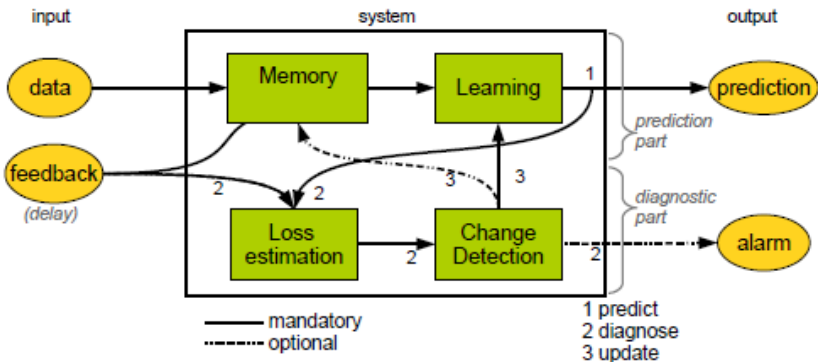


Adaptive Learning Algorithms



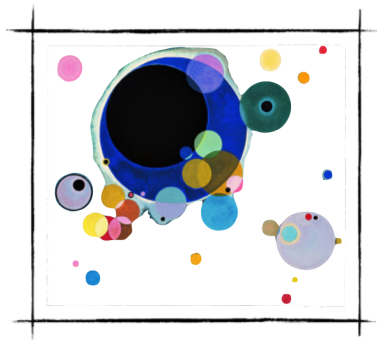
A generic schema for an online adaptive learning algorithm.

A survey on concept drift adaptation, Gama, Zliobaite, Bifet et al, ACM-CSUR 2014

Clustering

Clustering people or things into groups based on their attributes

- Customers segmentation
- Social network communities



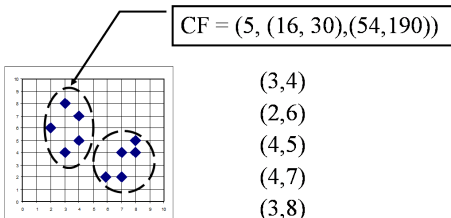


Cluster Feature Vector

Birch: Balanced Iterative Reducing and Clustering using Hierarchies, by Zhang, Ramakrishnan, Livny 1996

Cluster Feature Vector: $CF = (N, LS, SS)$

- N : Number of data points
- LS : $\sum_1^N \vec{x}_i$
- SS : $\sum_1^N (\vec{x}_i)^2$



Constant space irrespective to the number of examples!

Micro clusters

Given two micro-clusters CF_a and CF_b , a CF entry has sufficient information to calculate the norms:

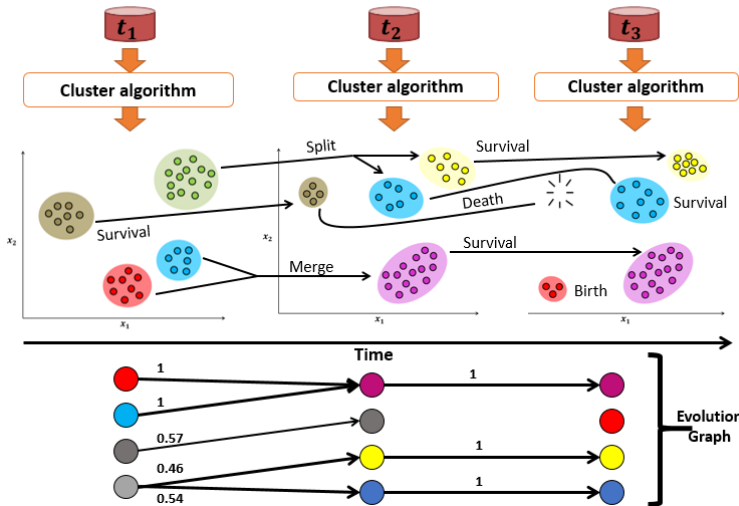
$$L_1 = \sum_{i=1}^n |LS_{a_i} - LS_{b_i}|$$

and

$$L_2 = \sqrt{\sum_{i=1}^n (LS_{a_i} - LS_{b_i})^2}$$

Cluster Evolution

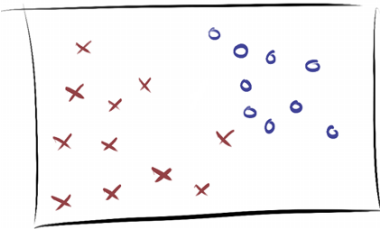
M.Oliveira, J.Gama, *A framework to monitor clusters evolution applied to economy and finance problems* Intell. Data Anal. 2012



Classification

Classifying people or things into groups by recognizing patterns

- Email spam filter
- Twitter sentiment analyzer



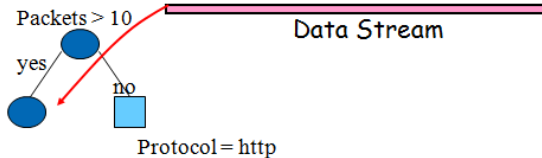
Very Fast Decision Trees

The base Idea

A small sample can often be enough to choose the optimal splitting attribute

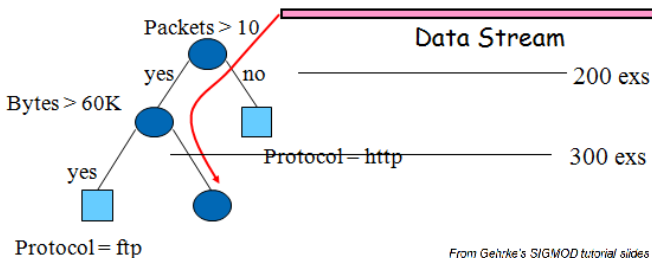
- Collect sufficient statistics from a small set of examples
- Estimate the merit of each attribute
- Suppose that after seeing n examples,
 $G(X_a) > G(X_b) > \dots > G(X_k)$
- Given a desired ϵ , the Hoeffding bound ensures that X_a is the correct choice, with probability $1 - \delta$, if $G(X_a) - G(X_b) > \epsilon$.
- If $G(X_a) - G(X_b) < \epsilon$, collect more examples

VFDT



$$H(\text{Bytes}) - H(\text{Packets}) > \epsilon$$

$$\epsilon = \sqrt{R^2 \log(1/\delta) / 2N}$$



From Gehrke's SIGMOD tutorial slides

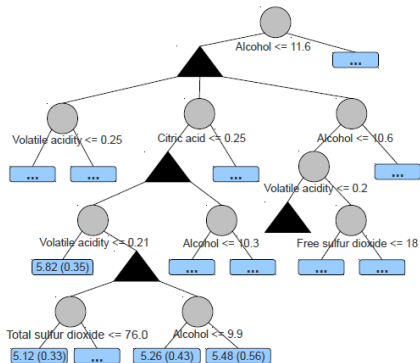
Option Trees

Speeding-Up Hoeffding-Based Regression Trees With Options, Ikonomovska, et al,

ICML 2011

Options nodes: OR nodes to encode alternatives

Use option nodes to solve ties

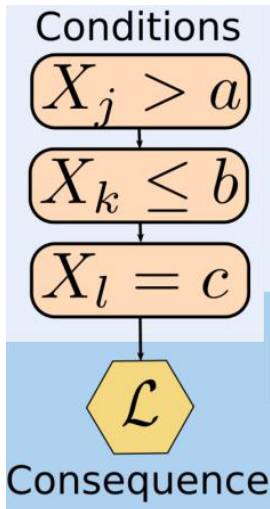


Decision and Regression Rules

Rules are one of the most expressive predictive models

- Rules are implications of the form Antecedent \rightarrow Consequent
- Antecedent: conjunction of conditions
- Consequent (\mathcal{L}) keeps sufficient statistics to:
 - make predictions
 - expand the rule
 - detect changes and anomalies

Rules are self-contained, modular, easier to interpret, no need to cover the universe

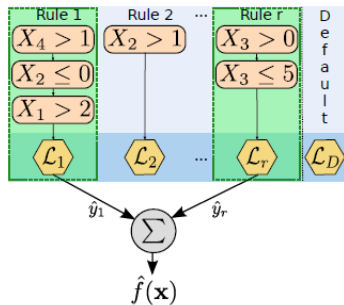


Adaptive Model Rules from Streams

Adaptive Model Rules from Data Streams, Almeida, Ferreira, Gama; ECML/PKDD

2013

- Ruleset: ensemble of rules
- Rule prediction: mean, linear model
- Ruleset prediction:
 - Ordered: only first rule covers instance
 - Unordered: weighted avg. of predictions of rules covering instance \mathbf{x}
 - Weights inversely proportional to error



E.g: $\mathbf{x} = [4, -1, 1, 2]$

$$\hat{f}(\mathbf{x}) = \sum_{R_l \in S(\mathbf{x}_i)} \theta_l \hat{y}_l,$$

AMRules Induction

- Rule creation: default rule expansion
- Rule expansion: split on attribute maximizing σ reduction
 - Hoeffding bound

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

$$\sigma_{1st}/\sigma_{2nd} < 1 - \epsilon$$
- Evict rule when P-H signals an alarm
- Detect and explain local anomalies

Algorithm 1: Training AMRules

Input: S : Stream of examples

begin
 $\mathcal{R} \leftarrow \{\}, D \leftarrow 0$
foreach $(x, y) \in S$ **do**
foreach $Rule\ r \in S(x)$ **do**
if $\neg IsAnomaly(x, r)$ **then**
if $PHTest(error_r, \lambda)$ **then**
 $\quad \perp$ Remove the rule from \mathcal{R}
else
 $\quad \perp$ Update sufficient statistics \mathcal{L}_r
 $\quad \perp$ $ExpandRule(r)$
if $S(x) = \emptyset$ **then**
 $\quad \perp$ Update \mathcal{L}_D
 $\quad \perp$ $ExpandRule(D)$
if D expanded **then**
 $\quad \perp$ $\mathcal{R} \leftarrow \mathcal{R} \cup D$
 $\quad \perp$ $D \leftarrow 0$
 \perp **return** $(\mathcal{R}, \mathcal{L}_D)$

Hoeffding Algorithms

- Classification:
Mining high-speed data streams, P. Domingos, G. Hulten, KDD, 2000
- Regression:
Learning model trees from evolving data streams; Ikonomovska, Gama, Dzeroski; Data Min. Knowl. Discov. 2011
- Decision Rules:
Learning Decision Rules from Data Streams, J. Gama, P. Kosina; IJCAI 2011
- Regression Rules
E. Almeida, C. Ferreira, J. Gama: Adaptive Model Rules from Data Streams. ECML/PKDD 2013
- Clustering:
Hierarchical Clustering of Time-Series Data Streams. Rodrigues, Gama, IEEE TKDE 20(5): 615-627 (2008)
- Multiple Models:
Ensembles of Restricted Hoeffding Trees. Bifet, Frank, Holmes, Pfahringer; ACM TIST; 2012
J. Duarte, J. Gama, Ensembles of Adaptive Model Rules from High-Speed Data Streams. BigMine 2014.
- ...

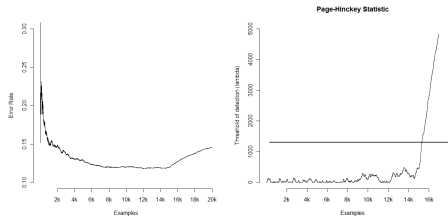
The Page-Hinckley Test

$$m_{t+1} = \sum_1^t (x_t - \bar{x}_t + \alpha)$$

- The minimum value of this variable is also computed with the following formula: $M_T = \min(m_t, t = 1 \dots T)$.
- The test monitors the difference between M_T and m_T :
 $PH_T = m_T - M_T$.
- When this difference is greater than a given threshold (λ) we alarm a change in the distribution.

Analysis

The threshold λ depends on the admissible false alarm rate. Increasing λ will entail fewer false alarms, but might miss some changes.



The left figure plots the on-line error rate of a learning algorithm. The right plot presents the evolution of the *PH* statistic.

Concept Drift

Gama, et. al, Learning with Drift Detection, SBIA 2004, Springer.

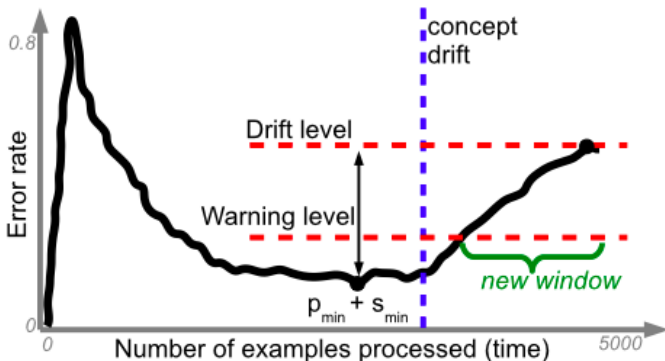
Learning from data streams is a continuous process. Monitor the quality of the learning process using quality control techniques.

The online error (e) of a learning algorithm is:

- **In-control** if $e < e_{min} + 2 \times s_{min}$
- **Out-control** if $e > e_{min} + 3 \times s_{min}$
- **Warning Level:** otherwise

Concept Drift

Statistical process control: monitor and control the learning process.



Algorithm ADaptive Sliding WINDow

A. Bifet, R Gavaldà: Learning from Time-Changing Data with Adaptive Windowing.
SDM 2007

Example

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

ADWIN: ADAPTIVE WINDOWING ALGORITHM

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

Algorithm ADaptive Sliding WINDOW

ADWIN using a Exponential Histogram Window Model,

- can provide the exact counts of 1's in $O(1)$ time per point.
- tries $O(\log W)$ cutpoints
- uses $O(\frac{1}{\epsilon} \log W)$ 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

Evaluation Methods

You cannot touch the same water twice.

Cross Validation and variants does not apply.

Two alternatives:

- Holdout if data is stationary.
- Sequential Sampling

What if the distribution is non-stationary ?

- The *prequential* approach.
 - For each example:
 - First: make a prediction
 - Second: update the model, whenever the target is available.
 - Evaluation over time-windows?



Prequential Evaluation

On evaluating stream learning algorithms Gama, Sebastião, Rodrigues, Machine Learning 2013

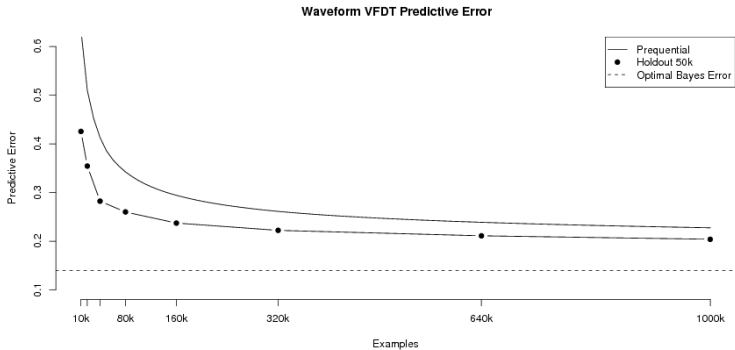
Definition: *The prequential error, computed at time i , is based on an accumulated sum of a loss function between the prediction and observed values:*

$$P_e(i) = \frac{1}{i} \sum_{k=1}^i L(y_k, \hat{y}_k) = \frac{1}{i} \sum_{k=1}^i e_k.$$

- 1 Provides a single number **at each time stamp**: a learning curve.
- 2 Pessimist estimator of accuracy.
- 3 Problematic to apply with algorithms with large testing time (k-NN).

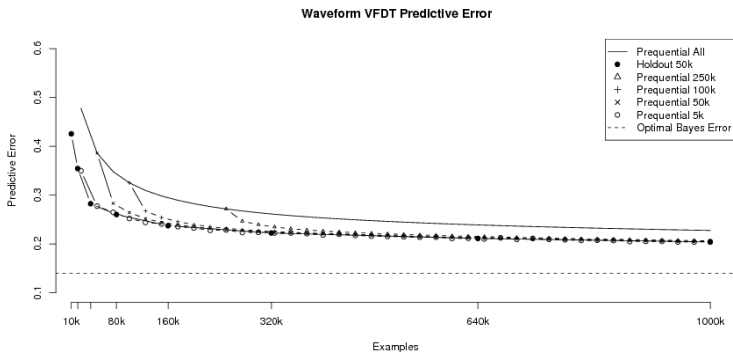
Prequential versus Holdout

Prequential is a pessimistic estimator.



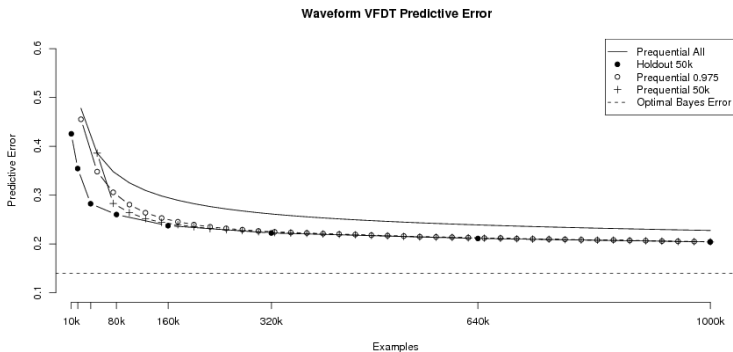
Prequential (sliding window) versus Holdout

Prequential over a sliding window converges to the holdout estimator.



Prequential (fading factor) versus Holdout

Prequential using fading factors converges to the holdout estimator.



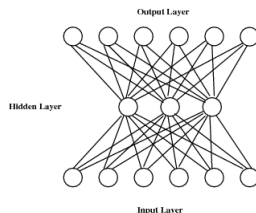
Definition

- Novelty Detection refers to the automatic identification of unforeseen phenomena embedded in a large amount of normal data.
- *Novelty* is a relative concept with regard to our current knowledge:
 - It must be defined in the context of a representation of our current knowledge.
- Specially useful when novel concepts represent abnormal or unexpected conditions
 - Expensive to obtain abnormal examples
 - Probably impossible to simulate all possible abnormal conditions

Autoassociator Networks

Concept-learning in the absence of counter-examples: an autoassociator-based approach Nathalie Japcowicz, 1999

- Three layer network
- The nr. of neurons in the output layer is equal to the input layer
- Train the network such that \vec{y} is equal to the \vec{x}
- The network is trained to reproduce the input at the output layer





Autoassociator Networks

To classify a test example \vec{x}

- Propagate \vec{x} through the network and let \vec{y} be the corresponding output;
- If $\sum_i^k (x_i - y_i)^2 < Threshold$ Then the example is considered from class *normal*;
- Otherwise, \vec{x} is a counter-example of the *normal* class.

Novelty detection

- Training set (Offline Phase)
 - $D_{tr} = (X_1, y_1), (X_2, y_2), \dots, (X_m, y_m)$
 - X_i : vector of input attributes for the i th example
 - y_i : target attribute
 - $y_i \in Y_{tr}$ where $Y_{tr} = c_1, c_2, \dots, c_L$
- When new data arrive (Online Phase)
 - Given a sequence of unlabelled examples X_{new}
 - Goal: Classify X_{new} in Y_{all} where $Y_{all} = c_1, c_2, \dots, c_L, \dots, c_K$ and $K > L$

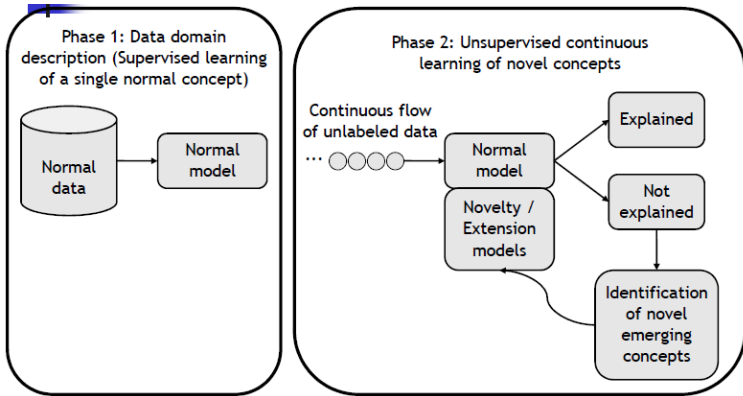
OLINDDA algorithm

OnLine Novelty and Drift Detection Algorithm

Spinosa, Carvalho, Gama: *OLINDDA: a cluster-based approach for detecting novelty and concept drift in data streams* SAC 2007

- Offline and Online phases
- Models: normal, extension and novelty
- Each model is represented by a set of clusters
- Not suitable for multi-class problem

OLLINDA



ECSMiner algorithm

Masud, Gao, Khan, Han, and Thuraisingham, *Classification and novel class detection in concept-drifting data streams under time constraints*, TKDE 2011

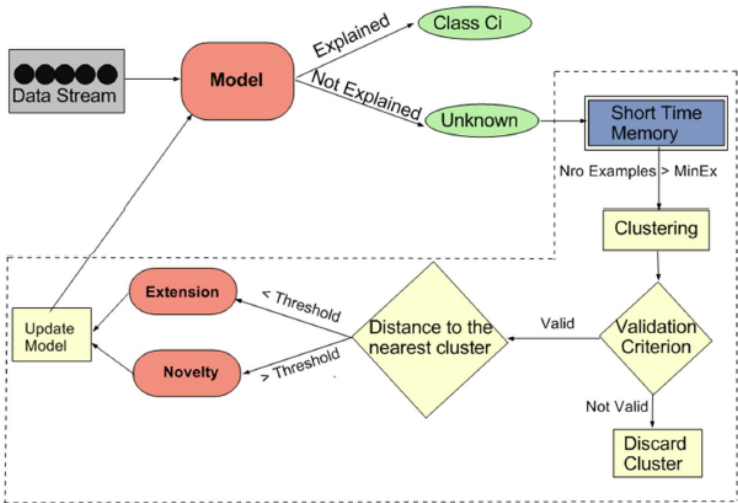
Supervised algorithm integrating novel concepts and concept drift

- Ensemble of classifiers
- Creates a new model when all examples in a chunk are labeled
 - Supposes that all examples in the stream will be labeled (after a delay of T_I time units)
 - An instance will be classified in until T_c time units of its arrival

MINAS - Offline phase

- Learns a decision model based on the known concept about the problem
KMeans or Clustream
- Run only once
- Each class is represent by a set of clusters (hyperspheres)

Minas



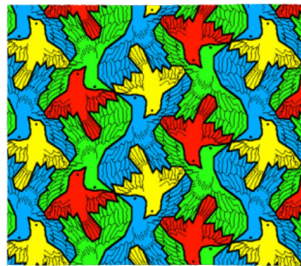
Novelty Detection Bibliography

- Masud, Gao, Khan, Han, and Thuruisingham, *Classification and novel class detection in concept-drifting data streams under time constraints*, TKDE 2011
- Spinoso, Carvalho, Gama: *OLINDDA: a cluster-based approach for detecting novelty and concept drift in data streams* SAC 2007
- MINAS: Multiclass Learning Algorithm for Novelty Detection in Data Streams, E. Faria, J. Gama, A. Carvalho, DAMI (to appear)
- P. Angelov and X. Zhou, *Evolving fuzzy-rule-based classifiers from data streams* Trans. Fuz Syst. 2008.
- D. Tax and R. Duin, *Growing a multi-class classifier with a reject option* Pattern Recognit. Lett., 2008.
- F. Denis, R. Gilleron, and F. Letouzey, *Learning from positive and unlabeled examples*, Theoretical Comput. Sci., 2005.
- D. Cardoso and F. França *A Bounded Neural Network for Open Set Recognition*, IJCNN 2015

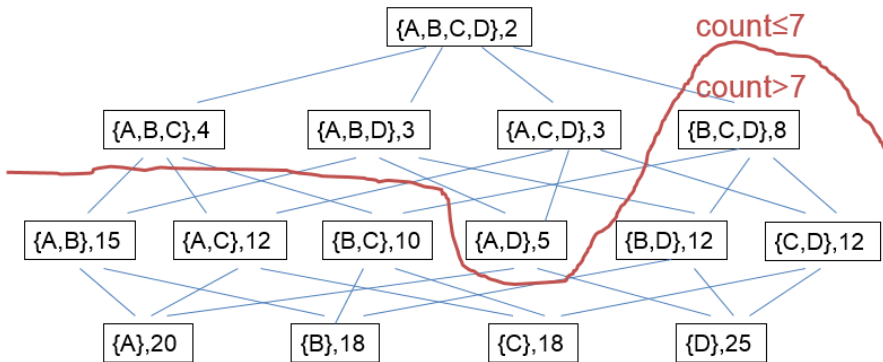
Frequent Pattern Mining

Given a collection of sets of items, find all the subsets that occur frequently

- Market basket mining
- Item recommendation



Key data structure: Lattice of patterns, with counts



Massive Online Analysis

Configure EvaluatePrequential -l trees.HoeffdingTree -s generators.WaveformGenerator Run

command	status	time elapsed	current activity	% complete
EvaluatePrequential -l trees.Ho...	running	10m11s	Evaluating learner...	21,22
EvaluatePrequential -l trees.Ho...	running	11m13s	Evaluating learner...	12,25

Pause Resume Cancel Delete

Preview (11m13s) Refresh Auto refresh: every second

84982E-7,8900000.0,84.6,76.90361458431755,8900000.0,-11239.0,3211.0,1606.0,1606.0,24.0,0.0,0.0,-Infinity
 84488E-7,9000000.0,83.8,75.6566322904254,9000000.0,-11330.0,3237.0,1619.0,1619.0,25.0,0.0,0.0,-Infinity
 67947E-7,9100000.0,86.0,78.92784895482129,9100000.0,-11505.0,3287.0,1644.0,1644.0,25.0,0.0,0.0,-Infinity
 17032E-7,9200000.0,86.1,79.14785222877956,9200000.0,-11589.0,3311.0,1656.0,1656.0,25.0,0.0,0.0,-Infinity
 81544E-7,9300000.0,85.39999999999999,78.11360239611082,9300000.0,-11757.0,3359.0,1680.0,1680.0,25.0,0.0,0.0,-Infinity
 92432E-7,9400000.0,85.0,77.47744365982531,9400000.0,-11841.0,3383.0,1692.0,1692.0,25.0,0.0,0.0,-Infinity
 08526E-7,9500000.0,85.3,77.89839274706439,9500000.0,-11960.0,3417.0,1709.0,1709.0,25.0,0.0,0.0,-Infinity
 92083E-7,9600000.0,84.89999999999999,77.34827846916366,9600000.0,-12100.0,3457.0,1729.0,1729.0,25.0,0.0,0.0,-Infinity
 5001E-7,9700000.0,85.1,77.62678779233454,9700000.0,-12219.0,3491.0,1746.0,1746.0,25.0,0.0,0.0,-Infinity

Export as .txt file...

Evaluation Values

Measure	Current	Mean
Accuracy	84,90	83,30
Kappa	77,30	74,91
Ram-Hours	0,00	0,00
Time	670,60	481,87
Memory	0,01	0,01

Plot

Zoom in Y
Zoom out Y
Zoom in X
Zoom out X



References III



Japkowicz, N., Myers, C., e Gluck, M. A. (1995).

A novelty detection approach to classification.

In *IJCAI*, pages 518–523. Morgan Kaufmann.



Kosina, P. e Gama, J. (2015).

Very fast decision rules for classification in data streams.

Data Min. Knowl. Discov., 29(1):168–202.



Krawczyk, B., Minku, L. L., Gama, J., Stefanowski, J., e Wozniak, M. (2017).

Ensemble learning for data stream analysis: A survey.

Information Fusion, 37:132–156.



Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., e Leisch, F. (2014).

e1071: Misc Functions of the Department of Statistics (e1071), TU Wien.

R package version 1.6-4.



Pereira, P., Ribeiro, R. P., e Gama, J. (2014).

Failure prediction - an application in the railway industry.

In *Discovery Science - 17th International Conference, DS 2014, Bled, Slovenia, October 8-10, 2014*.

Proceedings, pages 264–275.



R Core Team (2014).

R: A Language and Environment for Statistical Computing.

R Foundation for Statistical Computing, Vienna, Austria.



Ribeiro, R. P., Pereira, P. M., e Gama, J. (2016).

Sequential anomalies: a study in the railway industry.

Machine Learning, 105(1):127–153.

