**moa**

**Mining Big Data in Real Time**

Albert Bifet

**YAHOO! LABS**

Dortmund, 24 October 2013

# Albert Bifet



**2004-2009**
Ph. D. Degree
UPC-Barcelona Tech
Advisors: Ricard Gavaldà
and José L. Balcázar.

**2009-2012**
Post-Doctoral Researcher
University of Waikato,
Hamilton, New Zealand

**20011-2013**
Researcher
Yahoo! Research Barcelona

# Albert Bifet



**2004-2009**
Ph. D. Degree
UPC-Barcelona Tech
Advisors: Ricard Gavaldà
and José L. Balcázar.

**2009-2012**
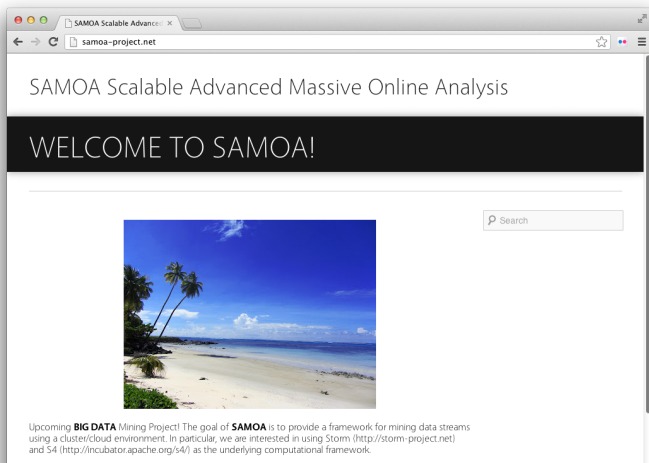Post-Doctoral Researcher
University of Waikato,
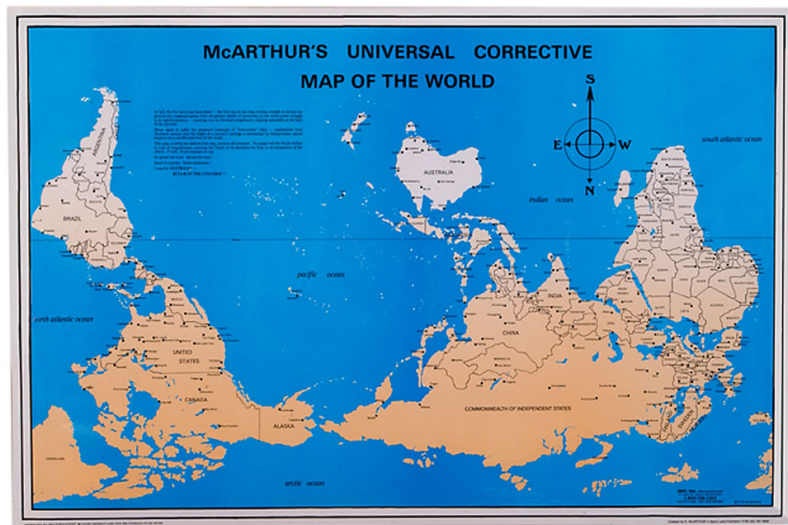Hamilton, New Zealand

**20011-2013**
Researcher
Yahoo! Research Barcelona

# MOA-SAMOA

# New Zealand

# Hamilton

# Outline

# Outline

Big Data & Real Time

# Big Data



McKinsey Global Institute (MGI) Report on Big Data, 2011.

**Big data** refers to datasets whose size is beyond the ability of typical database software tools to capture, store, manage, and analyze.

# BIG Data

- Volume
- Variety
- Velocity

$$3 \text{ Vs}$$

# BIG Data

- Volume
- Variety
- Velocity

- Variability
- Value
- Veracity

$$\boxed{\text{6 Vs}}$$

Distributed systems

# Methodology



*Paolo Boldi*
**Facebook** Four degrees of separation

Big Data does not need big machines,
it needs big **intelligence**

# Introduction: Data Streams

### Data Streams

- Sequence is potentially infinite
- High amount of data: sublinear space
- High speed of arrival: sublinear time per example
- Once an element from a data stream has been processed it is discarded or archived

### Example

Puzzle: Finding Missing Numbers

- Let $\pi$ be a permutation of $\{1,\ldots,n\}$.
- Let $\pi_{-1}$ be $\pi$ with one element missing.
- $\pi_{-1}[i]$ arrives in increasing order

Task: Determine the missing number

# Introduction: Data Streams

## Data Streams

- Sequence is potentially infinite
- High amount of data: sublinear space
- High speed of arrival: sublinear time per example
- Once an element from a data stream has been processed it is discarded or archived

## Example

Puzzle: Finding Missing Numbers

- Let $\pi$ be a permutation of $\{1, \ldots, n\}$.
- Let $\pi_{-1}$ be $\pi$ with one element missing.
- $\pi_{-1}[i]$ arrives in increasing order

Task: Determine the missing number

Use a $n$-bit vector to memorize all the numbers ($O(n)$ space)

# Introduction: Data Streams

## Data Streams

- Sequence is potentially infinite
- High amount of data: sublinear space
- High speed of arrival: sublinear time per example
- Once an element from a data stream has been processed it is discarded or archived

## Example

Puzzle: Finding Missing Numbers

- Let $\pi$ be a permutation of $\{1, \ldots, n\}$.
- Let $\pi_{-1}$ be $\pi$ with one element missing.
- $\pi_{-1}[i]$ arrives in increasing order

Task: Determine the missing number

Data Streams: $O(\log(n))$ space.

# Introduction: Data Streams

## Data Streams

- Sequence is potentially infinite
- High amount of data: sublinear space
- High speed of arrival: sublinear time per example
- Once an element from a data stream has been processed it is discarded or archived

## Example

Puzzle: Finding Missing Numbers

- Let $\pi$ be a permutation of $\{1, \ldots, n\}$.
- Let $\pi_{-1}$ be $\pi$ with one element missing.
- $\pi_{-1}[i]$ arrives in increasing order

Task: Determine the missing number

Data Streams: $O(\log(n))$ space. Store

$$\frac{n(n+1)}{2} - \sum_{j \leq i} \pi_{-1}[j].$$

# Data Streams

## Data Streams

- Sequence is potentially infinite
- High amount of data: sublinear space
- High speed of arrival: sublinear time per example
- Once an element from a data stream has been processed it is discarded or archived

## Tools:

- approximation
- randomization, sampling
- sketching

# Data Streams

## Data Streams

- Sequence is potentially infinite
- High amount of data: sublinear space
- High speed of arrival: sublinear time per example
- Once an element from a data stream has been processed it is discarded or archived

## Approximation algorithms

- Small error rate with high probability
- An algorithm $(\varepsilon, \delta)-$approximates $F$ if it outputs $\tilde{F}$ for which $\Pr[|\tilde{F} - F| > \varepsilon F] < \delta$.

# Data Streams Approximation Algorithms

$$1011000111\ \boxed{1010101}$$

## Sliding Window

We can maintain simple statistics over sliding windows, using $O(\frac{1}{\varepsilon} \log^2 N)$ space, where

- $N$ is the length of the sliding window
- $\varepsilon$ is the accuracy parameter

📄 M. Datar, A. Gionis, P. Indyk, and R. Motwani.
Maintaining stream statistics over sliding windows. 2002

# Data Streams Approximation Algorithms

10110001111 $\boxed{0101011}$

## Sliding Window

We can maintain simple statistics over sliding windows, using $O(\frac{1}{\varepsilon} \log^2 N)$ space, where

- $N$ is the length of the sliding window
- $\varepsilon$ is the accuracy parameter

📄 M. Datar, A. Gionis, P. Indyk, and R. Motwani.
Maintaining stream statistics over sliding windows. 2002

# Data Streams Approximation Algorithms

101100011110 | 1010111 |

## Sliding Window

We can maintain simple statistics over sliding windows, using $O(\frac{1}{\varepsilon} \log^2 N)$ space, where

- $N$ is the length of the sliding window
- $\varepsilon$ is the accuracy parameter

📄 M. Datar, A. Gionis, P. Indyk, and R. Motwani.
Maintaining stream statistics over sliding windows. 2002

# Data Streams Approximation Algorithms

1011000111101 $\boxed{0101110}$

## Sliding Window

We can maintain simple statistics over sliding windows, using $O(\frac{1}{\varepsilon} \log^2 N)$ space, where

- $N$ is the length of the sliding window
- $\varepsilon$ is the accuracy parameter

📄 M. Datar, A. Gionis, P. Indyk, and R. Motwani.
Maintaining stream statistics over sliding windows. 2002

# Data Streams Approximation Algorithms

$$10110001111010\ \boxed{1011101}$$

## Sliding Window

We can maintain simple statistics over sliding windows, using $O(\frac{1}{\varepsilon}\log^2 N)$ space, where

- $N$ is the length of the sliding window
- $\varepsilon$ is the accuracy parameter

📄 M. Datar, A. Gionis, P. Indyk, and R. Motwani.
Maintaining stream statistics over sliding windows. 2002

# Data Streams Approximation Algorithms

101100011110101 $\boxed{0111010}$

### Sliding Window

We can maintain simple statistics over sliding windows, using $O(\frac{1}{\varepsilon} \log^2 N)$ space, where

- $N$ is the length of the sliding window
- $\varepsilon$ is the accuracy parameter

📄 M. Datar, A. Gionis, P. Indyk, and R. Motwani.
Maintaining stream statistics over sliding windows. 2002

# What is MOA?

{M}assive {O}nline {A}nalysis is a framework for online learning from data streams.



- It is closely related to WEKA
- It includes a collection of offline and online as well as tools for evaluation:
  - classification
  - clustering
- Easy to extend
- Easy to design and run experiments

# History - timeline

## WEKA

- 1993 - WEKA : project starts (Ian Witten)
- 1996 - First public release of WEKA in C
- Early 1997 - decision was made to rewrite WEKA in Java
- Mid 1999 - WEKA 3 (100% Java) released

## MOA

- Nov 2007 - First public release of MOA: Richard Kirkby Thesis
- 2009 - MOA Concept Drift
- 2010 - MOA Clustering
- 2011 - MOA Graph Mining, Multi-label classification, Twitter Reader, Active Learning
- 2013 - MOA Outlier

# WEKA

- **W**aikato **E**nvironment for **K**nowledge **A**nalysis
- Collection of state-of-the-art machine learning algorithms and data processing tools implemented in Java
  - Released under the GPL
- Support for the whole process of experimental data mining
  - Preparation of input data
  - Statistical evaluation of learning schemes
  - Visualization of input data and the result of learning

- Used for education, research and applications
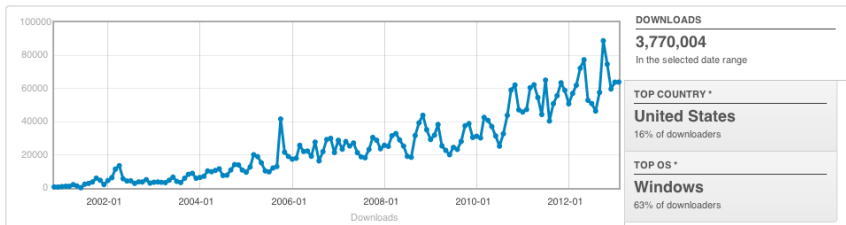- Complements "Data Mining" by Witten & Frank & Hall
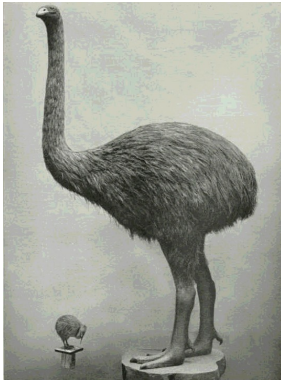
# WEKA: Impact Downloads

# MOA: the bird

The Moa (another native NZ bird) is not only flightless, like the Weka, but also extinct.

# MOA: the bird

The Moa (another native NZ bird) is not only flightless, like the Weka, but also extinct.

# MOA: the bird

The Moa (another native NZ bird) is not only flightless, like the Weka, but also extinct.

# Classification Experimental Setting

# Classification Experimental Setting

## Evaluation procedures for Data Streams

- Holdout
- Interleaved Test-Then-Train or Prequential

# Classification Experimental Setting

## Data Sources

- Random Tree Generator
- Random RBF Generator
- LED Generator
- Waveform Generator
- Hyperplane
- SEA Generator
- STAGGER Generator

# Classification Experimental Setting

## Classifiers

- Naive Bayes
- Decision stumps
- Hoeffding Tree
- Hoeffding Option Tree
- Bagging and Boosting
- ADWIN Bagging and Leveraging Bagging

## Prediction strategies

- Majority class
- Naive Bayes Leaves
- Adaptive Hybrid

# RAM-Hours

## RAM-Hour

Every GB of RAM deployed for 1 hour

## Cloud Computing Rental Cost Options

# Clustering Experimental Setting

# Clustering Experimental Setting

| Internal measures | External measures |
|---|---|
| Gamma | Rand statistic |
| C Index | Jaccard coefficient |
| Point-Biserial | Folkes and Mallow Index |
| Log Likelihood | Hubert Γ statistics |
| Dunn's Index | Minkowski score |
| Tau | Purity |
| Tau <u>A</u> | van Dongen criterion |
| Tau <u>C</u> | V-measure |
| Somer's Gamma | Completeness |
| Ratio of Repetition | Homogeneity |
| Modified Ratio of Repetition | Variation of information |
| Adjusted Ratio of Clustering | Mutual information |
| Fagan's Index | Class-based entropy |
| Deviation Index | Cluster-based entropy |
| <u>Z</u>-Score Index | Precision |
| <u>D</u> Index | Recall |
| Silhouette coefficient | F-measure |

Table : Internal and external clustering evaluation measures.

# Clustering Experimental Setting

## Clusterers

- StreamKM++
- CluStream
- ClusTree
- Den-Stream
- D-Stream
- CobWeb

# Web

http://www.moa.cms.waikato.ac.nz

# Easy Design of a MOA classifier



- `void resetLearningImpl ()`
- `void trainOnInstanceImpl (Instance inst)`
- `double[] getVotesForInstance (Instance i)`

# Easy Design of a MOA clusterer



- void resetLearningImpl ()
- void trainOnInstanceImpl (Instance inst)
- Clustering getClusteringResult()

# Extensions of MOA



- Multi-label Classification
- Active Learning
- Regression
- Closed Frequent Graph Mining
- Twitter Sentiment Analysis

## Challenges for bigger data streams

Sampling and distributed systems (Map-Reduce, Hadoop, S4)

# MOA: Impact Downloads

# Outline

# Outline

# Hoeffding Trees

## Hoeffding Tree : VFDT

📄 Pedro Domingos and Geoff Hulten.
Mining high-speed data streams. 2000

- With high probability, constructs an identical model that a traditional (greedy) method would learn
- With theoretical guarantees on the error rate

# Hoeffding Naive Bayes Tree

## Hoeffding Tree

Majority Class learner at leaves

## Hoeffding Naive Bayes Tree

📄 G. Holmes, R. Kirkby, and B. Pfahringer.
   Stress-testing Hoeffding trees, 2005.

- monitors accuracy of a Majority Class learner
- monitors accuracy of a Naive Bayes learner
- predicts using the most accurate method

# Bagging



Figure : Poisson(1) Distribution.

Bagging builds a set of *M* base models, with a bootstrap sample created by drawing random samples with replacement.

# Bagging



Figure : Poisson(1) Distribution.

Each base model's training set contains each of the original training example $K$ times where $P(K = k)$ follows a binomial distribution.

# Oza and Russell's *Online Bagging* for *M* models

1: Initialize base models $h_m$ for all $m \in \{1, 2, ..., M\}$
2: **for all** training examples **do**
3:    **for** $m = 1, 2, ..., M$ **do**
4:       Set $w = $ *Poisson*(1)
5:       Update $h_m$ with the current example with weight $w$

6: **anytime output:**
7: **return** hypothesis: $h_{fin}(x) = \arg\max_{y \in Y} \sum_{t=1}^{T} I(h_t(x) = y)$

# Data Mining Algorithms with Concept Drift



No Concept Drift

DM Algorithm

input → output

Counter$_5$

Counter$_4$

Counter$_3$

Counter$_2$

Counter$_1$

Concept Drift

DM Algorithm

Static Model

input → output

Change Detect.

# Data Mining Algorithms with Concept Drift



No Concept Drift

DM Algorithm
input → output
$Counter_5$
$Counter_4$
$Counter_3$
$Counter_2$
$Counter_1$

Concept Drift

DM Algorithm
input → output
$Estimator_5$
$Estimator_4$
$Estimator_3$
$Estimator_2$
$Estimator_1$

# Optimal Change Detector and Predictor

- High accuracy
- Low false positives and false negatives ratios
- Theoretical guarantees

- Fast detection of change
- Low computational cost: minimum space and time needed

- No parameters needed

# Algorithm ADaptive Sliding WINdow

### Example

$W =$ | 101010110111111 |

### ADWIN: ADAPTIVE WINDOWING ALGORITHM

1 Initialize Window $W$
2 **for** each $t > 0$
3     **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4         **repeat** Drop elements from the tail of $W$
5           **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6             for every split of $W$ into $W = W_0 \cdot W_1$
7         Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W =$ | 101010110111111 |

$W_0 =$ | 1 |   $W_1 =$ | 01010110111111 |

---

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1  Initialize Window $W$
2  **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4          **repeat** Drop elements from the tail of $W$
5              **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6                  for every split of $W$ into $W = W_0 \cdot W_1$
7          Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W=$ | 101010110111111 |

$W_0=$ | 10 | $W_1 =$ | 1010110111111 |

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1   Initialize Window $W$
2   **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4        **repeat** Drop elements from the tail of $W$
5         **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6          for every split of $W$ into $W = W_0 \cdot W_1$
7        Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

## Example

$W =$ | 101010110111111 |

$W_0 =$ | 101 |   $W_1 =$ | 010110111111 |

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1  Initialize Window $W$
2  **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4          **repeat** Drop elements from the tail of $W$
5              **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6                  for every split of $W$ into $W = W_0 \cdot W_1$
7          Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W =$ | 101010110111111 |

$W_0 =$ | 1010 |   $W_1 =$ | 10110111111 |

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1   Initialize Window $W$
2   **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4        **repeat** Drop elements from the tail of $W$
5          **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6           for every split of $W$ into $W = W_0 \cdot W_1$
7        Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W =$ | 101010110111111 |

$W_0 =$ | 10101 |   $W_1 =$ | 0110111111 |

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1  Initialize Window $W$
2  **for** each $t > 0$
3     **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4       **repeat** Drop elements from the tail of $W$
5         **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6           for every split of $W$ into $W = W_0 \cdot W_1$
7       Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W=$ | 101010110111111 |

$W_0=$ | 101010 |    $W_1 =$ | 110111111 |

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1   Initialize Window $W$
2   **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4         **repeat** Drop elements from the tail of $W$
5           **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6             for every split of $W$ into $W = W_0 \cdot W_1$
7         Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W = $ | 101010110111111 |

$W_0 = $ | 1010101 | $W_1 = $ | 10111111 |

---

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1  Initialize Window $W$
2  **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4         **repeat** Drop elements from the tail of $W$
5           **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6             for every split of $W$ into $W = W_0 \cdot W_1$
7         Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W=$ | 101010110111111 |

$W_0=$ | 10101011 |   $W_1 =$ | 0111111 |

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1  Initialize Window $W$
2  **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4          **repeat** Drop elements from the tail of $W$
5              **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6                  for every split of $W$ into $W = W_0 \cdot W_1$
7          Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W=$ `101010110111111`  $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ : CHANGE DET.!
$W_0=$ `101010110`  $W_1 =$ `111111`

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1  Initialize Window $W$
2  **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4          **repeat** Drop elements from the tail of $W$
5            **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6              for every split of $W$ into $W = W_0 \cdot W_1$
7          Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W$ = | 1̲01010110111111 |   Drop elements from the tail of $W$

$W_0$ = | 101010110 |  $W_1$ = | 111111 |

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1  Initialize Window $W$
2  **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4          **repeat** Drop elements from the tail of $W$
5              **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6                  for every split of $W$ into $W = W_0 \cdot W_1$
7              Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Example

$W =$ | 01010110111111 |  Drop elements from the tail of $W$

$W_0 =$ | 101010110 |  $W_1 =$ | 111111 |

ADWIN: ADAPTIVE WINDOWING ALGORITHM

1  Initialize Window $W$
2  **for** each $t > 0$
3      **do** $W \leftarrow W \cup \{x_t\}$ (i.e., add $x_t$ to the head of $W$)
4          **repeat** Drop elements from the tail of $W$
5              **until** $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \varepsilon_c$ holds
6                  for every split of $W$ into $W = W_0 \cdot W_1$
7              Output $\hat{\mu}_W$

# Algorithm ADaptive Sliding WINdow

### Theorem

*At every time step we have:*

1. (False positive rate bound). *If $\mu_t$ remains constant within $W$, the probability that ADWIN shrinks the window at this step is at most $\delta$.*

2. (False negative rate bound). *Suppose that for some partition of $W$ in two parts $W_0 W_1$ (where $W_1$ contains the most recent items) we have $|\mu_{W_0} - \mu_{W_1}| > 2\varepsilon_c$. Then with probability $1 - \delta$ ADWIN shrinks $W$ to $W_1$, or shorter.*

ADWIN tunes itself to the data stream at hand, with no need for the user to hardwire or precompute parameters.

# Algorithm ADaptive Sliding WINdow

ADWIN using a Data Stream Sliding Window Model,

- can provide the exact counts of 1's in $O(1)$ time per point.
- tries $O(\log W)$ cutpoints
- uses $O(\frac{1}{\varepsilon} \log W)$ memory words
- the processing time per example is $O(\log W)$ (amortized and worst-case).

### Sliding Window Model

| 1010101 | 101 | 11 | 1 | 1 |
|---------|-----|----|----|----|

| | | | | | |
|---|---|---|---|---|---|
| Content: | 4 | 2 | 2 | 1 | 1 |
| Capacity: | 7 | 3 | 2 | 1 | 1 |

# Decision Trees: Hoeffding Adaptive Tree

### Hoeffding Adaptive Tree:

- replace frequency statistics counters by estimators
  - don't need a window to store examples, due to the fact that we maintain the statistics data needed with estimators
- change the way of checking the substitution of alternate subtrees, using a change detector with theoretical guarantees

### Advantages over CVFDT:

1. Theoretical guarantees
2. No Parameters

# ADWIN Bagging (KDD'09)

### ADWIN
An adaptive sliding window whose size is recomputed online according to the rate of change observed.

### ADWIN has rigorous guarantees (theorems)
- On ratio of false positives and negatives
- On the relation of the size of the current window and change rates

### ADWIN Bagging
When a change is detected, the worst classifier is removed and a new classifier is added.

# Leveraging Bagging for Evolving Data Streams

Randomization as a powerful tool to increase accuracy and diversity

There are three ways of using randomization:

- Manipulating the input data
- Manipulating the classifier algorithms
- Manipulating the output targets

# Leveraging Bagging for Evolving Data Streams

## Leveraging Bagging

- Using $Poisson(\lambda)$

## Leveraging Bagging MC

- Using $Poisson(\lambda)$ and Random Output Codes

## Fast Leveraging Bagging ME

- if an instance is misclassified: weight = 1
- if not: weight = $e_T/(1 - e_T)$,

# Empirical evaluation

|  | Accuracy | RAM-Hours |
|---|---|---|
| Hoeffding Tree | 74.03% | 0.01 |
| Online Bagging | 77.15% | 2.98 |
| ADWIN Bagging | 79.24% | 1.48 |
| Leveraging Bagging | 85.54% | 20.17 |
| Leveraging Bagging MC | 85.37% | 22.04 |
| Leveraging Bagging ME | 80.77% | 0.87 |

## Leveraging Bagging

- Leveraging Bagging
  - Using *Poisson*($\lambda$)
- Leveraging Bagging MC
  - Using *Poisson*($\lambda$) and Random Output Codes
- Leveraging Bagging ME
  - Using weight 1 if misclassified, otherwise $e_T/(1 - e_T)$

# Outline

# Active Learning

📄 I. **Z**liobaitė, A. Bifet, B. Pfahringer, G. Holmes
Active learning with evolving streaming data

ACTIVE LEARNING FRAMEWORK

Input: labeling budget $B$ and strategy parameters

```
1  for each Xt - incoming instance,
2      do if ACTIVE LEARNING STRATEGY(Xt, B, ...) = true
3          then request the true label yt of instance Xt
4              train classifier L with (Xt, yt)
5              if Ln exists then train classifier Ln with (Xt, yt)
6              if change warning is signaled
7                  then start a new classifier Ln
8              if change is detected
9                  then replace classifier L with Ln
```

# Active Learning

I. **Z**liobaitė, A. Bifet, B. Pfahringer, G. Holmes
Active learning with evolving streaming data

|  | Controlling Budget | Instance space Coverage |
|---|---|---|
| Random | present | full |
| Fixed uncertainty | no | fragment |
| Variable uncertainty | handled | fragment |
| Randomized uncertainty | handled | full |

Table : Summary of strategies.

# Outline

# Multi-label classification



- Binary Classification: e.g. is this a beach? $\in \{$No, Yes$\}$
- Multi-class Classification: e.g. what is this?
  $\in \{$Beach, Forest, City, People$\}$
- Multi-label Classification: e.g. which of these?
  $\subseteq \{$Beach, Forest, City, People $\}$

# Methods for Multi-label Classification

Problem Transformation: Using off-the-shelf binary / multi-class classifiers for multi-label learning.

- **Binary Relevance method (BR)**
  - One binary classifier for each label:
    - simple; flexible; fast but does not explicitly model label dependencies

- **Label Powerset method (LP)**
  - One multi-class classifier; one class for each labelset

# Data Streams Multi-label Classification

- **Adaptive Ensembles of Classifier Chains (ECC)**
  - Hoeffding trees as base-classifiers reset classifiers based on current performance / concept drift

- **Multi-label Hoeffding Tree**
  - Label Powerset method (LP) at the leaves an ensemble strategy to deal with concept drift

- **MOA Multi-label Setting**
  - generating synthetic multi-label data streams
  - setting a benchmark on real-world and synthetic data

# Outline

# Pattern Mining

## Dataset Example

| Document | Patterns | Class |
|----------|----------|-------|
| d1 | abce | yes |
| d2 | cde | no |
| d3 | abce | yes |
| d4 | acde | no |
| d5 | abcde | no |
| d6 | bcd | yes |

Classification using patterns: mapping from patterns to vectors of attributes

# Itemset Mining

| d1 | abce |
|----|------|
| d2 | cde |
| d3 | abce |
| d4 | acde |
| d5 | abcde |
| d6 | bcd |

| Support | Frequent |
|---------|----------|
| 6 | c |
| 5 | e,ce |
| 4 | a,ac,ae,ace |
| 4 | b,bc |
| 4 | d,cd |
| 3 | ab,abc,abe |
|   | be,bce,abce |
| 3 | de,cde |

# Itemset Mining

| | |
|---|---|
| d1 | abce |
| d2 | cde |
| d3 | abce |
| d4 | acde |
| d5 | abcde |
| d6 | bcd |

| Support | Frequent | Gen | Closed |
|---------|----------|-----|--------|
| 6 | c | c | c |
| 5 | e,ce | e | ce |
| 4 | a,ac,ae,ace | a | ace |
| 4 | b,bc | b | bc |
| 4 | d,cd | d | cd |
| 3 | ab,abc,abe | ab | |
| | be,bce,abce | be | abce |
| 3 | de,cde | de | cde |

# Itemset Mining

| | |
|---|---|
| d1 | abce |
| d2 | cde |
| d3 | abce |
| d4 | acde |
| d5 | abcde |
| d6 | bcd |

| Support | Frequent | Gen | Closed | Max |
|---|---|---|---|---|
| 6 | c | c | c | |
| 5 | e,ce | e | ce | |
| 4 | a,ac,ae,ace | a | ace | |
| 4 | b,bc | b | bc | |
| 4 | d,cd | d | cd | |
| 3 | ab,abc,abe | ab | | |
| | be,bce,abce | be | abce | abce |
| 3 | de,cde | de | cde | cde |

# Graph Coresets KDD'11

## Coreset of a set *P* with respect to some problem

Small subset that approximates the original set *P*.

- Solving the problem for the coreset provides an approximate solution for the problem on *P*.

## $\delta$-tolerance Closed Graph

A graph *g* is $\delta$-*tolerance closed* if none of its proper frequent supergraphs has a weighted support $\leq (1 - \delta) \cdot support(g)$.

- Maximal graph: 1-tolerance closed graph
- Closed graph: 0-tolerance closed graph.

# Graph Coresets KDD'11

### Relative support of a closed graph

Support of a graph minus the relative support of its closed supergraphs.

- The sum of the closed supergraphs' relative supports of a graph is equal to its own support.

### $(s, \delta)$-coreset for the problem of computing closed graphs

Weighted multiset of frequent $\delta$-tolerance closed graphs with minimum support $s$ using their relative support as a weight.

# Graph Dataset

| Transaction Id | Graph | Weight |
|---|---|---|
| 1 | O<br>‖<br>C - C - S - N<br>‖<br>O | 1 |
| 2 | O<br>‖<br>C - C - S - N<br>C | 1 |
| 3 | O<br>‖<br>C - S - N<br>C | 1 |
| 4 | N<br>‖<br>C - C - S - N | 1 |

# Graph Coresets

| Graph | Relative Support | Support |
|-------|:----------------:|:-------:|
| C - C - S - N | 3 | 3 |
| O<br>C - S - N | 3 | 3 |
| N<br>‖<br>C - S | 3 | 3 |

Table : Example of a coreset with minimum support 50% and $\delta = 1$

# Graph Coresets



Figure : Number of graphs in $(40\%, \delta)$ for NCI.

# ChemDB dataset



**Memory ChemDB Dataset**

# ChemDB dataset



**Time ChemDB Dataset**

Legend: IncGraphMiner — IncGraphMiner-C - - MoSS — closeGraph

x-axis: **Instances** (10.000, 240.000, 470.000, 700.000, 930.000, 1.160.000, 1.390.000, 1.620.000, 1.850.000, 2.080.000, 2.310.000, 2.540.000, 2.770.000, 3.000.000, 3.230.000, 3.460.000, 3.690.000, 3.920.000)

y-axis: **Seconds** (0, 500, 1000, 1500, 2000, 2500, 3000, 3500, 4000, 4500, 5000)

# Outline

# New Techniques: Distributed Systems



Hadoop, S4 and Storm

Hadoop

# Hadoop



Hadoop architecture

# Apache Mahout



Mahout: open source framework

# **S4** *distributed stream computing platform*

Apache S4

# Apache S4



| EV | Quote |
|---|---|
| KEY | null |
| VAL | quote="I ..." |

A keyless event (EV) arrives at PE1 with quote:
"*I meant what I said and I said what I meant.*", Dr. Seuss

**QuoteSplitterPE** (PE1) counts unique words in Quote and emits events for each word.

| EV | WordEvent |
|---|---|
| KEY | word="said" |
| VAL | count=2 |

| EV | WordEvent |
|---|---|
| KEY | word="i" |
| VAL | count=4 |

**WordCountPE** (PE2–4) keeps total counts for each word across all quotes. Emits an event any time a count is updated.

| EV | UpdatedCountEv |
|---|---|
| KEY | sortID=2 |
| VAL | word=said count=9 |

**SortPE** (PE5–7) continuously sorts partial lists. Emits lists at periodic intervals

| EV | UpdatedCountEv |
|---|---|
| KEY | sortID=9 |
| VAL | word="i" count=35 |

| EV | PartialTopKEv |
|---|---|
| KEY | topk=1234 |
| VAL | words={w:cnt} |

**MergePE** (PE8) combines partial TopK lists and outputs final TopK list.

| PE ID | PE Name | Key Tuple |
|---|---|---|
| PE1 | QuoteSplitterPE | null |
| PE2 | WordCountPE | word="said" |
| PE4 | WordCountPE | word="i" |
| PE5 | SortPE | sortID=2 |
| PE7 | SortPE | sortID=9 |
| PE8 | MergePE | topK=1234 |

# Storm



Storm from Twitter

# Storm



Stream, Spout, Bolt, Topology

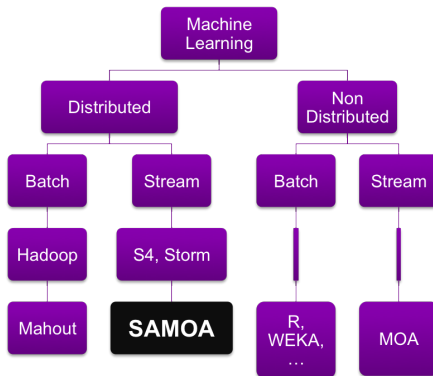# Storm

Tools



"Lambda Architecture"

Runaway complexity in Big Data
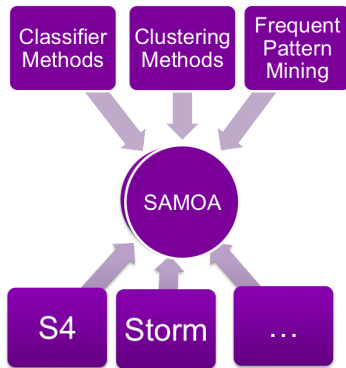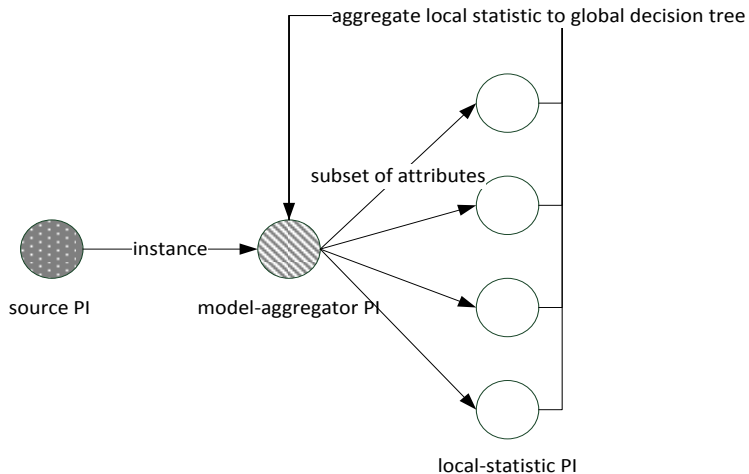Nathan Marz, 2012

# SAMOA



SAMOA: MOA + S4/Storm

# SAMOA

# SAMOA

# Vertical Hoeffding Tree

# Summary

{M}assive {O}nline {A}nalysis is a framework for online learning from data streams.



**http://moa.cs.waikato.ac.nz**

- It is closely related to WEKA
- It includes a collection of offline and online as well as tools for evaluation:
    - classification
    - clustering
    - frequent pattern mining
- MOA deals with evolving data streams
- MOA is easy to use and extend

# SAMOA



SAMOA: MOA + S4/Storm

# Thanks!