with the current timestamp, and increment the counter TOTAL.
- Given a parameter, e, if there are |1/e|/2 + 2 buckets of the same size, merge the oldest two of these same-size buckets into a single bucket of double size.
- The larger timestamp of the two buckets is then used as the timestamp of the newly created bucket.
- If the last bucket gets merged, we update the size of the merged bucket to the counter LAST.

Exponential Histograms State-of-the-Art in Data Stream Mining Whenever we want to estimate the moving sum: (Part I) João Gama Check if the oldest bucket is within the sliding window. If not, we drop that bucket: subtract its size from the variable TOTAL and update the size of the current oldest bucket to the variable Basic Methods LAST. Repeat the procedure until all the buckets with Predictive timestamps outside of the sliding window are dropped. The estimate of 1's in the sliding window is TOTAL-LAST/2. Decision Trees Neural Networks ◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ Exponential Histograms: Analysis State-of-the-Art in Data Stream Mining (Part I) The size of the buckets grows exponentially: João Gama $2^0, 2^1, 2^2 \dots 2^h$ Need only O(logN) buckets. ■ It is shown that, for N 1's in the sliding window, we only Basic Methods need $O((log N)/\epsilon)$ buckets to maintain the moving sum and the error of estimating Learning The error in the oldest bucket only. The moving sum is proven to be bounded within a given relative error. ϵ . Neural Networks <ロ> <四> <四> <三</td> $\mathfrak{I} \mathfrak{Q} \mathfrak{Q}$

Exponential Histograms: Example State-of-the-Art in Data Stream Time 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 Mining (Part I) 1 Element 1 1 1 0 1 0 1 1 1 1 1 1 1 0 João Gama Time **Buckets** Total Last T1 $\mathbf{1}_1$ 1 1 T2 2 1 $1_1, 1_2$ Т3 3 $1_1, 1_2, 1_3$ 1 (merge) $2_2, 1_3$ 3 1 Window length=10 Basic Methods Τ4 3 2 $2_2, 1_3, 1_4$ Relative Error=0.5 . . . Merge if 3 buckets of the T11 9 4 Predictive $4_4, 2_8, 2_{10}, 1_{11}$ same size: |1/0.5|/2/2 T12 $4_4, 2_8, 2_{10}, 1_{11}, 1_{12}$ 10 4 T13 $4_4, 4_{10}, 2_{12}, 1_{13}$ 11 4 T14 $4_4, 4_{10}, 2_{12}, 1_{13}, 1_{14}$ 12 4 (Removing out-of-date) T15 4 $4_{10}, 2_{12}, 1_{13}, 1_{14}$ 8 Decision Trees Neural Networks ◆□▶ ◆□▶ ◆国▶ ◆国▶ -21 590 Current Research on Data Streams

Introduction

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive

Learning Clustering

Predictive Models from

Data Streams Decision Trees Neural Networks Data flows continuously over time *Dynamic Environments*. Some characteristic properties of the problem can change over time.

Machine Learning algorithms assume:

- Instances are generated at random according to some probability distribution D.
- Instances are independent and identically distributed
- \blacksquare It is required that ${\mathcal D}$ is stationary

Examples:

- e-commerce, user modelling
- Spam emails
- Fraud Detection, Intrusion detection

Introduction

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Streams

Predictive Models from Data Streams Decision Trees Neural Networks **Concept drift** means that the concept about which data is obtained may shift from time to time, each time after some minimum permanence.

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

JAC.

Any change in the distribution underlying the data

Context: a set of examples from the data stream where the underlying distribution is stationary

A Framework based on Statistical Quality Control

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection

Predictive Learning

Clustering

Predictive Models from Data Streams Decision Trees Neural Networks Suppose a sequence of examples in the form $\langle \vec{x_i}, y_i \rangle$ The actual decision model classifies each example in the sequence

In the 0-1 loss function, predictions are either True or False The predictions of the learning algorithm are sequences: $T, F, T, F, T, F, T, T, T, F, \dots$

The Error is a random variable from *Bernoulli* trials.

The Binomial distribution gives the general form of the probability of observing a F:

 $p_i = (F/i)$ and $s_i = \sqrt{p_i(1-p_i)/i}$ where *i* is the number of trials.

The P-chart Algorithm

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Stream

Predictive Models from Data Streams Decision Trees Neural Networks The algorithm maintains two registers: P_{min} and S_{min} such that $P_{min} + S_{min} = min(p_i + s_i)$

Minimum of the error rate taking into account the variance of the estimator.

At example *j*:

The error of the learning algorithm will be

- Out-control if $p_j + s_j > p_{min} + \alpha * s_{min}$
- **In-control** if $p_j + s_j < p_{min} + \beta * s_{min}$
- Warning Level: if

 $p_{min} + \alpha * s_{min} > p_j + s_j > p_{min} + \beta * s_{min}$

The constants α and β depend on the desired confidence level. Admissible values are $\beta = 2$ and $\alpha = 3$.

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

Main Characteristics in Change Detection

State-of-the-Art in Data Stream Mining (Part I) Data management João Gama Characterizes the information about training examples stored in memory. Detection methods Characterizes the techniques and mechanisms for drift Basic Methods detection Adaptation methods Predictive Learning Adaptation of the decision model to the current distribution Decision model management Decision Trees Neural Networks ◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ Decision model management State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Basic Methods

Change Detection Predictive Learning

Clustering Data Stream

Models from Data Streams Decision Trees Neural Networks Model management characterize the number of decision models needed to maintain in memory.

The key issue here is the assumption that data generated comes from multiple distributions,

- at least in the transition between contexts.
- Instead of maintaining a single decision model several authors propose the use of multiple decision models.

<ロ> <四> <四> <三</td>

JAC.

Dynamic Weighted Majority

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detectior

Predictive Learning

Clustering Data Stream

Predictive Models from Data Streams Decision Trees Neural Networks A seminal work, is the system presented by Kolter and Maloof (ICDM03, ICML05).

The Dynamic Weighted Majority algorithm (DWM) is an ensemble method for tracking concept drift.

- Maintains an ensemble of base learners,
- Predicts using a weighted-majority vote of these *experts*.

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

<ロ> <四> <四> <三</td>

JAC.

 Dynamically creates and deletes experts in response to changes in performance.

Granularity of Decision Models

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Stream

Predictive Models from Data Streams Decision Trees Neural Networks Occurrences of drift can have impact in part of the instance space.

- Global models: Require the reconstruction of all the decision model. (like naive Bayes, SVM, etc)
- Granular decision models: Require the reconstruction of parts of the decision model (like decision rules, decision trees)

Online Divisive-Agglomerative Clustering

State-of-the-Art in Data Stream Mining Key Concept – Diameter of a cluster: the maximum (Part I) distance between two variables. João Gama Incremental system to monitor clusters' diameters Performs hierarchical clustering of first-order differences Can detect changes in the clustering structure Basic Methods Two Operators: Predictive Splitting: expand the structure Agglomeration: contract the structure Clustering Data Streams Splitting and agglomerative criteria are supported by a confidence level given by the Hoeffding bounds. Decision Trees Neural Networks ◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ Main Algorithm [Rodrigues, Gama, 2006] State-of-the-Art in Data Stream Mining (Part I) João Gama ForEver Read Next Example Compute first order differences For all the clusters Basic Methods Update the sufficient statistics Time to Time Predictive Learning Verify Merge Clusters Clustering Verify Expand Cluster Data Streams Neural Networks ◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Properties of ODAC

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Streams

Predictive Models from Data Streams Decision Trees Neural Networks

 For stationary data the cluster's diameters monotonically decrease.

Constant update time/memory consumption with respect to the number of examples!

- Every time a **split** is reported
 - the time to process the next example decreases, and
 - the **space** used by the new leaves is **less than** that used by the parent.

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

A snapshot - 1 year data, 2500 variables

	Outline
State-of-the- Art in Data Stream Mining (Part I) João Gama	1 Motivation
Outline Motivation	2 Data Streams
Data Streams Basic Methods Change Detection	3 Change Detection
Predictive Learning Clustering Data Streams	4 Clustering Data Streams
Predictive Models from Data Streams Decision Trees Neural Networks	5 Predictive Models from Data Streams
	うどの 州 へがた へゆ くち く
	Desirable properties:
State-of-the- Art in Data Stream Mining (Part I) João Gama Outline Motivation Data Streams Basic Methods Change Detection Predictive Learning Clustering Data Streams Predictive Models from Data Streams Decision Trees Neural Networks	 Processing each example: Small constant time Fixed amount of main memory Single scan of the data Without (or reduced) revisit old records. Eventually using a sliding window of more recent examples Processing examples at the speed they arrive Classifiers at anytime Ideally, produce a model equivalent to the one that would be obtained by a batch data-mining algorithm Ability to detect and react to concept drift

Very Fast Decision Trees

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive

Learning

Data Streams

Predictive Models from

Decision Trees Neural Networks

Mining High-Speed Data Streams, P. Domingos, G. Hulten; KDD00

The base Idea:

A small sample can often be enough to choose the optimal splitting attribute $% \left({{{\mathbf{x}}_{i}}} \right)$

- Collect sufficient statistics from a small set of examples
- Estimate the merit of each attribute
- Use Hoeffding bound to guarantee that the best attribute is really the *best*.
 - Statistical evidence that it is better than the second best

Very Fast Decision Trees: Main Algorithm

State-of-the-Art in Data Stream Mining (Part I)

João Gama

Outline

Motivation

Data Streams Basic Methods

Change Detection Predictive Learning

Clustering Data Streams

Predictive Models from Data Streams

Decision Trees Neural Networks

- **Input:** δ desired probability level.
- **Output:** T A decision Tree
- Init: $\mathcal{T} \leftarrow \text{Empty Leaf (Root)}$
- While (TRUE)

- Read next Example
- Propagate Example through the Tree from the Root till a leaf
- Update Sufficient Statistics at leaf
- If leaf (#examples) > N_{min}
 - Evaluate the merit of each attribute
 - Let A_1 the best attribute and A_2 the second best
 - Let $\epsilon = \sqrt{R^2 \ln(1/\delta)/(2n)}$
 - If $G(A_1) G(A_2) > \epsilon$
 - Install a splitting test based on A_1
 - Expand the tree with two descendant leaves

・ロト・日本・日本・日本・日本

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

VFDT: Analysis State-of-the-Art in Data Stream Mining (Part I) João Gama Low variance models: Stable decisions with statistical support. Low overfiting: Examples are processed only once. Basic Methods **Convergence**: VFDT becomes asymptotically close to Predictive that of a batch learner. The expected disagreement is δ/p ; Clustering where p is the probability that an example fall into a leaf. Decision Trees Neural Networks ◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ Neural-Nets and Data Streams State-of-the-Art in Data Stream Mining Multilayer Neural Networks (Part I) A general Function approximation method; João Gama A 3 layer ANN can approximate any continuous function with arbitrary precision; Fast Train and Prediction: Each example is propagated once The Error is back-propagated once Learning No overfitting First: Prediction Second: Update the Model Smoothly adjust to gradual changes Neural Networks <ロ> <四> <四> <三</td> $\mathfrak{I} \mathfrak{Q} \mathfrak{Q}$

State-of-the-art in Data Stream Mining (Part II)

Mohamed Medhat Gaber Tasmanian CSIRO ICT Centre Mail: GPO Box 1538, Hobart, TAS 7001, Australia E-mail: Mohamed.Gaber@csiro.au

Outline

- Frequent Pattern Mining in Data Streams
- Time Series Analysis in Data Streams
- Data Stream Mining Systems
- Applications of Mining Data Streams
- Future Directions
- Open Issues
- Future Vision
- Resources

Outline

- Frequent Pattern Mining in Data Streams
- Time Series Analysis in Data Streams
- Data Stream Mining Systems
- Applications of Mining Data Streams
- Future Directions
- Open Issues
- Future Vision
- Resources

Introduction to Frequent Pattern Mining

- Frequent pattern mining refers to finding patterns that occur greater than a prespecified threshold value.
- Patterns refer to items, itemsets, or sequences.
- Threshold refers to the percentage of the pattern occurrences to the total number of transactions. It is termed as <u>Support</u>

Introduction to Frequent Pattern Mining (Cont'd)

- Finding frequent patterns is the first step for the discovery of association rules in the form of $A \rightarrow B$.
- Apriori algorithm represents a pioneering work for association rules discovery
 - R Agrawal and R Srikant, Fast Algorithms for Mining Association Rules. In Proc. of the 20th International Conference on Very Large Databases, Santiago, Chile, September 1994
- An important step towards improving the performance of association rules discovery was FP-Growth
 - J. Han, J. Pei, and Y. Yin. Mining Frequent Patterns without Candidate Generation. In Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data (SIGMOD'00), Dallas, TX, May 2000.

Introduction to Frequent Pattern Mining (Cont'd)

- Many measurements have been proposed for finding the strength of the rules.
- The very frequently used measure is confidence.
- Confidence refers to the probability that set B exists given that A already exists in a transaction.
 - □ Confidence $(A \rightarrow B)$ = Support (AB) / Support (A)

Outline

- Frequent Pattern Mining in Data Streams
- Time Series Analysis in Data Streams
- Data Stream Mining Systems
- Applications of Mining Data Streams
- Future Directions
- Open Issues
- Future Vision
- Resources

Introduction to Time Series Analysis

- Time Series Analysis refers to applying different data analysis techniques on measurements acquired over temporal basis.
- Data analysis techniques recently applied on time series include clustering, classification, indexing, and association rules.
- The focus of classical time series analysis was on forecasting and pattern identification

Introduction to Time Series Analysis (Cont'd)

- Similarity measures over time series data represent the main step in time series analysis.
- Euclidean and dynamic time warping represent the major similarity measures used in time series.
- Longer time series could be represent computationally hard for the analysis tasks.
- Different time series representations have been proposed to reduce the length of a time series.

Time Series Analysis in Data Streams

- When data elements (records) in a data stream are processed based on their temporal dimension, we consider the process as time series analysis.
- Time series analysis in data streams are different in two aspects:
 - Several data points are considered to be an entry.
 - The analysis is done in real-time as opposed to traditional time series analysis.

Symbolic ApproXimation (SAX)

- SAX is a fast symbolic approximation of time series.
 - J. Lin, E. Keogh, S. Lonardi, and B. Chiu, A Symbolic Representation of Time Series, with Implications for Streaming Algorithms, in proceedings of the 8th ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery, San Diego, CA. June 13, 2003.
- It allows a time series with a length n to be transformed to an approximated time series with an arbitrarily length w, where w <<n.
- SAX follows three main steps:
 - Piecewise Aggregate Approximation (PAA)
 - Symbolic Discretization
 - Distance measurement
- SAX is generic and could be applied to any time series analysis technique.

Piecewise Aggregate Approximation (PAA)

 A time series with size n is approximated using PAA to a time series with size w using the following equation.

Where \bar{c}_i is the ith element in the approximated time series

Symbolic Discretization

- Breakpoints are calculated that produce equal areas from one point to another under Gaussian distribution.
 - □ A lookup table could be used.
- According to the output of PAA
 - If a point is less than the smallest breakpoint, then it is denoted as "a".
 - Otherwise and if the point is greater than the smallest breakpoint and less than the next larger one, then it is denoted as "b".
 - etc.

Hot SAX SAX has been used to discover discords in time series. The technique is termed as Hot SAX. Keogh, E., Lin, J. and Fu, A., HOT SAX: Efficiently Finding the Most Unusual Time Series Subsequence. In the 5th IEEE International Conference on Data Mining, New Orleans, LA. Nov 27-30, 2005. Discords are the time series subsequences that are maximally different from the rest of the time series subsequences. It is 3 to 4 times faster than brute force technique. This makes it a candidate for data streaming applications

Hot SAX (Cont'd)

- The process starts with sliding widows of a fixed size over the whole time series to generate subsequence
- Each generated subsequence is approximated using SAX
- The approximated subsequence is then inserted in an array indexed according to its position in the original time series
- The number of occurrences of each SAX word is also inserted in the array.

Data Stream Mining Systems

- Diamond Eye
 - The aim of the project is to enable remote systems as well as scientists to extract patterns from spatial objects in real time image streams.
 - The system uses a high performance computational facility for processing the data mining request
 - The scientist uses a web interface that uses java applets to connect to the server that requests that images to perform the image mining process.

M. Burl, Ch. Fowlkes, J. Roden, A. Stechert, and S. Mukhtar, Diamond Eye: A distributed architecture for image data mining, in SPIE DMKD, Orlando, April 1999, pp. 197-206

Data Stream Mining Systems (Cont'd)

- MobiMine
 - It is a client/server PDA-based distributed data mining application for financial data streams.
 - The system prototype has been developed using a single data source and multiple mobile clients; however the system is designed to handle multiple data sources.
 - The server functionalities in the proposed system are data collection from different financial web sites and storage, selection of active stocks using common statistics methods, and applying online data mining techniques to the stock data.

Kargupta, H., Park, B., Pittie, S., Liu, L., Kushraj, D. and Sarkar, K, MobiMine: Monitoring the Stock Market from a PDA. ACM SIGKDD Explorations. January 2002. Volume 3, Issue 2. Pages 37--46. ACM Press

- VEDAS
 - It stands for Vehicle Data Stream Mining System
 - It is a ubiquitous data stream mining system that allows continuous monitoring and pattern extraction from data streams generated on-board a moving vehicle.
 - The mining component is located on a PDA placed onboard the vehicle.
 - VEDAS uses online incremental clustering for modelling of driving behaviour.
 - Hillol Kargupta, Ruchita Bhargava, Kun Liu, Michael Powers, Patrick Blair, Samuel Bushra, James Dull, Kakali Sarkar, Martin Klein, Mitesh Vasa, and David Handy,

VEDAS: A Mobile and Distributed Data Stream Mining System for Real-Time Vehicle Monitoring, Proceedings of SIAM International Conference on Data Mining 2004

Data Stream Mining Systems (Cont'd)

EVE

- □ It stands for EnVironment for On-Board Processing
- □ It is used for astronomical data stream mining.
- Data streams are generated from measurements of different on-board sensors.
- Only interesting patterns are sent to the ground stations for further analysis preserving the limited bandwidth.

S. Tanner, M. Alshayeb, E. Criswell, M. Iyer, A. McDowell, M. McEniry, K. Regner, EVE: On-Board Process Planning and Execution, Earth Science Technology Conference, Pasadena, CA, Jun. 11 - 14, 2002

Data Stream Mining Systems (Cont'd)

MAIDS

- It stands for Mining Alarming Incidents of Data Streams.
- The system can classify, cluster, count frequency and query over data streams.
- It is a generic system as opposed to the other data stream mining systems that are application-based.

Y. D. Cai, D. Clutter, G. Pape, J. Han, M. Welge, and L. Auvil, MAIDS: Mining Alarming Incidents from Data Streams, (system demonstration), Proc. 2004 ACM-SIGMOD Int. Conf. Management of Data (SIGMOD'04), Paris, France, June 2004

Outline

- Frequent Pattern Mining in Data Streams
- Time Series Analysis in Data Streams
- Data Stream Mining Systems
- Applications of Mining Data Streams
- Future Directions
- Open Issues
- Future Vision
- Resources

Applications of Mining Data Streams

- Analysis of biosensor measurements around a city for security reasons
- Analysis of simulation results and on-board sensors in scientific laboratories and spacecrafts has its potential in changing the mission plan or the experimental settings in real time
- Analysis of web logs and web clickstreams

Outline

- Frequent Pattern Mining in Data Streams
- Time Series Analysis in Data Streams
- Data Stream Mining Systems
- Applications of Mining Data Streams
- Future Directions
- Open Issues
- Future Vision
- Resources

Future Directions

- Developing analysis algorithms for sensor networks to serve a number of real-time critical applications. SenosrNet (<u>www.sensornet.gov</u>) is one example in this direction.
- Online medical, scientific and biological analysis using data generated from medical, biological instruments and various tools employed in scientific laboratories.

- Frequent Pattern Mining in Data Streams
- Time Series Analysis in Data Streams
- Data Stream Mining Systems
- Applications of Mining Data Streams
- Future Directions
- Open Issues
- Future Vision
- Resources

Open Issues

- Interactive mining environment to satisfy user requirements
- The integration between data stream management systems and the ubiquitous data stream mining approaches
- Matching techniques with real world applications
- Data stream pre-processing

Open Issues (Cont'd)

- Model overfitting
- Data stream mining technology
- Real-time accuracy evaluation
- Theoretical foundations of data stream computing

Outline

- Frequent Pattern Mining in Data Streams
- Time Series Analysis in Data Streams
- Data Stream Mining Systems
- Applications of Mining Data Streams
- Future Directions
- Open Issues
- Future Vision
- Resources

Future Vision

- Wireless Sensor Networks provide environmental information.
- Building data mining models from this information according to the current context would contribute to build smart environments.
- Context-aware computing, data stream querying/mining, and wireless sensor networks will bring together the potential of research in this direction
- Examples include: Smart marketplace, smart workplace, smart vehicle and smart house.

The Future of Machine Learning

State-of-the-Art in Data Stream Mining (Part I)

João Gama and Mohamed Gaber

Conclusions and Open Issues Learning from small datasets: emphasis in variance reduction. Whats about large datasets?

- Increasing data = Variance reduction. Stable statistics estimators
- Learning from large datasets may be more effective using algorithms that places greater emphasis on bias management
- Solutions to these problems require
 - New Sampling and Randomize Techniques,
 - New Approximate, Incremental Algorithms,
 - Management the cost of Model's update and the Gains in Performance.
 - Incorporation of Change Detection Algorithms inside the Learning Process.

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ○ □ ● ○ ○ ○ ○

- First International Workshop on Knowledge Discovery from Data Streams (IWKDDS) at ECML/PKDD 2004 on September 24th, 2004, in Pisa, Italy.
 - Organized by:
 - Joao Gama, University of Porto, Portugal
 - Jesus S. Aguilar-Ruiz, University of Seville, Spain
 - □ Web: http://www.lsi.us.es/~aguilar/ecml2004/
- Second International Workshop on Knowledge Discovery from Data Streams (IWKDDS) at ECML/PKDD 2005 on October 10th, 2005, in Porto, Portugal.
 - Organized by:
 - Jesus S. Aguilar-Ruiz, University of Seville, Spain
 - Joao Gama, University of Porto, Portugal
 - Web: http://www.niaad.liacc.up.pt/~jgama/IWKDDS/

Resources (Cont'd)

- Third International Workshop on Knowledge Discovery from Data Streams (IWKDDS) at ICML 2006 on June 29th, 2006, at Carnegie Mellon University (CMU) in Pittsburgh, PA, USA.
 - Organized by:
 - Joao Gama, University of Porto, Portugal
 - Jesús S. Aguilar-Ruiz, University of Pablo de Olavide, Spain
 - Josep Roure, Carnegie Mellon University, US
 - Web: http://www.cs.cmu.edu/~jroure/iwkdds/iwkdds_icml06.html
- ECML/PKDD 2006 Workshop on Knowledge Discovery from Data Streams
 - Organized by:
 - João Gama, University of Porto, Portugal
 - Jesus S. Aguilar-Ruiz, University of Seville / University of Pablo de Olavide, Spain
 - Ralf Klinkenberg, University of Dortmund, Germany
 - Web: http://www.machine-learning.eu/iwkdds-2006/

Master References

Books

- Data Streams: Algorithms and Applications (Foundations and Trends in Theoretical Computer Science,) by S. Muthukrishnan (Now Publishers)
- Data Streams: Models and Algorithms (Advances in Database Systems) by Charu C. Aggarwal (Ed) (Springer)
- Learning from Data Streams: Processing Techniques in Sensor Networks by Joao Gama and Mohamed Medhat Gaber (Eds) (Springer)
- Seminal Surveys
 - B. Babcock, S. Babu, M. Datar, R. Motwani, and J. Widom. Models and Issues in Data Stream Systems, in Proceedings of PODS, 2002.
 - Gaber, M, M., Zaslavsky, A., and Krishnaswamy, S., Mining Data Streams: A Review, in ACM SIGMOD Record, Vol. 34, No. 1, March 2005, ISSN: 0163-5808
 - S. Muthukrishnan, Data streams: Algorithms and Applications. Proceedings of the fourteenth annual ACM-SIAM symposium on discrete algorithms, 2003

Researchers

- Charu Aggarwal
- Jesús S. Aguilar-Ruiz
- Yun Chi
- Graham Cormode
- Pedro Domingos
- Wei Fan
- João Gama
- Venkatesh Ganti
- Minos N. Garofalakis
- Johannes Gehrke
- Sudipto Guha
- Jiawei Han
- Geoff Hulten

Researchers (Cont'd)

- Hillol Kargupta
- Eamonn Keogh
- Ralf Klinkenberg
- Nikos Koudas
- Jessica Lin
- Nina Mishra
- Rajeev Motwani
- Muthu Muthukrishnan
- Olfa Nasraoui
- Rajeev Rastogi
- Haixun Wang
- Qian Weining
- Philip S. Yu

Thanks for your attention!

State-of-the-Art in Data Stream Mining (Part I)

João Gama and Mohamed Gaber

Conclusions and Open Issues

More information:

Sensors	J. Gama, R. Pederson; <i>Predictive Learning from Sensory Data</i> , Learning from Data Streams – Processing Techniques in Sensor Networks, Springer Verlag, 2007.
Streams	Learning from Data Streams – Processing Techniques in Sensor Networks, Editores J. Gama and M. Gaber, Springer Verlag, 2007.
Streams	S. Muthukrishnan, <i>Data Streams: Algorithms and Applications</i> , Now Publishers, 2003.
VFDT	P. Domingos, G. Hulten; <i>Learning from Infinite Data in Finite Time</i> , Advances in Neural Information Processing Systems 14. Cambridge, MA: MIT Press, 2002
VFDT	J. Gama, R. Fernandes, R. Rocha, <i>Decision Trees for Mining Data Streams</i> Intelligent Data Analysis, Vol. 10, Number 1, IOS Press, 2006.
ODAC	P. P. Rodrigues, J. Gama and J. P. Pedroso. <i>ODAC: Hierarchical Clustering of Time Series Data Streams</i> . In <i>Proceedings of the Sixth SIAM International Conference on Data Mining</i> , 2006.

<□> <□> <□> <□> <=> <=> <=> <<

Thanks for your attention! State-of-the-Art in Data Stream Mining (Part I) João Gama Mohamed Medhat Gaber (Eds.) João Gama and Mohamed Gaber (Eds. Conclusions and Open Learning from Issues Learning from Data Streams Learning from Data Streams Data Streams Processing Techniques in Sensor Networks 🖄 Springer 5900