

Data Stream Mining

Leszek Rutkowski

Data stream mining



in cooperation with
Piotr Duda, Maciej Jaworski, Lena Pietruczuk

Czestochowa University of Technology, Poland



Data stream mining

- content

- Data streams – introduction to the topic
- Concept drift
- Various strategies of learning
- How to deal with concept drift?
- Data stream classification methods – short overview
- Decision trees for data streams (including new results 2016)
- Ensemble methods for data streams (including new results 2016)
- Probabilistic neural networks for stream data mining (including new results 2016)
- Final remarks and challenging problems
- References

Overview of the research in the Institute of Computational Intelligence at CUT

**flexible
neuro-fuzzy systems**

**relational
neuro-fuzzy systems**

**rough
neuro-fuzzy systems**

**neuro-fuzzy systems
of type 2**

**dynamic
neuro-fuzzy systems**

**ensembles of
neuro-fuzzy systems**

**gradient and nature-
inspired algorithms**

**soft computing
techniques**

**interpretability of
neuro-fuzzy systems**

**stream data mining
algorithms**

**algorithms for
intelligent control**

**algorithms for
computed tomography**

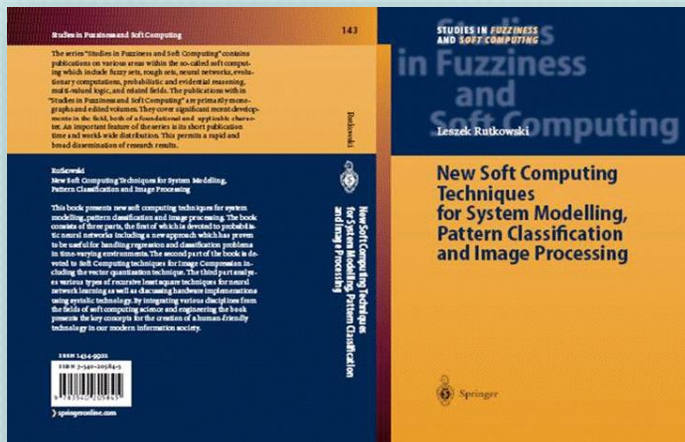
**behavioural
biometric systems**

**content based image
retrieval algorithms**

other areas

Overview of the research in the Institute of Computational Intelligence at CUT

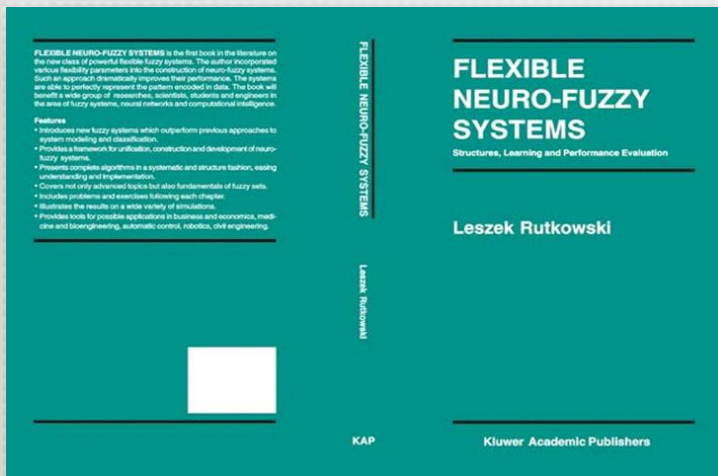
Rutkowski L., New Soft Computing Techniques for System Modelling, Pattern Classification and Image Processing, Springer, 2004.



- (1) Probabilistic neural networks to solve various problems of system modelling and classification in a non-stationary environment.
- (2) Image compression methods.
- (3) RLS learning algorithms for multilayer neural networks.
- (4) Systolic architectures of multilayer neural networks.

Overview of the research in the Institute of Computational Intelligence at CUT

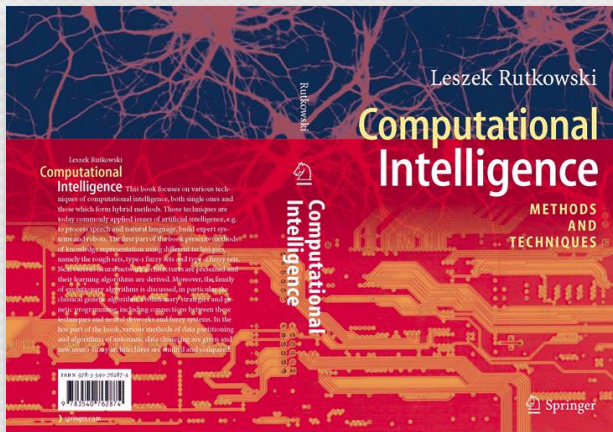
Rutkowski L., Flexible Neuro-Fuzzy Systems. Structures, learning and performance evaluation, Kluwer, 2004.



- (1) New fuzzy systems which outperform previous approaches to system modelling and classification.
- (2) Framework for unification, construction and development of neuro-fuzzy systems.
- (3) Complete algorithms in a systematic and structured fashion, easing understanding and implementation.

Overview of the research in the Institute of Computational Intelligence at CUT

Rutkowski L., Computational intelligence,
Springer, 2008.



- (1) Selected issues of artificial intelligence.
- (2) Methods of knowledge representation using rough sets.
- (3) Methods of knowledge representation using type-1 and type-2 fuzzy sets.
- (4) Neural networks and their learning algorithms.
- (5) Evolutionary algorithms.
- (6) Data clustering methods.
- (7) Neuro-fuzzy systems of Mamdani, logical and Takagi-Sugeno type.
- (8) Flexible neuro-fuzzy systems.

Overview of the research in the Institute of Computational Intelligence at CUT

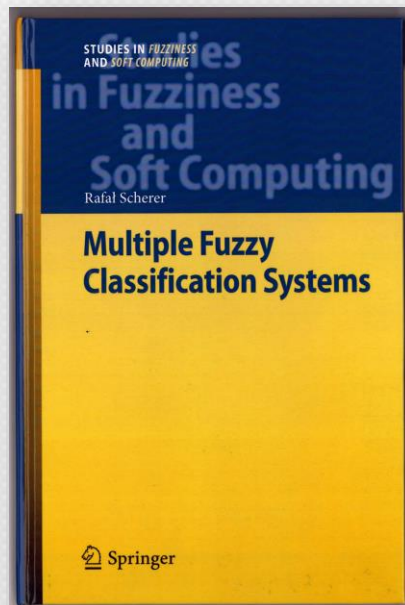
Rutkowski L., Duda P., Jaworski M., Pietruczuk L., Stream Data Mining: Algorithms and Their Probabilistic Properties, Studies in Big Data, Springer, 2017.

This book shows methods and algorithms which are mathematically justified.

- (1) It shows how to adopt the static decision trees, like ID.3 or CART, to deal with data streams.
- (2) A new technique, based on the McDiarmid bound, is developed.
- (3) New decision trees are designed, by a proper combination of the Gini index and misclassification error impurity measures, leading to the original concept of the hybrid decision trees.
- (4) The problem of designing ensembles and automatic choosing their sizes is described and solved.
- (5) Nonparametric techniques based on the Parzen – kernels and orthogonal series, are adopted to deal with concept drift in the problem of non-stationary regressions and classification in time-varying environment. Nonparametric procedures are developed and their probabilistic properties are investigated.

Overview of the research in the Institute of Computational Intelligence at CUT

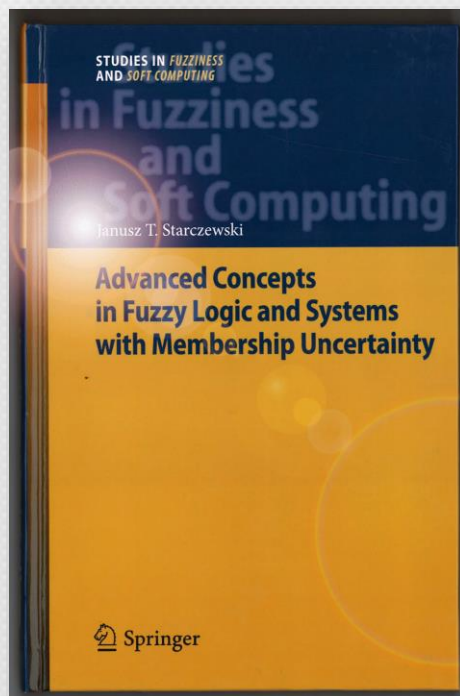
Scherer R., Multiple Fuzzy Classification Systems, Springer, 2012.



- (1) Ensemble techniques (bagging, boosting, negative correlation learning etc.).
- (2) Relational modular fuzzy systems.
- (3) Ensembles of the Mamdani, logical and Takagi-Sugeno fuzzy systems.
- (4) Rough–neuro–fuzzy ensembles for classification with missing data.

Overview of the research in the Institute of Computational Intelligence at CUT

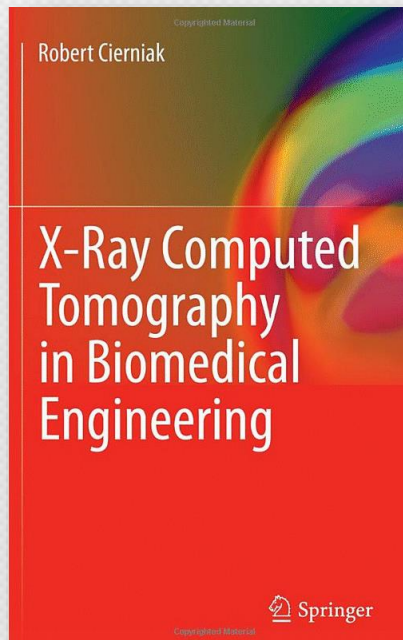
Starczewski J., Fuzzy Logic and Systems with Membership Uncertainty, Springer, 2013.



- (1) Algebraic operations on fuzzy valued fuzzy sets.
- (2) Defuzzification of uncertain fuzzy sets. Generalized uncertain fuzzy logic systems.
- (3) Uncertainty generation in uncertain fuzzy logic systems.
- (4) Designing uncertain fuzzy logic systems.

Overview of the research in the Institute of Computational Intelligence at CUT

Cierniak R., X-Ray Computed Tomography in Biomedical Engineering, Springer, 2011



- (1) Discusses X-ray CT tomography from a historical point of view, the design and physical operating principles of computed tomography apparatus, the algorithms of image reconstruction and the quality assessment criteria of tomography scanners.
- (2) Algorithms of image reconstruction from projections, a crucial problem in medical imaging, are considered in depth.

Data stream mining

- content

- **Data streams – introduction to the topic**
- Concept drift
- Various strategies of learning
- How to deal with concept drift?
- Data stream classification methods – short overview
- Decision trees for data streams (including new results 2016)
- Ensemble methods for data streams (including new results 2016)
- Probabilistic neural networks for stream data mining (including new results 2016)
- Final remarks and challenging problems
- References

Stream Data

Stream of data:

- huge volumes of continuous data,
- possibly infinite,
- multidimensional features,
- often fast changing,
- requiring fast, real-time responses.

Examples of data streams

- *Telecommunication calling records,*
 - *Credit card transaction flows,*
- *Network monitoring and traffic engineering,*
 - *Financial market,*
- *Audio and video recording of various processes,*
 - *Computer network information flow,*
 - *Web logs and Web page click streams,*
 - *Satelite data flow.*



Stream Data

Static data	Stream of data
Fixed number of data elements (unlimited memory usage)	Potentially infinite number of data elements (memory limitation problem)
Stationary distribution of data	Changing data distribution (concept drift)
All data available at any time – multiple passes of data	One pass of data
Unlimited processing time	Processing time depends on rate of incoming data

Typical problems in data mining

$$X_i = [x_i^1, \dots, x_i^d, y_i] = [\mathbf{x}_i, y_i], x_i^j \in A_j, y_i \in B$$

$A_1 \times \dots \times A_d$ - space of attributes values

X_1, \dots, X_N from probability distribution $q(\mathbf{x}, y)$

Two types of supervised learning:

- 1) **Regression** - B continuous
- 2) **Classification** - B nominal

Typical problems in data mining

- Regression

X_1, \dots, X_N derived from $\varrho(\mathbf{x}, y)$, $B \subset \mathbb{R}$

Aim: Construct a mapping function:

$$f: A_1 \times \dots \times A_d \rightarrow B,$$

such that for any $\tilde{X} = [\tilde{\mathbf{x}}, \tilde{y}]$ from $\varrho(\mathbf{x}, y)$:

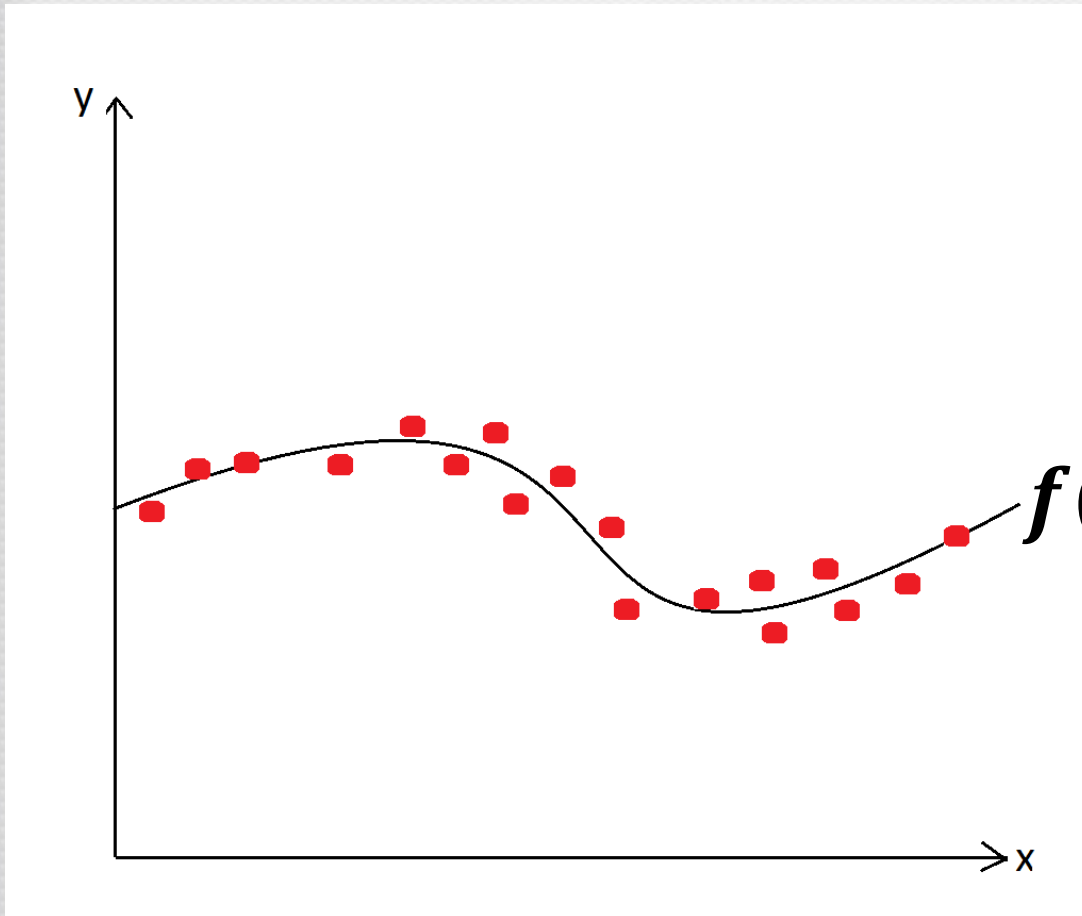
$$f = \underset{f^*}{\operatorname{argmin}} \{E[L(f^*(\tilde{\mathbf{x}}), \tilde{y})]\},$$

where L is a loss function, e.g.

$$L(f^*(\tilde{\mathbf{x}}), \tilde{y}) = (f^*(\tilde{\mathbf{x}}) - \tilde{y})^2$$

Typical problems in data mining

- Regression



one-dimensional
regression problem

Typical problems in data mining

- Classification

X_1, \dots, X_N derived from $\varrho(\mathbf{x}, y)$, $B = \{C_1, \dots, C_K\}$

Aim: Construct a mapping function:

$$f: A_1 \times \dots \times A_d \rightarrow B,$$

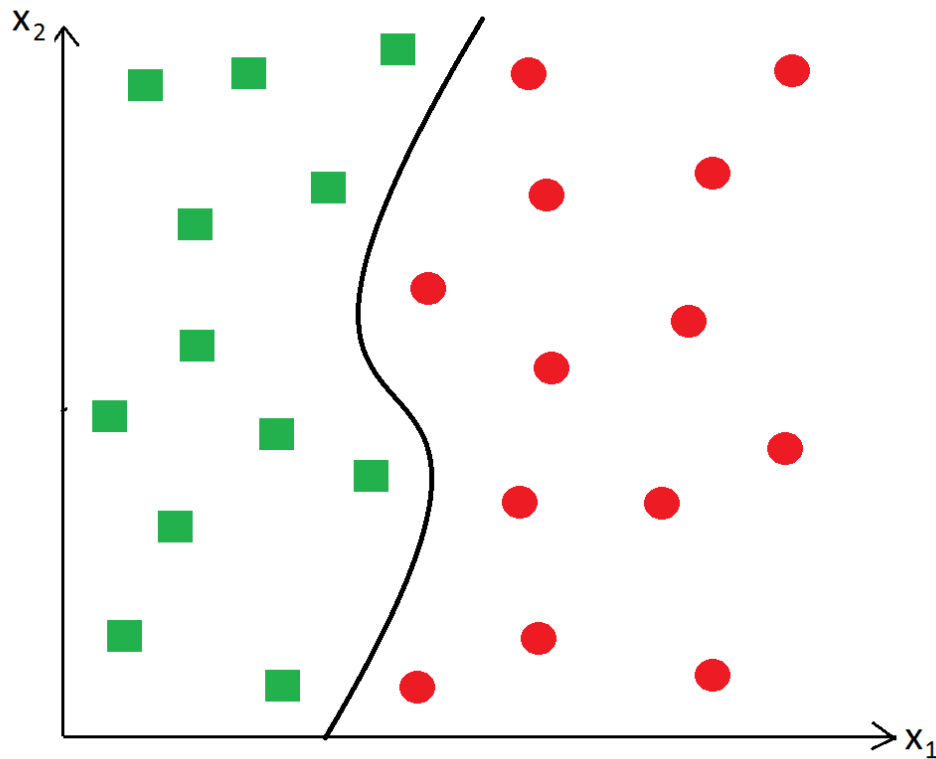
such that for any $\tilde{X} = [\tilde{\mathbf{x}}, \tilde{y}]$ from $\varrho(\mathbf{x}, y)$:

$$f = \underset{f^*}{\operatorname{argmin}} \{E[L(f^*(\tilde{\mathbf{x}}), \tilde{y})]\},$$

$$\text{where } L(f^*(\tilde{\mathbf{x}}), \tilde{y}) = \begin{cases} 0, & f^*(\tilde{\mathbf{x}}) = \tilde{y} \\ 1, & f^*(\tilde{\mathbf{x}}) \neq \tilde{y} \end{cases}$$

Supervised learning: Classification

$$g(x) = 0$$



$$f(x) = \begin{cases} \bullet, & \text{sgn}(g(x)) \geq 0 \\ \blacksquare, & \text{sgn}(g(x)) < 0 \end{cases}$$

two-dimensional
classification problem

Data stream mining

- content

- Data streams – introduction to the topic
- **Concept drift**
- Various strategies of learning
- How to deal with concept drift?
- Data stream classification methods – short overview
- Decision trees for data streams (including new results 2016)
- Ensemble methods for data streams (including new results 2016)
- Probabilistic neural networks for stream data mining (including new results 2016)
- Final remarks and challenging problems
- References

Stream Data

- concept drift

CONCEPT DRIFT - it is a phenomenon describing the change in data distribution, character or their meaning, e.g. assigning e-mails to the "spam" category.



Stream Data

- concept drift

In data streams the joint probability distribution $q(\mathbf{x}, y)$ can vary over time:

X_1 from $q_1(\mathbf{x}, y)$;
 X_2 from $q_2(\mathbf{x}, y)$;
..... ;
 X_N from $q_N(\mathbf{x}, y)$;

Stream Data

- concept drift

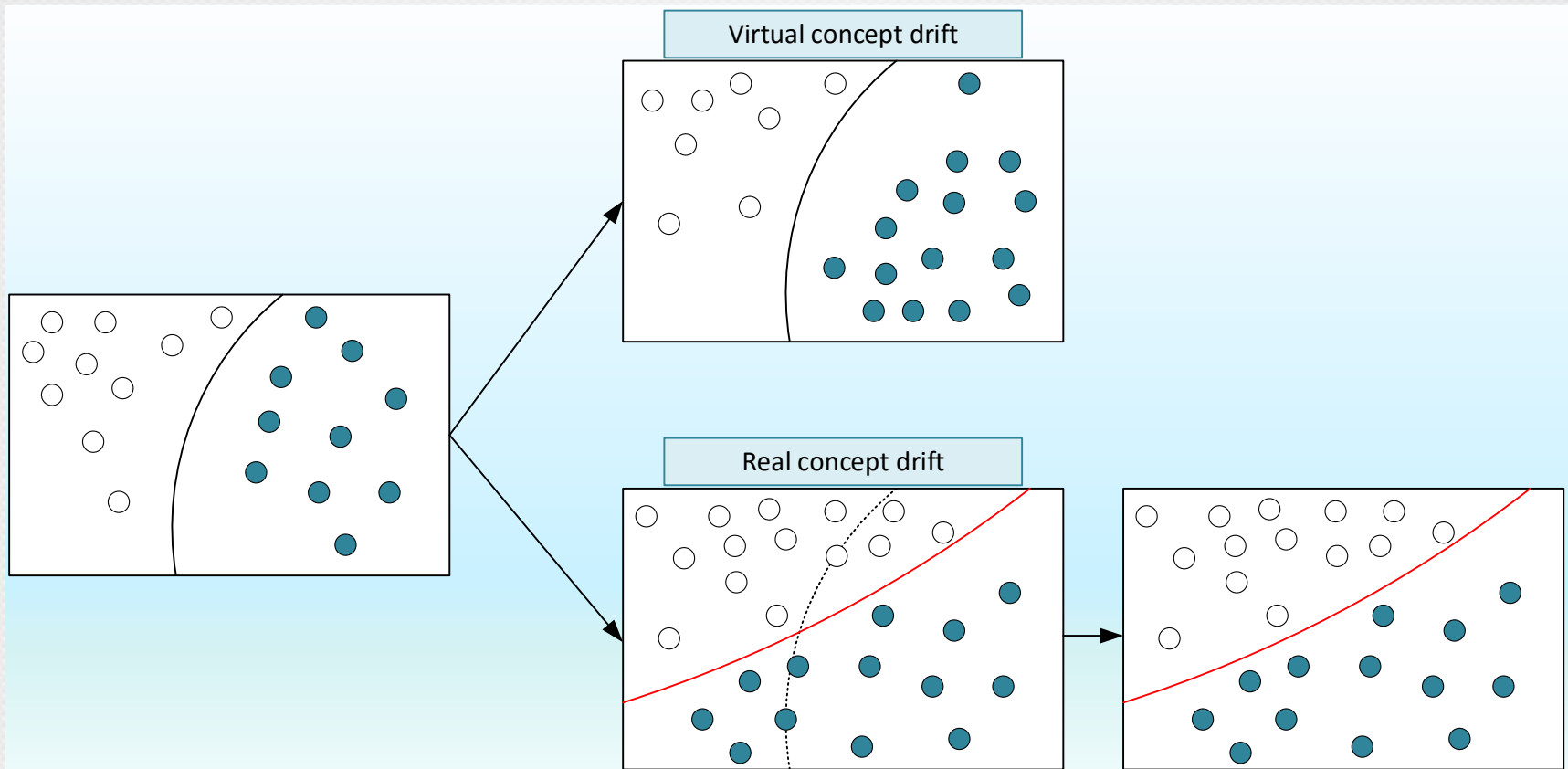
$q_i(\mathbf{x}, y)$ - joint probability distribution

$$q_i(\mathbf{x}, y) = \hat{q}_i(y|\mathbf{x})\tilde{q}_i(\mathbf{x})$$

Type of change:	drift:
$\hat{q}_i(y \mathbf{x})$	real
$\tilde{q}_i(\mathbf{x})$	virtual

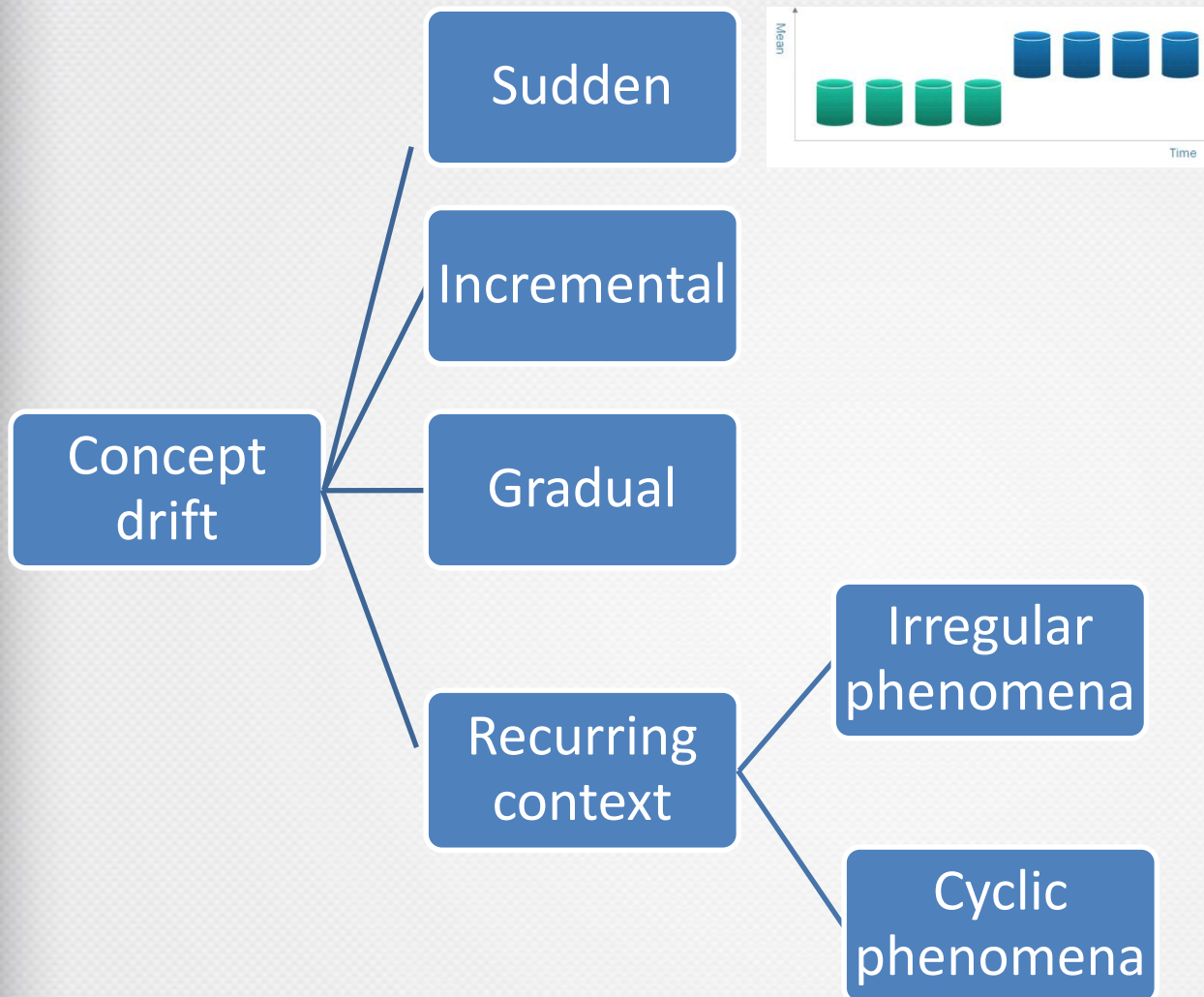
Stream Data

- concept drift



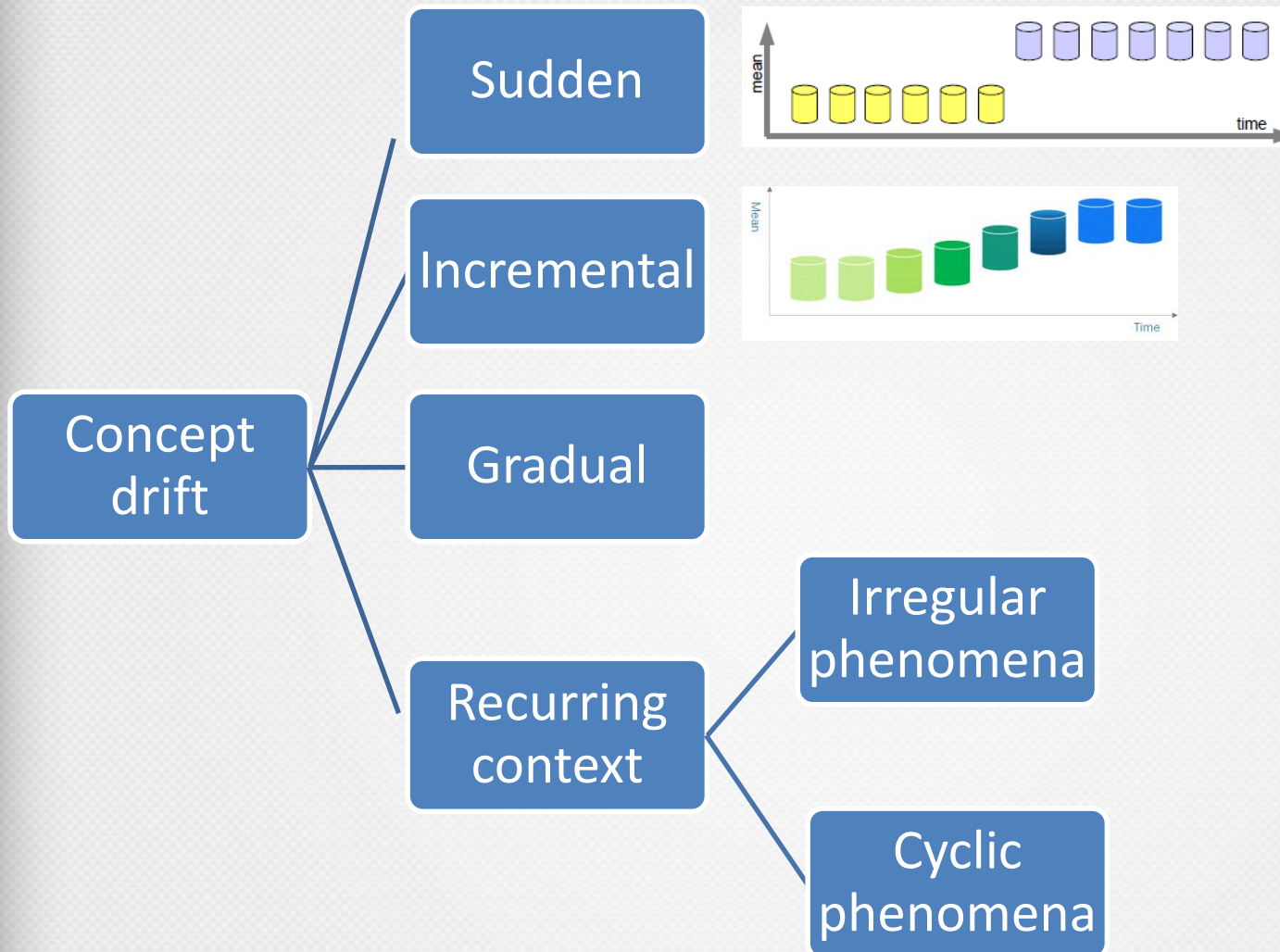
Stream Data

- concept drift



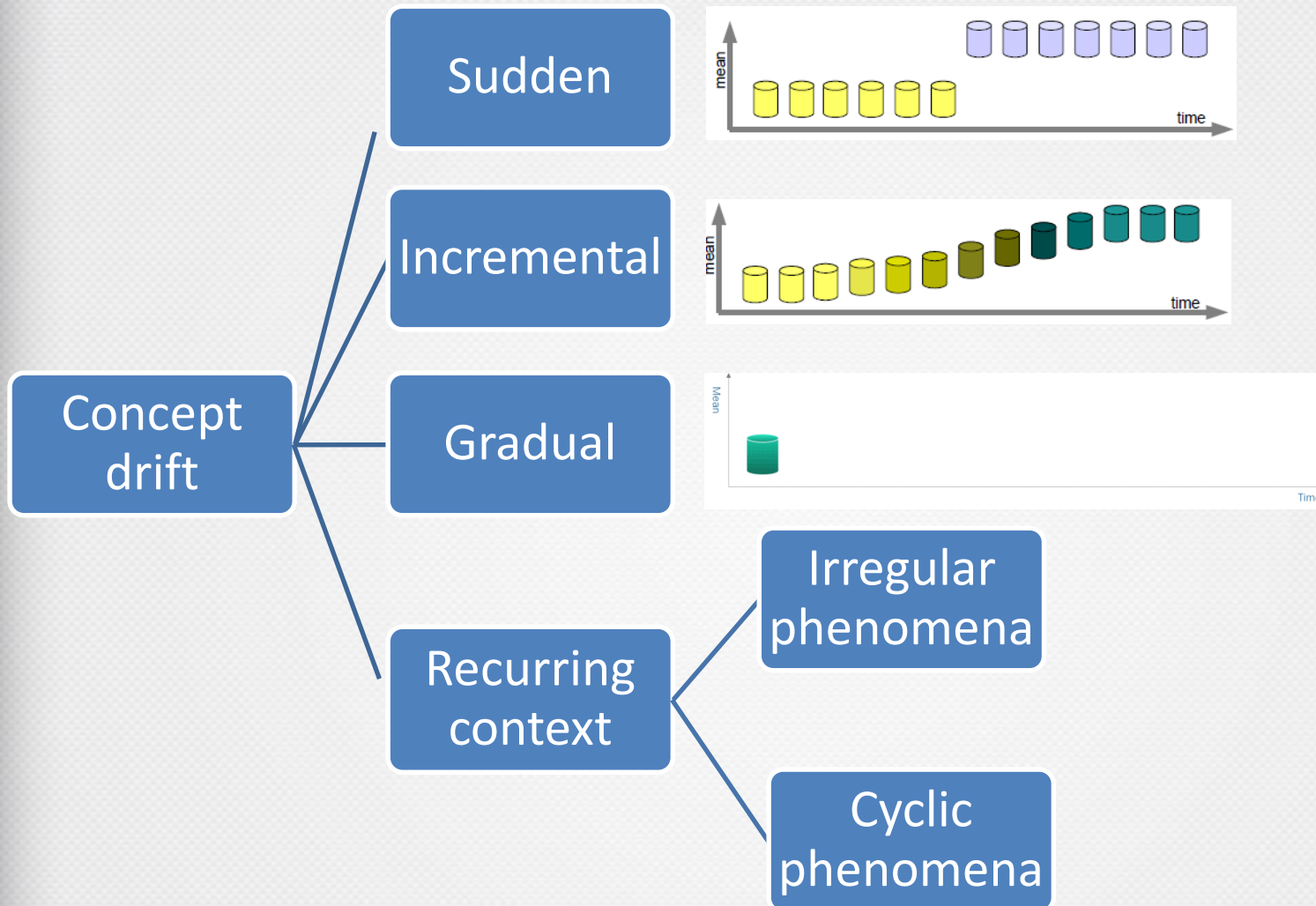
Stream Data

- concept drift



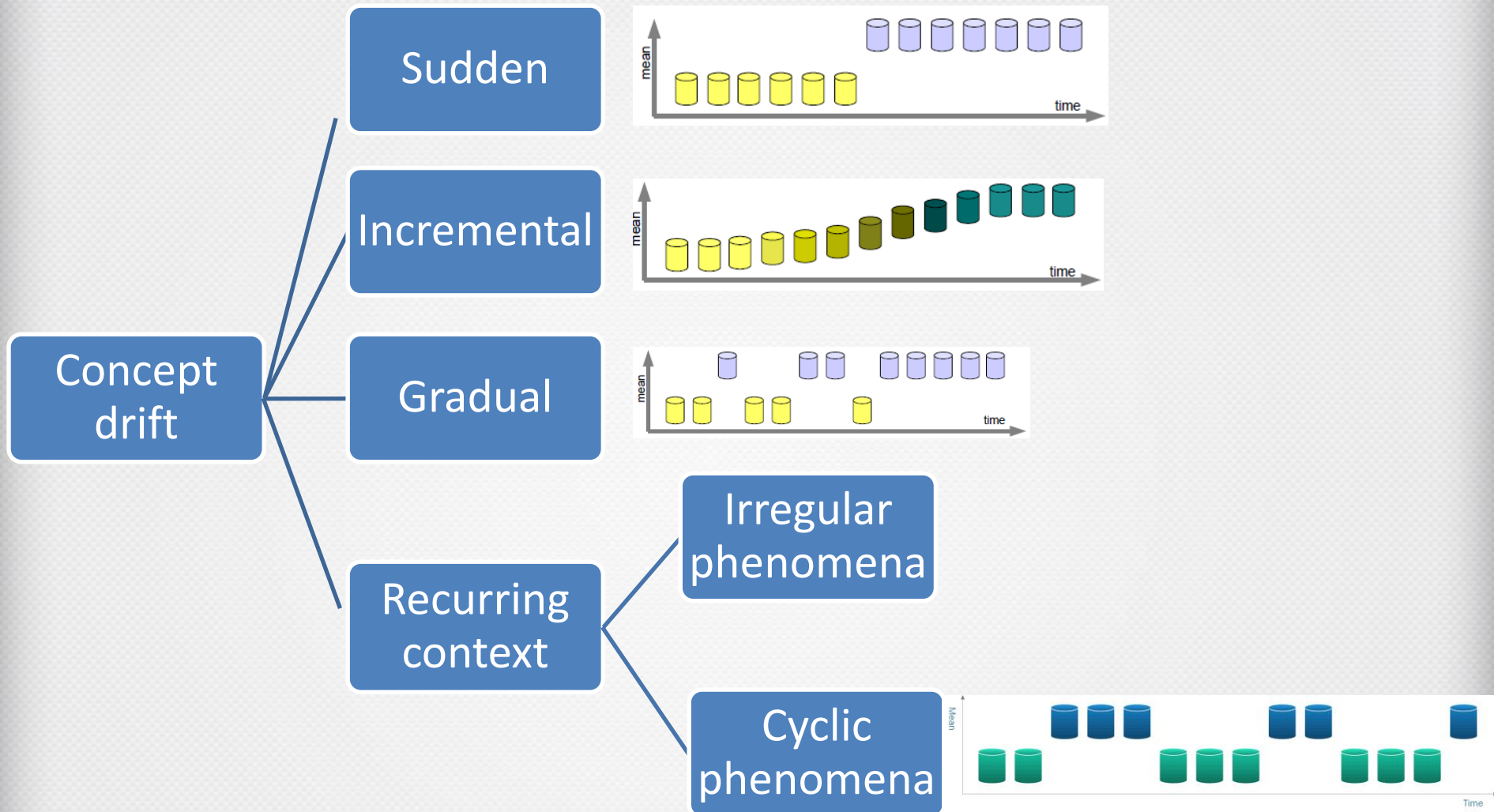
Stream Data

- concept drift



Stream Data

- concept drift



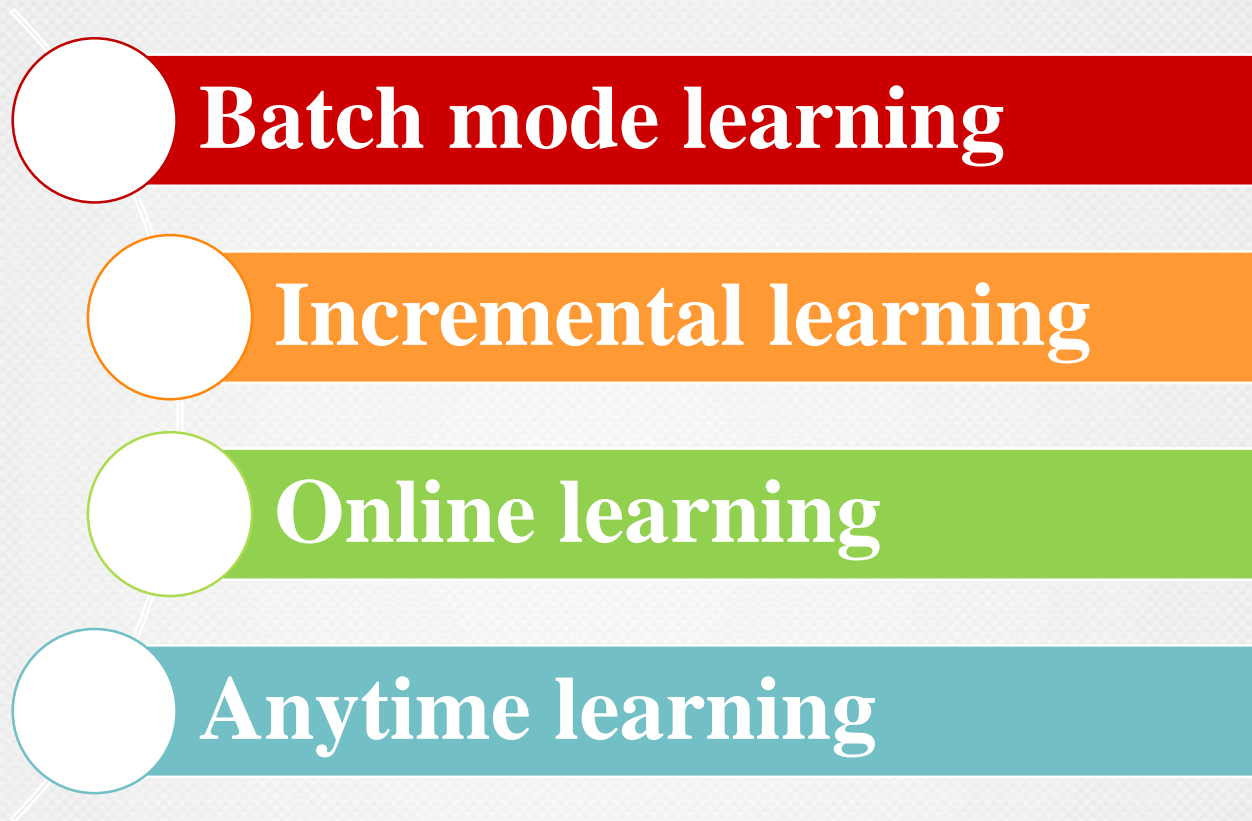
Data stream mining

- content

- Data streams – introduction to the topic
- Concept drift
- **Various strategies of learning**
- How to deal with concept drift?
- Data stream classification methods – short overview
- Decision trees for data streams (including new results 2016)
- Ensemble methods for data streams (including new results 2016)
- Probabilistic neural networks for stream data mining (including new results 2016)
- Final remarks and challenging problems
- References

Different kinds of learning^{*}

(regarding time constraints and examples availability)



^{*} V. Lemaire, Ch. Slaperwyck, A. Bondu, *A Survey on Supervised Classification on Data Streams*, Lecture Notes in Business Information Processing, vol. 205, pp. 88-125, Springer, 2015

Batch mode learning

Batch mode learning:

consists of learning a model from a representative dataset which is fully available at the beginning of the learning stage. This type of algorithm is not appropriate for stream data.

Incremental learning

Incremental learning:

consists of receiving and integrating new examples without the need to perform a full learning phase from scratch. For any examples X_1, X_2, \dots , it generates the hypotheses f_1, f_2, \dots , such that f_{n+1} depends only on f_n and the current example x_{n+1} .

Online learning

Online learning:

an incremental learning for which the examples continuously arrive from data stream. The requirements in terms of time complexity are stronger than for the incremental learning. Concept drift must be managed by algorithms of this type.

Anytime learning

Anytime learning:

algorithm is able to maximize the quality of the learned model (with respect to a particular evaluation criterion) until an interruption (which may be the arrival of a new example).

Learning on data streams

Desired properties*

	(1)	(2)	(3)
incremental	x	x	
read data only once	x	x	x
memory management	x	x	x
anytime	x	x	x
deal with concept drift		x	

(1) Fayyad, U.M. et al.: *Advances in Knowledge Discovery and Data Mining*. American Association for Artificial Intelligence, Menlo Park, CA, USA (1996);

(2) Hulten, G., Spencer, L., Domingos, P.: Mining time-changing data streams. In: *Proceedings of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, NY, USA (2001) 97-106;

(3) Stonebraker, M. Çetintemel, U., Zdonik, S.: The 8 requirements of real-time stream processing. *ACM SIGMOD Record* 34(4) (December 2005) 42-47;

* V. Lemaire, Ch. Slaperwyck, A. Bondu, *A Survey on Supervised Classification on Data Streams*, *Lecture Notes in Business Information Processing*, vol. 205, pp. 88-125, Springer, 2015

Data stream mining

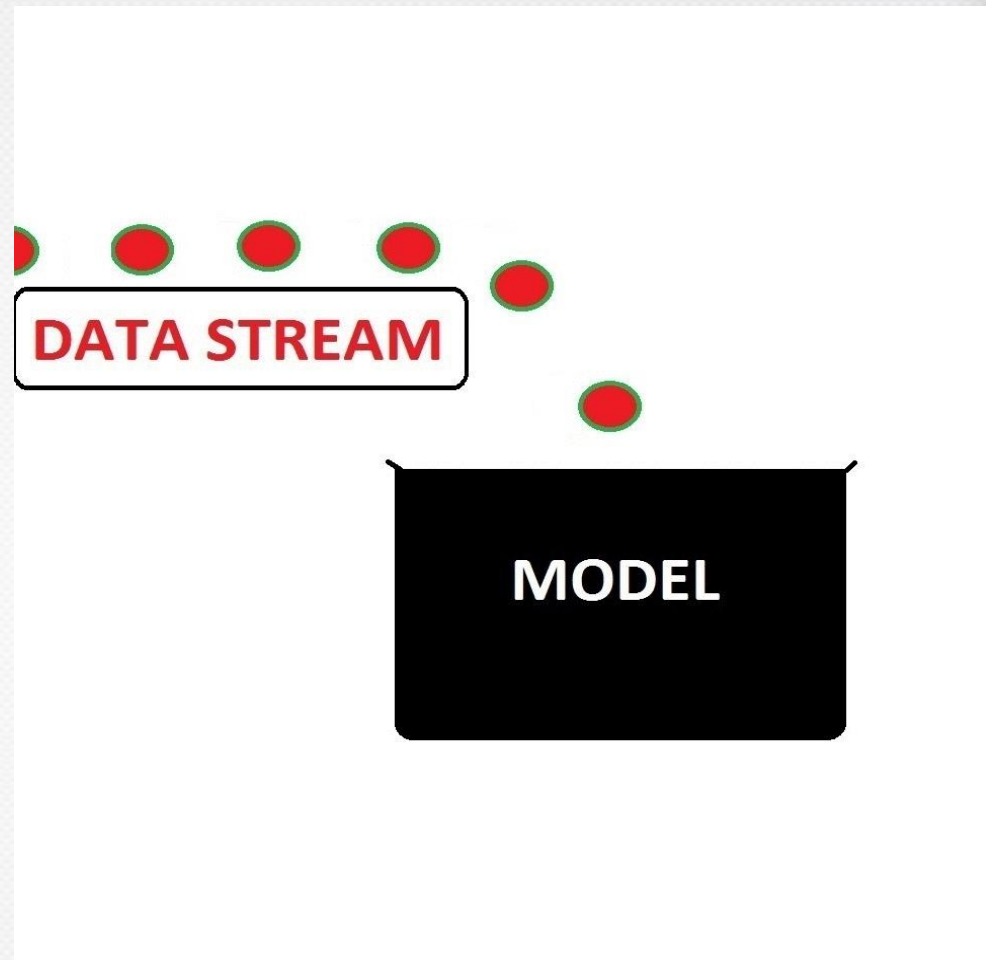
- content

- Data streams – introduction to the topic
- Concept drift
- Various strategies of learning
- **How to deal with concept drift?**
- Data stream classification methods – short overview
- Decision trees for data streams (including new results 2016)
- Ensemble methods for data streams (including new results 2016)
- Probabilistic neural networks for stream data mining (including new results 2016)
- Final remarks and challenging problems
- References

How to process stream data?

Online processing (instant incremental)

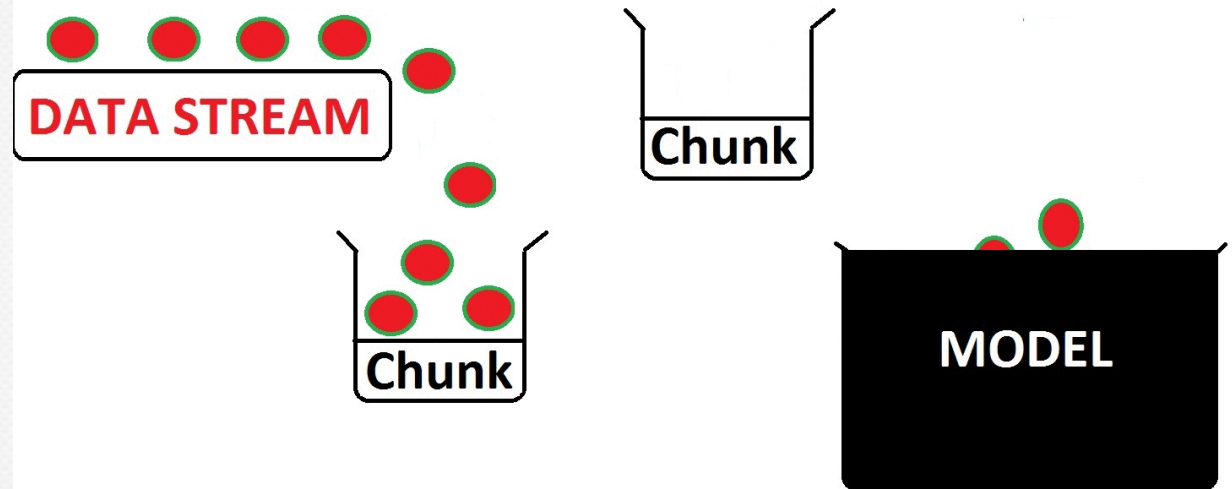
- The model (classifier) is updated after reading every single data element
- Updating must be completed before the next data element arrives



How to process stream data?

Block (chunk) processing

- The model (classifier) is updated after reading every chunk of data
- Updating must be completed before the next chunk of data is collected



How to deal with concept drift?

How to deal with concept drift?

Sliding windows

Goal: limit the number of training examples to the most recent ones, consequently eliminate examples coming from an old concept

Two approaches:

- window of a fixed sized
 - classifier based on a small window reacts quickly to concept drifts but loses an accuracy in periods of stationary data
 - Classifier based on a large window fails to react quickly to concept drift
- window size adjusted heuristically (e.g. ADWIN)

* Alexey Tsymbal, Mykola Pechenizkiy, Padraig Cunningham, and Seppo Puuronen. Dynamic integration of classifiers for handling concept drift. Information Fusion, vol. 9, no 1, pp. 56–68, 2008.

How to deal with concept drift?

Sliding windows

Algorithm Basic windowing algorithm

Require: S_∞ - data stream, W - window of examples

- 1: initialize window W
 - 2: **for all** examples X_i in S_∞ ($i = 1, 2, \dots$) **do**
 - 3: $W \leftarrow W \cup \{X_i\}$
 - 4: if necessary remove outdated examples from W
 - 5: rebuild/update C using W
 - 6: **end for**
 - 7: **return** a classifier built on examples in window W
-

How to deal with concept drift?

Weighted windows

The oldest examples are discarded by using a decay function assigning a weight to each example, e.g.

$$w_1(t) = e^{-\lambda t}, \lambda > 0,$$

$$w_2(t) = t^{-\alpha}, \alpha > 0,$$

$$w_3(t) = 1 - \frac{t}{|W|},$$

where t represents the age of an example ($t = 0$ for a new example and $t = |W| - 1$ for the last one).

* Edith Cohen, Martin J. Strauss, Maintaining time-decaying stream aggregates, J. Algorithms, vol. 59, no. 1, pp. 19–36, 2006.

How to deal with concept drift?

Weighted windows

Algorithm Weighted windows

Require: S_∞ - data stream, d - window size, $w(\cdot)$ weight function

```
1: for all examples  $X_i$  in  $S_\infty$  ( $i = 1, 2, \dots$ ) do
2:   if  $|W| = d$  then
3:     remove the oldest example from  $W$ 
4:   end if
5:    $W \leftarrow W \cup \{X_i\}$ 
6:   for all examples  $X^j$  in  $S_\infty$  ( $j = 1, 2, \dots, d$ ) do
7:     calculate example's weight  $w(X^j)$ 
8:   end for
9: end for
10: return  $W$  - a window of examples
```

How to deal with concept drift?

(a family of algorithms called FISH)

The FISH* algorithms are family of methods that take advantage of similarities between data elements both in time and space. Based on distances in space $d^{(s)}_{ij}$ and in time $d^{(t)}_{ij}$ we calculate the distance between X_i and X_j as follows:

$$D_{ij} = a_1 d^{(s)}_{ij} + a_2 d^{(t)}_{ij}$$

where a_1 and a_2 are the weight coefficients.

* Indre Žliobaite, Combining time and space similarity for small size learning under concept drift, In Foundations of Intelligent Systems, volume 5722 of Lecture Notes in Computer Science, pp.412–421, Springer Berlin Heidelberg, 2009.

** Indre Žliobaite, Adaptive Training Set Formation, PhD thesis, Vilnius University, Lithuania, 2010.

How to deal with concept drift?

(a family of algorithms called FISH)

In the first version of the algorithm, the number of data instances in the created training sample is determined by the user, while in FISH2** this number is not fixed. Moreover, FISH* builds separate data sets for each class, while the second version takes data from all classes based only on the closeness of the data elements. In the third version of this algorithm (FISH3**), a search for the best weight coefficients a_1 and a_2 is being conducted. In this case, the size of the data window is established dynamically.

* Indre Žliobaite, Combining time and space similarity for small size learning under concept drift, In Foundations of Intelligent Systems, volume 5722 of Lecture Notes in Computer Science, pp.412–421, Springer Berlin Heidelberg, 2009.

** Indre Žliobaite, Adaptive Training Set Formation, PhD thesis, Vilnius University, Lithuania, 2010.

How to deal with concept drift?

(a family of algorithms called FISH)

Algorithm The Instance Selection Algorithm (FISH2)

Require: S - set of data elements X_1, X_2, \dots, X_t , target observation x_{t+1} with unknown label, neighbourhood size k , time/space similarity weight A

- 1: **for** $j \in \{1, \dots, t\}$ **do**
- 2: calculate distance $D_{t+1,j}$
- 3: **end for**
- 4: sort the distances from minimum to maximum: $D_{z1,t+1} < D_{z2,t+1} < \dots < D_{zt,t+1}$
- 5: **for** $s = k : step : t$ **do**
- 6: select s instances having the smallest distance D
- 7: using cross-validation* build a classifier τ_s using the instances indexed $\{i_1, \dots, i_s\}$ and test it on the k nearest neighbours indexed $\{i_1, \dots, i_k\}$, record testing error e_s
- 8: **end for**
- 9: find the minimum error classifier C_L , where $L = \arg \min_{L=k}^t (e_L)$
- 10: output the instances $\{i_1, \dots, i_L\}$

*when test on the instance X_{zk} , this data is excluded from the validation set

How to deal with concept drift?

(a family of algorithms called FISH)

Algorithm The Instance Selection Algorithm (FISH3)

Require: S - set of data elements X_1, X_2, \dots, X_t , target observation x_{t+1} with unknown label, neighbourhood size k

- 1: **for** $J = 0 : \text{step 1}, a_1 = j, a_2 = 1 - a_1$ every time and space proportion **do**
- 2: calculate distance $D_{ij,t+1}^j = a_1 d_{i,t+1}^s + a_2 d_{i,t+1}^t$ for $i = 1, \dots, t$
- 3: sort the distances from minimum to maximum: $D_{jz1,t+1}^j < D_{jz2,t+1}^j < \dots < D_{jzt,t+1}^j$
- 4: **for** $N = k : \text{step 2} : t$ select the training set size **do**
- 5: pick N instances having the smallest distances D^j
- 6: using cross-validation[†] build a classifier τ_s^{jN} using the instances $(X_{jz1}, X_{jz2}, \dots, X_{jzN})$ and the training set
- 7: test τ_s^{jN} on the k nearest neighbours $(X_{jz1}, X_{jz2}, \dots, X_{jzk})$, record testing error e_N^j
- 8: **end for**
- 9: **end for**
- 10: find the minimum error classifier τ_s^{jN*} , where $jN* = \arg \min_{j=0}^1 \min_{N=k}^t (e_N^j)$
- 11: output the indexes $\{jz1, \dots, jzN*\}$

[†]when test on the instance X_{jzk} , this data is excluded from the validation set

How to deal with concept drift?

- the ADWIN

This algorithm was obtained to regulate the size of data elements windows. It uses hypothesis that is based on differences of mean value of windows W_0 and W_1 compared with the threshold value α . If we assume that $|W_0|$ and $|W_1|$ denote the number of data elements in windows W_0 and W_1 respectively, then we have to find ϵ_{cut} such that $|\mu\hat{W}_0 - \mu\hat{W}_1| < \epsilon_{cut}$. The value of ϵ_{cut} is determined as follows:

$$\epsilon_{cut} = \sqrt{\frac{1}{2m} \ln \frac{4}{\delta'}}$$
$$m = \frac{1}{\frac{1}{|W_0|} + \frac{1}{|W_1|}}$$
$$\delta' = \frac{\delta}{|W_0| + |W_1|}$$

In this method we compare every possible split of entire window into windows W_0 and W_1 .

* Bifet A., Avalda R., Kalman filters and adaptive windows for learning in data streams. In Proc. of the 9th int. conf. on Discovery science, DS. 29–40, 2006.

How to deal with concept drift?

- the ADWIN

Algorithm ADWIN: Adaptive Windowing Algorithm

```
1: initialize window  $W$ 
2: for all  $t > 0$  do
3:    $W \leftarrow W \cup \{X_t\}$  (i.e., add  $X_t$  to the head of  $W$ )
4:   repeat
5:     drop elements from the tail of  $W$ 
6:   until  $|\hat{\mu}_{W_0} - \hat{\mu}_{W_1}| \geq \epsilon_{cut}$  holds
7:   for every split of  $W$  into  $W = W_0 \cup W_1$  output  $\hat{\mu}_W$ 
8: end for
```

* Bifet A., Avalda R., Kalman filters and adaptive windows for learning in data streams. In Proc. of the 9th int. conf. on Discovery science, DS. 29–40, 2006.

How to deal with concept drift?

- the Drift Detection Method (DDM)

In the DDM algorithm the authors noticed that the class assigned by a classifier can be either true or false. Therefore they model the number of classification errors with a binomial distribution. Let p_i denote the probability of a false prediction, then the standard deviation is calculated as follows:

$$s_i = \sqrt{\frac{p_i(1 - p_i)}{i}}$$

For a sufficiently large number of examples ($n > 30$), the binomial distribution can be approximated by a Gaussian distribution with the same mean and variance. The error rate is monitored by updating two registers: p_{min} and s_{min} .

- **warning level**

$$p_i + s_i \geq p_{min} + \alpha \cdot s_{min}$$

- **alarm level**

$$p_i + s_i \geq p_{min} + \beta \cdot s_{min}$$

How to deal with concept drift?

- the Drift Detection Method (DDM)

Algorithm DDM: Drift Detection Method

Require: S_∞ - stream of data elements X_1, X_2, X_3, \dots , values of parameters α and β

```
1: initialize  $p_i, s_i, p_{min}, s_{min}, p_{min} + s_{min}, W \leftarrow \emptyset, W' \leftarrow \emptyset$ ,
2: for all new data elements  $X_i$  do
3:    $W \leftarrow W \cup \{X_i\}$ 
4:   update the value of  $p_i$  and  $s_i$  for  $W$ 
5:   if  $i > 30$  then
6:     if  $p_i + s_i < p_{min} + s_{min}$  then
7:       update  $p_{min}$  and  $s_{min}$ 
8:     end if
9:     if  $p_i + s_i \geq p_{min} + \alpha s_{min}$  then
10:       $W' \leftarrow W' \cup \{X_i\}$ 
11:      if  $p_i + s_i \geq p_{min} + \beta s_{min}$  then
12:         $W \leftarrow W'$ 
13:         $W' \leftarrow \emptyset$ 
14:      end if
15:    end if
16:  end if
17: end for
```

How to deal with concept drift?

- the Early Drift Detection Method (EDDM)

The main difference between DDM* and EDDM** is the definition of parameter that is being investigated. In this method p'_i defines the distance between two errors and s'_i is its standard deviation. It is expected that with increasing accuracy of the system this distance will increase.

- **warning level**

$$\frac{p'_i + 2s'_i}{p'_{max} + 2s'_{max}} < \alpha$$

- **alarm level**

$$\frac{p'_i + 2s'_i}{p'_{max} + 2s'_{max}} < \beta$$

In the article authors proposed to set values of α and β to 2 and 3, respectively, which represents 95% and 90% of the distribution.

*João Gama, Pedro Medas, Gladys Castillo, and Pedro Rodrigues, Learning with drift detection, In AnaL.C. Bazzan and Sofiane Labidi, editors, Advances in Artificial Intelligence – SBIA 2004, volume 3171 of Lecture Notes in Computer Science, pp. 286–295, Springer Berlin Heidelberg, 2004.

** Manuel Baena-García, José del Campo-Ávila, Raúl Fidalgo, Albert Bifet, Ricard Gavaldá, and Rafael Morales-Bueno. Early drift detection method, In Fourth International Workshop on Knowledge Discovery from Data Streams, 2006.

How to deal with concept drift?

- the Early Drift Detection Method (EDDM)

Algorithm EDDM: Early Drift Detection Method

Require: S_∞ - stream of data elements X_1, X_2, X_3, \dots , value of parameters α and β

```
1: initialize  $p'_i, s'_i, p'_{min}, s'_{min}, p'_{min}, s'_{min}, p'_{min} + s'_{min}, W \leftarrow \emptyset, W' \leftarrow \emptyset$ ,
2: for all new data elements  $X_i$  do
3:    $W \leftarrow W \cup \{X_i\}$ 
4:   update the value of  $p'_i$  and  $s'_i$  for  $W$ 
5:   if more than 30 errors are spotted then
6:     if  $p'_i + s'_i < p'_{min} + s'_{min}$  then
7:       update  $p_{min}$  and  $s_{min}$ 
8:     end if
9:     if  $(p'_i + 2s'_i)/(p'_{max} + 2s'_{max}) < \alpha$  then
10:       $W' \leftarrow W' \cup \{X_i\}$ 
11:      if  $(p'_i + 2s'_i)/(p'_{max} + 2s'_{max}) < \alpha$  then
12:         $W \leftarrow W'$ 
13:         $W' \leftarrow \emptyset$ 
14:      end if
15:    end if
16:  end if
17: end for
```

* Manuel Baena-García, José del Campo-Ávila, Raúl Fidalgo, Albert Bifet, Ricard Gavaldá, and Rafael Morales-Bueno. Early drift detection method, In Fourth International Workshop on Knowledge Discovery from Data Streams, 2006.

How to deal with concept drift?

- the Page-Hinkley test

The Page-Hinkley test was originally used as a sequential analysis technique for change detection in signal processing. Recently it has been proposed as a drift detector*. It allows to efficiently detect changes in the normal behavior of a process established by a model. The cumulative variable U_T of this test is defined as the cumulative difference between the observed values X_i and their mean up to the current moment in time:

$$U_T = \sum_{i=1}^T (X_i - \bar{X}_T - \delta)$$

where $\bar{X}_T = 1/t \sum_{i=1}^t X_i$ and δ corresponds to the magnitude of changes that are allowed. In the drift detection we treat the classifier's error rate as the observed value.

* João Gama, Raquel Sebastião, Pedro P. Rodrigues, On evaluating stream learning algorithms, Machine Learning, vol. 90, no. 3, pp. 317–346, 2013

How to deal with concept drift?

- the Page-Hinkley test

The minimal U_T is defined as

$$U_T^{min} = \min\{U_T; i = 1, \dots, t\}.$$

The PH test calculates the difference between U_T^{min} and U_T

$$PH_T = U_T - U_T^{min}.$$

If this difference is higher than a user specified threshold λ , a change is flagged. The threshold λ depends on the admissible false alarm rate. Increasing its value will entail fewer false alarms, but might miss or delay some changes. Controlling this detection threshold parameter makes it possible to establish a trade-off between the false alarms and the miss detections.

How to deal with concept drift?

- the Welch's t -test

This test applies on two samples of size n_1 and n_2 and is an adaptation of the Student's t test. This test is used to statistically test the null hypothesis that the means of two populations \bar{X}_1 and \bar{X}_2 , with unequal variances (s_1^2 and s_2^2), are equal. The formula of this test is:

$$p\text{-value} = \frac{\bar{X}_1 - \bar{X}_2}{\left(\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right)}$$

The null hypothesis can be rejected depending on the p -value.

How to deal with concept drift?

- the Kolmogorov-Smirnov test

The Kolmogorov–Smirnov statistic determines a distance between the empirical distribution function of the sample and the reference distribution or between the empirical distribution functions of two samples. In the latter case we calculate

$$D_{n,n'} = \sup_x |F_{1,n}(x) - F_{2,n'}(x)|,$$

where $F_{1,n}(x)$ and $F_{2,n'}(x)$ are the empirical distribution functions of the first and the second sample respectively, and \sup is the supremum function.

The null hypothesis that the samples are drawn from the same distribution is rejected with a confidence α if

$$D_{n,n'} > c(\alpha) \sqrt{\frac{n + n'}{nn'}}.$$

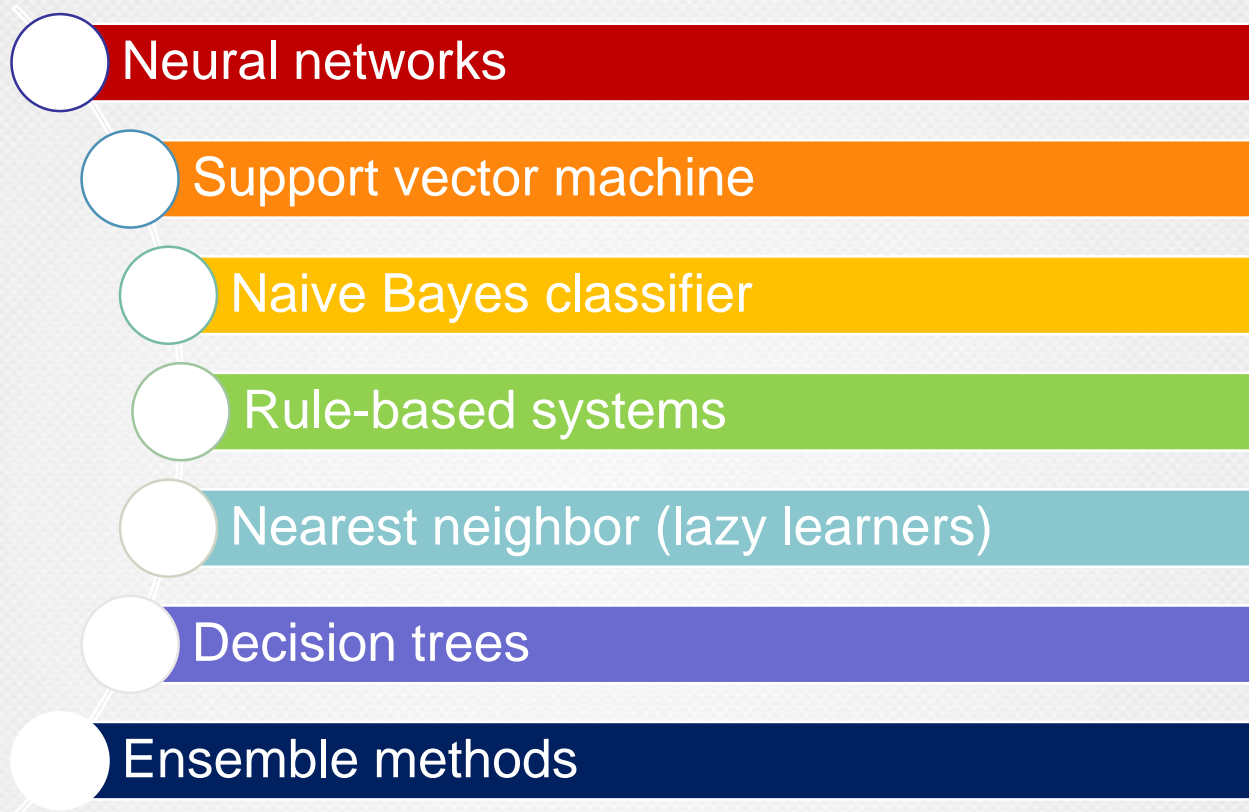
The value of $c(\alpha)$ can be found in the Kolmogorov-Smirnov table.

Data stream mining

- content

- Data streams – introduction to the topic
- Concept drift
- Various strategies of learning
- How to deal with concept drift?
- **Data stream classification methods – short overview**
- Decision trees for data streams (including new results 2016)
- Ensemble methods for data streams (including new results 2016)
- Probabilistic neural networks for stream data mining (including new results 2016)
- Final remarks and challenging problems
- References

Data stream classification methods



Data stream classification methods



Neural networks

- João Gama, Pedro P. Rodrigues, Stream-based electricity load forecast, In Proceedings of the 11th European Conference on Principles and Practice of Knowledge Discovery in Databases, vol. 4702, Lecture Notes in Computer Science, pp 446–453, Springer, 2007.
- Daniel Leite, Pyramo Costa Jr., Fernando Gomide, Evolving granular neural network for semi-supervised data stream classification, In Proceedings of the 2010 International Joint Conference on Neural Networks, pp. 1–8, IEEE, 2010.
- Daniel Leite, Pyramo Costa Jr., Fernando Gomide, Evolving granular neural networks from fuzzy data streams, Neural Networks, vol. 38, pp. 1–16, 2013.

Data stream classification methods

Support vector machine

- Domeniconi C., Gunopulos D., Incremental support vector machine construction, In ICDM, pp. 589–592, 2001.
- Syed N., Liu H., Sung K., Handling concept drifts in incremental learning with supportvector machines, Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining, ACM New York, NY, USA pp. 317–321, 1999.
- Fung G., Mangasarian O., Incremental support vector machine classification, Proceed-ings of the Second SIAM International Conference on Data Mining, Arlington, Virginia, pp. 247–260, 2002.
- Bordes A., Bottou L., The Huller: a simple and efficient online SVM, Proceedings of the16th European Conference on Machine Learning (ECML2005), 2005.
- Bordes A., Ertekin, S., Weston J., Bottou L., Fast kernel classiffiers with online and activelearning, Journal of Machine Learning Research 6, pp. 1579–1619, 2005.

Data stream classification methods

Naive Bayes classifier

- Remco R. Bouckaert, Voting massive collections of bayesian network classifiers for data streams. In Proceedings of the 19th Australian Joint Conference on Artificial Intelligence, vol. 4304, Lecture Notes in Computer Science, pp. 243–252, Springer, 2006.
- João Gama, Ricardo Rocha, Pedro Medas, Accurate decision trees for mining high-speed data streams, In Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 523–528, ACM, 2003.
- João Gama, Pedro P. Rodrigues, Stream-based electricity load forecast, In Proceedings of the 11th European Conference on Principles and Practice of Knowledge Discovery in Databases, vol. 4702, Lecture Notes in Computer Science, pp. 446–453, Springer, 2007.
- Richard Kirkby, Improving Hoeffding Trees, PhD thesis, Department of Computer Science, University of Waikato, 2007.

Data stream classification methods



Rule-based systems

- Magdalena Deckert, Jerzy Stefanowski, RILL: algorithm for learning rules from streaming data with concept drift, In Proceedings of the 21st ISMIS International Symposium on Foundations of Intelligent Systems, vol. 8502, Lecture Notes in Computer Science, pp. 20–29, Springer, 2014.
- Francisco J. Ferrer-Troyano, Jesús S. Aguilar-Ruiz, José Cristóbal Riquelme Santos, Discovering decision rules from numerical data streams, In Proceedings of the 2004 ACM Symposium on Applied Computing, pp. 649–653, ACM, 2004.
- Petr Kosina, João Gama, Handling time changing data with adaptive very fast decision rules, In Proceedings of the 2012 European Conference on Machine Learning and Knowledge Discovery in Databases, vol. 7523, Lecture Notes in Computer Science, pp. 827–842, Springer, 2012.
- Petr Kosina, João Gama, Very fast decision rules for classification in data streams, Data Min. Knowl. Discov., vol. 29, no. 1, pp. 168–202, 2015.
- Gerhard Widmer, M. Kubat, Learning in the presence of concept drift and hidden contexts, In Machine Learning, pp. 69–101, 1996.

Data stream classification methods



Nearest neighbor (lazy learners)

- Law Y., Zaniolo C., An adaptive nearest neighbor classification algorithm for data streams, Lecture notes in computer science, vol. 3721, pp. 108-120, 2005
- Beringer, J., Hüllermeier, E., Efficient instance-based learning on data streams, IntelligentData Analysis, vol. 11, no 6, pp. 627–650, 2007
- Shaker, A., Hüllermeier, E., Iblstreams: a system for instance-based classification and re-gression on data streams, Evolving Systems, vol. 3, no. 4, pp. 235–249, 2012
- Yan-Nei Law, Carlo Zaniolo, An adaptive nearest neighbor classification algorithm for data streams, In Proceedings of the 9th European Conference on Principles and Practice of Knowledge Discovery in Databases, vol. 3721, Lecture Notes in Computer Science, pp. 108–120, Springer, 2005.

Data stream classification methods



Decision trees

- P. Domingos and G. Hulten, Mining high-speed data streams, Proc. 6th ACM SIGKDD Internat. Conf. on Knowledge Discovery and Data Mining, pp. 71-80, 2000.
- G. Hulten, L. Spencer, P. Domingos, Mining time-changing data streams, In Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '01), ACM, New York, NY, USA, pp. 97-106, 2001
- Leszek Rutkowski, Maciej Jaworski, Lena Pietruczuk, Piotr Duda, A new method for data stream mining based on the misclassification error, IEEE Transaction on Neural Networks and Learning Systems, vol. 26, pp. 1048-1059, no. 5, 2015.
- Leszek Rutkowski, Lena Pietruczuk, Piotr Duda, Maciej Jaworski, Decision trees for mining data streams based on the McDiarmid's bound, IEEE Transactions on Knowledge and Data Engineering, vol. 25, no. 6, pp. 1272-1279, June 2013.

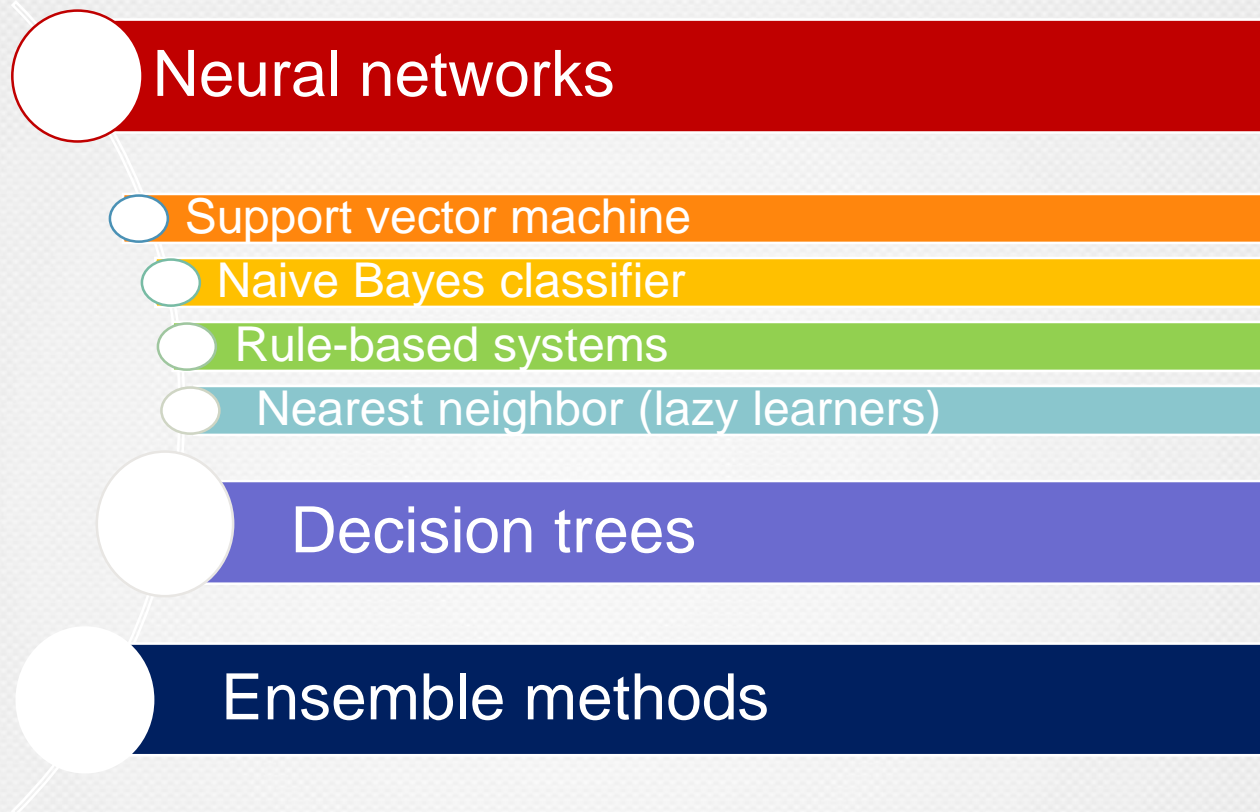
Data stream classification methods



Ensemble methods

- Brzezinski D., Stefanowski J., Combining block-based and online methods in learning ensembles from concept drifting data streams, Information Sciences, vol. 265, pp. 50-67, 2014.
- Dewan Md. Farid, Li Zhang, Alamgir Hossain, Chowdhury Mofizur Rahman, Rebecca Strachan, Graham Sexton, Keshav Dahal, An adaptive ensemble classifier for mining concept drifting data streams, Expert Systems with Applications, vol. 40, no. 15, pp. 5895-5906, 2013.
- Ryan Elwell, Robi Polikar, Incremental learning of concept drift in nonstationary environments, IEEE Transactions on Neural Networks, vol. 22, no. 10, pp. 1517-1531, Oct 2011.

Data stream classification methods



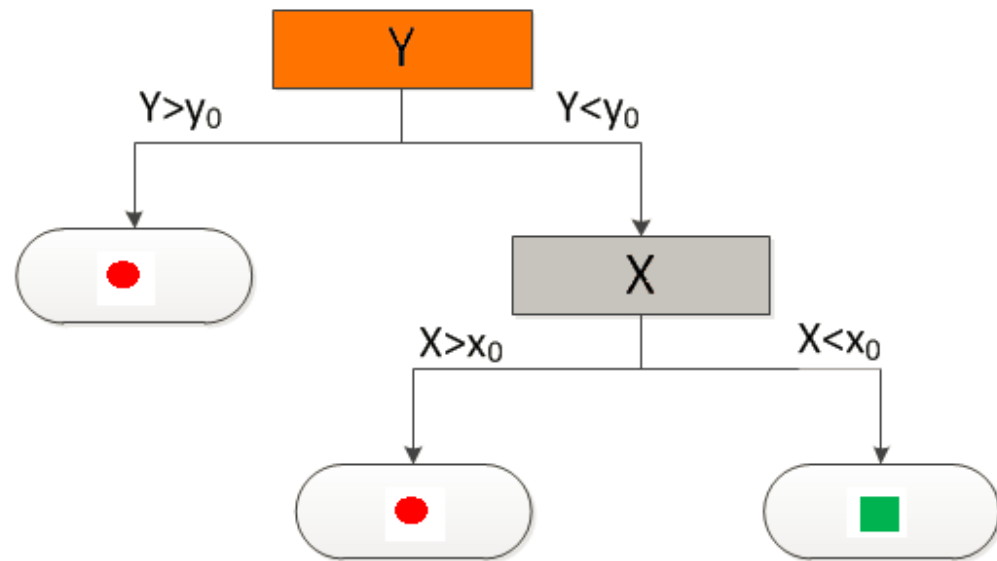
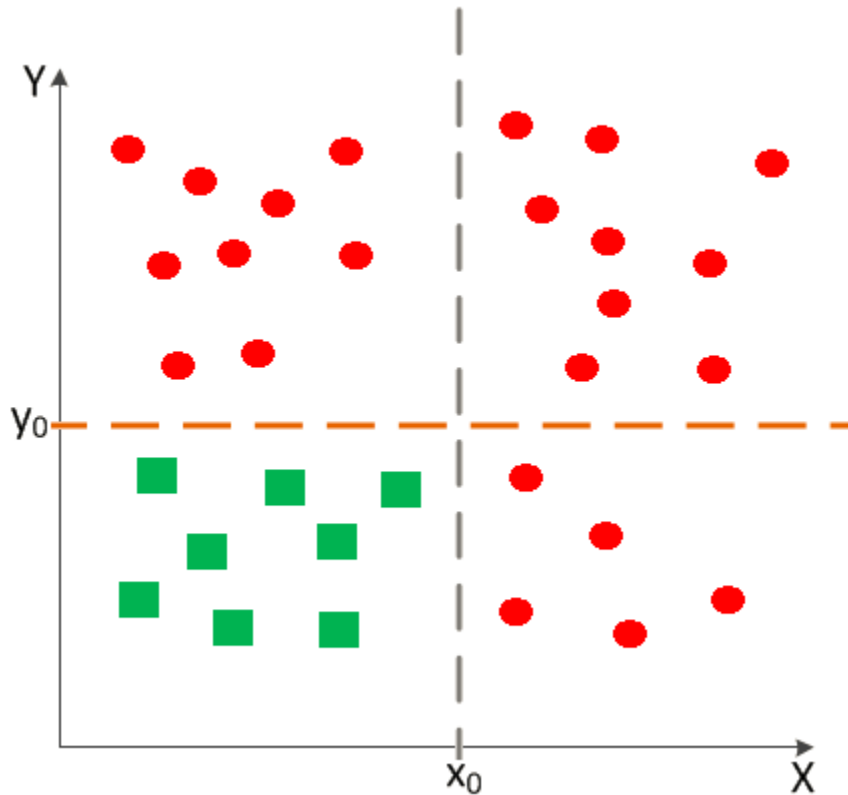
Data stream mining

- content

- Data streams – introduction to the topic
- Concept drift
- Various strategies of learning
- How to deal with concept drift?
- Data stream classification methods – short overview
- **Decision trees for data streams (including new results 2016)**
- Ensemble methods for data streams (including new results 2016)
- Probabilistic neural networks for stream data mining (including new results 2016)
- Final remarks and challenging problems
- References

Decision trees for data mining (non-stream data)

DECISION TREES



DECISION TREES

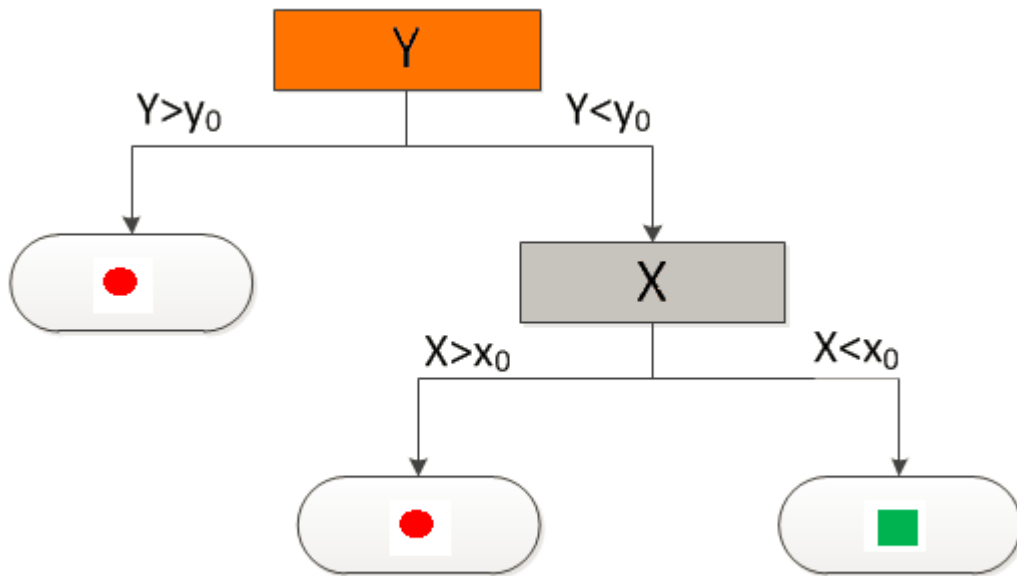
Commonly known decision tree algorithms:

❖ **ID3 algorithm** – J. R. Quinlan, “Induction of decision tree,” Machine Learning, Vol. 1, No. 1, pp. 81-106, 1986.

❖ **C4.5 algorithm** – J. R. Quinlan, „C4.5: Programs for Machine Learning,” Morgan Kaufmann Publishers, 1993

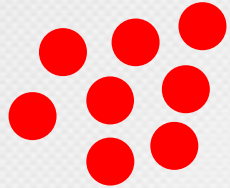
❖ **CART algorithm** – L. Breiman, J. H. Friedman, R. A. Olshen, & C. J. Stone, „Classification and regression trees.” Monterey, CA: Wadsworth & Brooks/Cole, 1984. Advanced Books & Software.

DECISION TREES

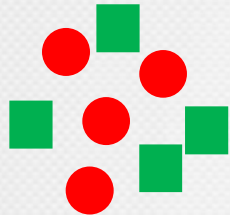


- The most critical point is the choice of a splitting attribute in each node;
- In ID3, C4.5, CART – the choice based on an impurity measure.

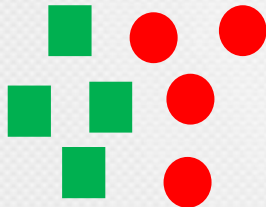
Impurity Measure



Lowest possible (0) impurity measure – maximally „pure” set



Highest possible impurity measure



Highest possible impurity measure

Impurity measure – information entropy

The entropy of S is defined as follows

$$H(S) = - \sum_{i=1}^K p_i \log(p_i) ,$$

S – training set (X_1, \dots, X_N)

K – number of classes

p_i – probability that element belongs to the i -th class
(proportion of elements in S belonging to the i -th class)

Impurity measure – information entropy

$$H(S) = - \sum_{i=1}^K p_i \log(p_i) ,$$

- If $p_i = \frac{1}{K}, i = 1, \dots, K$, then $H(S) = \log_2 K$ - maximal value
- If $p_j = 1, p_{i \neq j} = 0$, then $H(S) = 0$ - minimal value

Split measure function

For each attribute the quality of potential split is measured by the **SPLIT MEASURE FUNCTION**

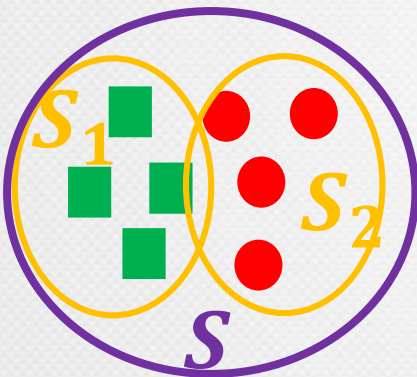
Split measure function — a reduction of an impurity measure

Split measure function

Split measure:

$$f(S) = g(S) - \frac{|S_1|}{|S|} g(S_1) - \frac{|S_2|}{|S|} g(S_2)$$

Example



S – high impurity measure $g(S)$

S_1 - low impurity measure $g(S_1)$

S_2 - low impurity measure $g(S_2)$

Split measure – information gain

$$H(S|a) = H(S) - \sum_{i=1}^{|a|} w_i H(S_i),$$

where

a – an attribute with values in set $\{a_1, \dots, a_{|a|}\}$

$|a|$ – number of different values of attribute a

w_i - fraction of elements with value a_i

(probability that elements in S take value a_i for attribute a)

S_i - subset of set S , with data elements for which the value of attribute a is equal a_i

Information entropy and information gain – an example

married	age	uses computer at work	buy computer
Yes	young	Yes	Yes
Yes	young	Yes	Yes
No	young	Yes	Yes
Yes	middle	Yes	Yes
No	middle	Yes	Yes
No	old	Yes	Yes
Yes	old	No	No
Yes	old	No	No
Yes	young	No	Yes
Yes	middle	No	No

$$H(S) = - \left(\frac{7}{10} \log_2 \left(\frac{7}{10} \right) + \frac{3}{10} \log_2 \left(\frac{3}{10} \right) \right) = 0,88129$$

Partition according to „uses computer at work”

S_{Yes}			
is married	age	uses computer at work	by computer
Yes	young	Yes	Yes
Yes	young	Yes	Yes
No	young	Yes	Yes
Yes	middle	Yes	Yes
No	middle	Yes	Yes
No	old	Yes	Yes

S_{No}			
is married	age	uses computer at work	by computer
Yes	old	No	No
Yes	old	No	No
Yes	young	No	Yes
Yes	middle	No	No

Information gain: $H(S|\text{use computer at work})$

S_{Yes}				S_{No}			
is married	age	uses computer at work	by computer	is married	age	uses computer at work	by computer
Yes	young	Yes	Yes				
Yes	young	Yes	Yes	Yes	old	No	No
No	young	Yes	Yes	Yes	old	No	No
Yes	middle	Yes	Yes	Yes	young	No	Yes
No	middle	Yes	Yes	Yes	middle	No	No
No	old	Yes	Yes				

$$P(a = Yes) = 6/10 = w_1$$

$$P(a = No) = 4/10 = w_2$$

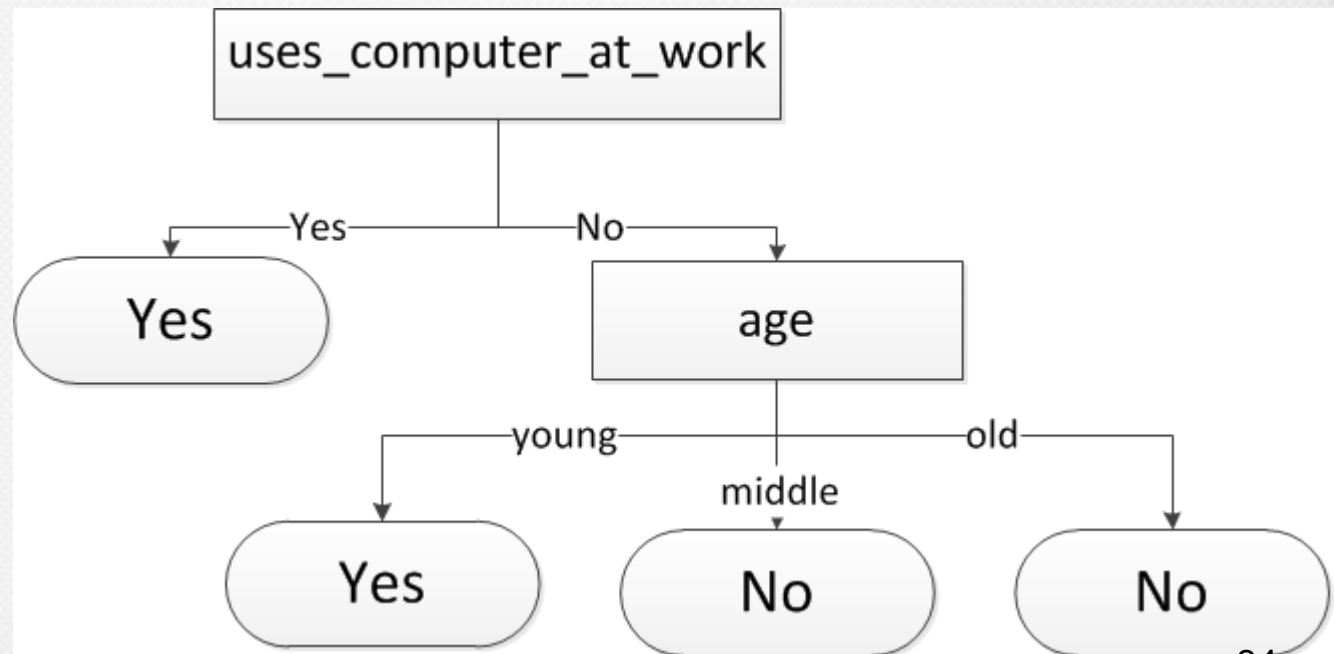
$$H(S_{Yes}) = 0$$

$$H(S_{No}) = 0,81128$$

$$\begin{aligned}
 H(S|a) &= H(S) - 6/10 * 0 - 4/10 * 0,81128 \\
 &= 0,55678
 \end{aligned}$$

Obtained decision tree

<i>attribute a</i>	<i>H(S/a)</i>
is married	0,191631
age	0,330313
uses computer at work	0,55678



Other impurity measures

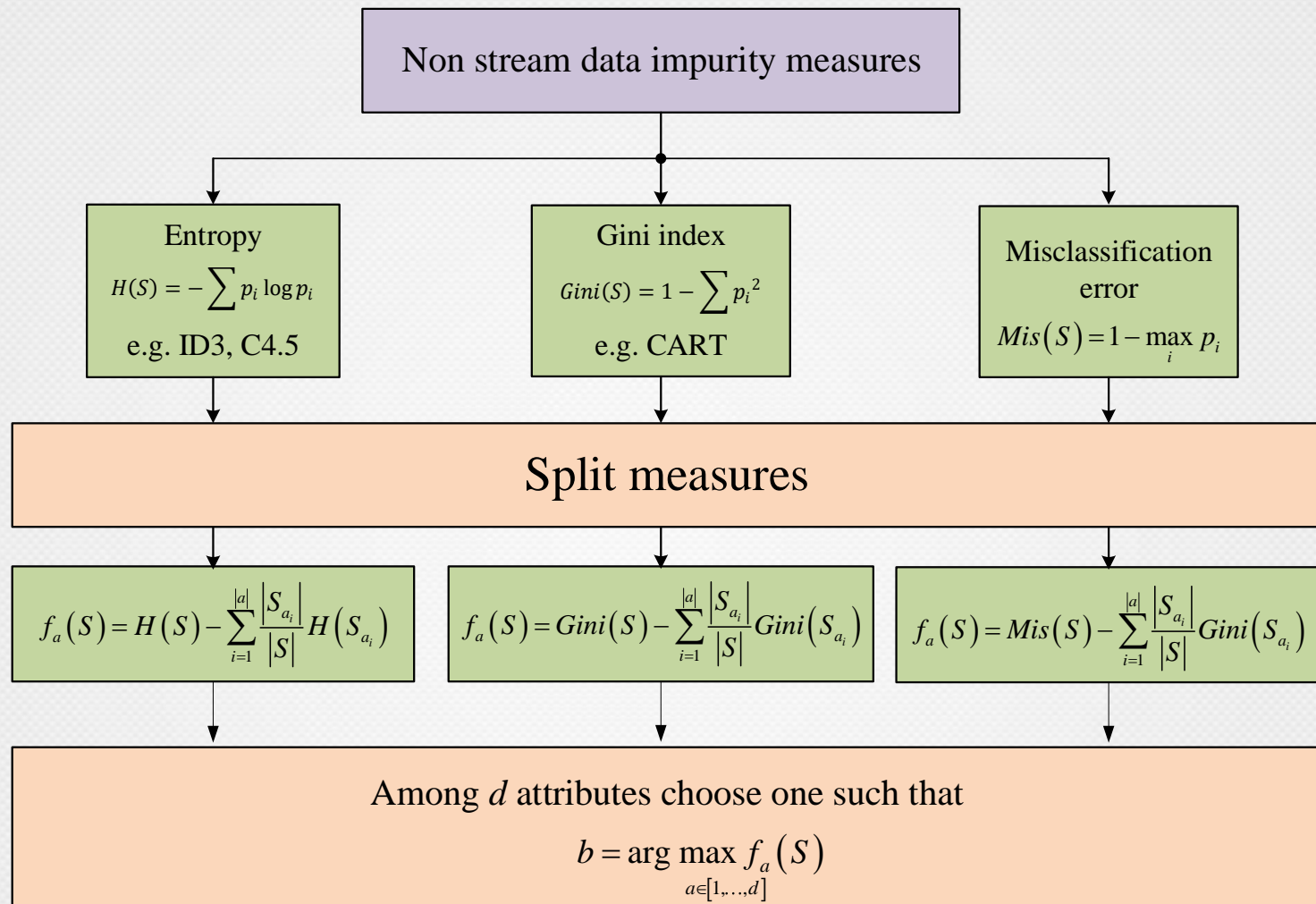
- **Gini index:**

$$\text{Gini}(\mathcal{S}) = 1 - \sum_{i=1}^K (p_i)^2 ,$$

- **Misclassification error:**

$$\text{Mis}(\mathcal{S}) = 1 - \max_{i \in \{1, \dots, K\}} p_i$$

Decision Trees for data mining (non-stream data)



Decision Trees for data streams

Decision tree for data stream

In the case of data streams data arrive continuously to the considered node

Decision tree for static data:

- 1) Which attribute to choose?

Decision tree for stream data:

- 1) Which attribute to choose?
- 2) When to make a split?

Decision tree for data stream

The aim is to design a decision tree learning system, applied to data streams, which provides an output nearly identical to that induced by a conventional learner.

Decision tree for data stream

Hoeffding trees

P. Domingos and G. Hulten, "Mining high-speed data streams", Proc. 6th ACM SIGKDD Internat. Conf. on Knowledge Discovery and Data Mining, pp. 71-80, 2000.

G. Hulten, L. Spencer, P. Domingos, Mining time-changing data streams, In Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '01), ACM, New York, NY, USA, pp. 97-106, 2001

Hoeffding Tree Algorithm - 2000

Main ideas applied in the Hoeffding tree algorithm:

- 1) **Sufficient statistics**
- 2) **Splitting criterion**

P. Domingos and G. Hulten, "Mining high-speed data streams", Proc. 6th ACM SIGKDD Internat. Conf. on Knowledge Discovery and Data Mining, pp. 71-80, 2000.

Hoeffding tree algorithm: sufficient statistics

In each node information collected in a form of „sufficient statistics”:

N_{ijk} - numer of data elements from the **k-th class**
which take the **j-th value** of the **i-th attribute**

Required memory per node: $M = K \sum_{i=1}^d v_i$.

For two-class binary problem: $M = 4d$.

Hoeffding tree algorithm: splitting criterion

- 1) Compute the split measure function $f_i(S)$ for each attribute $i = 1, \dots, d$ based on the currently collected data sample $S = X_1, \dots, X_N$
- 2) Find two attributes with the highest values of f :
 $f_{a_{max1}}(S), \quad f_{a_{max2}}(S)$

Hoeffding tree algorithm: splitting criterion

3) **If**

$$f_{a_{max1}}(S) - f_{a_{max2}}(S) > \varepsilon = \varepsilon(N, \delta),$$

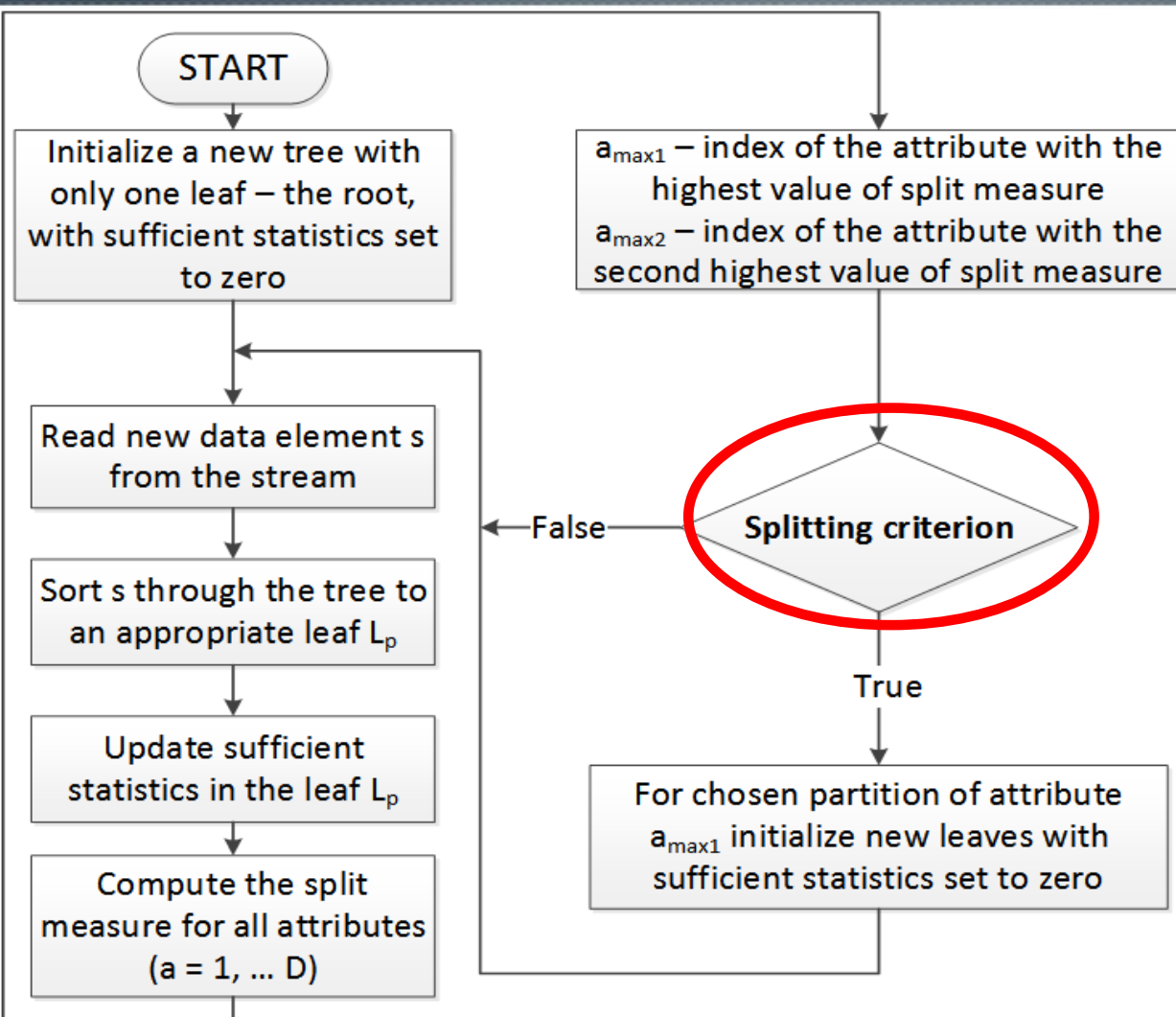
then with probability $1 - \delta$

$$E[f_{a_{max1}}(S)] > E[f_{a_{max2}}(S)]$$

and choose a_{max1} as a splitting attribute

Else, wait for more data elements in this node

Online Decision Tree – general algorithm



Splitting criterion:

$$f_{a_{\max 1}} - f_{a_{\max 2}} > \epsilon$$

Hoeffding tree algorithm: splitting criterion

Splitting criterion:

$$f_{a_{max1}}(S) - f_{a_{max2}}(S) > \varepsilon = \varepsilon(N, \delta),$$

Challenge: How to find formula for $\varepsilon(N, \delta)$???

Hoeffding's inequality - 1963

The commonly known algorithm called 'Hoeffdings Tree' was introduced by P. Domingos and G. Hulten in [1]. The main mathematical tool used in this algorithm was the Hoeffding's inequality [2]:

Theorem: *If X_1, X_2, \dots, X_N are independent random variables and $a_i \leq X_i \leq b_i$ ($i = 1, 2, \dots, N$), then for $\epsilon > 0$*

$$P\{\bar{X} - E[\bar{X}] \geq \epsilon\} \leq e^{-2N^2\epsilon^2 / \sum_{i=1}^N (b_i - a_i)^2}$$

where

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i \text{ and } E[\bar{X}] \text{ is expected value of } \bar{X}.$$

[1] P. Domingos and G. Hulten, "Mining high-speed data streams", Proc. 6th ACM SIGKDD Internat. Conf. on Knowledge Discovery and Data Mining, pp. 71-80, 2000.

[2] W. Hoeffding, "Probability inequalities for sums of bounded random variables", Journal of the American Statistical Association, vol. 58, issue 301, pp. 13-30, March 1963.

Hoeffding's inequality - 1963

Assuming that $P\{\bar{X} - E[\bar{X}] \geq \epsilon\} \leq \delta$
and $b_i - a_i = R, i = 1, \dots, N$,
the Hoeffding's inequality is equivalent to:

$$P\left\{\bar{X} - E[\bar{X}] \leq \sqrt{\frac{R^2 \ln 1/\delta}{2N}}\right\} \geq 1 - \delta$$

[1] W. Hoeffding, "Probability inequalities for sums of bounded random variables", Journal of the American Statistical Association, vol. 58, issue 301, pp. 13-30, March 1963.

[2] P. Domingos and G. Hulten, "Mining high-speed data streams", Proc. 6th ACM SIGKDD Internat. Conf. on Knowledge Discovery and Data Mining, pp. 71-80, 2000.

Hoeffding's trees

P. Domingos and G. Hulten claimed that owing to the Hoeffding's inequality the form of ε is given by:

$$\varepsilon = \varepsilon(N, \delta) = \sqrt{\frac{R^2 \ln(1/\delta)}{2N}},$$

where R is a range of values of the applied split measure, e.g. $R = \log_2 K$ for information gain

[1] P. Domingos and G. Hulten, "Mining high-speed data streams", Proc. 6th ACM SIGKDD Internat. Conf. on Knowledge Discovery and Data Mining, pp. 71-80, 2000.

Algorithms based on the Hoeffding's bound - 2000

Very Fast Decision Tree (VFDT)

The VFDT (Very Fast Decision Tree) algorithm makes several modifications to the Hoeffding tree algorithm to improve both speed and memory utilization. The modifications include:

- breaking near-ties during attribute selection more aggressively,
- computing function f_a after a number of training examples,
- deactivating the least promising leaves whenever memory running low,
- dropping poor splitting attributes,
- improving the initialization method.

[1] P. Domingos and G. Hulten, "Mining high-speed data streams", Proc. 6th ACM SIGKDD Internat. Conf. on Knowledge Discovery and Data Mining, pp. 71-80, 2000.

Algorithms based on the Hoeffding's bound - 2001

Concept-adapting Very Fast Decision Tree (CVFDT) for handling concept drift

CVFDT uses a sliding window approach, the main features are the following:

- it does not construct a new model from scratch each time,
- it updates statistics at the nodes by incrementing the counts associated with new examples and decrementing the counts associated with old ones,

G. Hulten, L. Spencer, P. Domingos, Mining time-changing data streams, In Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '01), ACM, New York, NY, USA, pp. 97-106, 2001

Algorithms based on the Hoeffding's bound - 2001

Concept-adapting Very Fast Decision Tree (CVFDT) for handling concept drift

- If there is a concept drift, some nodes may no longer pass the Hoeffding bound. When this happens, an alternate subtree will be grown, with the new best splitting attribute at the root. As new examples stream in, the alternate subtree will continue to develop, without yet being used for classification. Once the alternate subtree becomes more accurate than the existing one, the old subtree is replaced.

G. Hulten, L. Spencer, P. Domingos, Mining time-changing data streams, In Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '01), ACM, New York, NY, USA, pp. 97-106, 2001

The Hoeffding's bound

The Hoeffding's bound is a wrong tool to solve the problem of choosing the best attribute to make a split in the node!!!

The Hoeffding's bound

Hoeffding's inequality is applicable only for sums (or arithmetic averages) of random variables.

Nonlinear impurity measures, like information entropy

$$H(S) = - \sum_{i=1}^K p_i(\mathbf{S}) \log_2 p_i(\mathbf{S})$$

or Gini index

$$Gini(S) = 1 - \sum_{i=1}^K (p_i(\mathbf{S}))^2,$$

can not be presented as a sum of elements.

The Hoeffding's bound

The idea presented in [1] violates the assumptions of the Hoeffding's theorem (see [2]) and the concept of Hoeffding Trees has no theoretical justification.

[1] P. Domingos and G. Hulten, "Mining high-speed data streams", Proc. 6th ACM SIGKDD Internat. Conf. on Knowledge Discovery and Data Mining, pp. 71-80, 2000.

[2] W. Hoeffding, "Probability inequalities for sums of bounded random variables", Journal of the American Statistical Association, vol. 58, issue 301, pp. 13-30, March 1963.

Challenge in Stream Data Mining

Find an appropriate, mathematically justified, form of the bound $\varepsilon = \varepsilon(N, \delta)$ in the general splitting criterion:

$$f_{a_{max1}}(S) - f_{a_{max2}}(S) > \varepsilon = \varepsilon(N, \delta)$$

What will be shown on the next slides???

- a) We will study split measures based on the following impurity measures:
- information entropy,
 - Gini index,
 - misclassification error.
- b) We will propose two different techniques to solve the problem:
- the McDiarmid's inequality,
 - the Gaussian approximation.

What will be shown on the next slides???

- c) We will propose the decision trees algorithms such that if

$$f_{a_{max_1}} - f_{a_{max_2}} > \varepsilon(N, \delta)$$

then

$$E \left[f_{a_{max_1}} \right] > E \left[f_{a_{max_2}} \right]$$

with probability $(1 - \delta)^{d-1}$, where d is the number of attributes.

WHAT IS THE VALUE OF ε ???

Papers concerning the issue

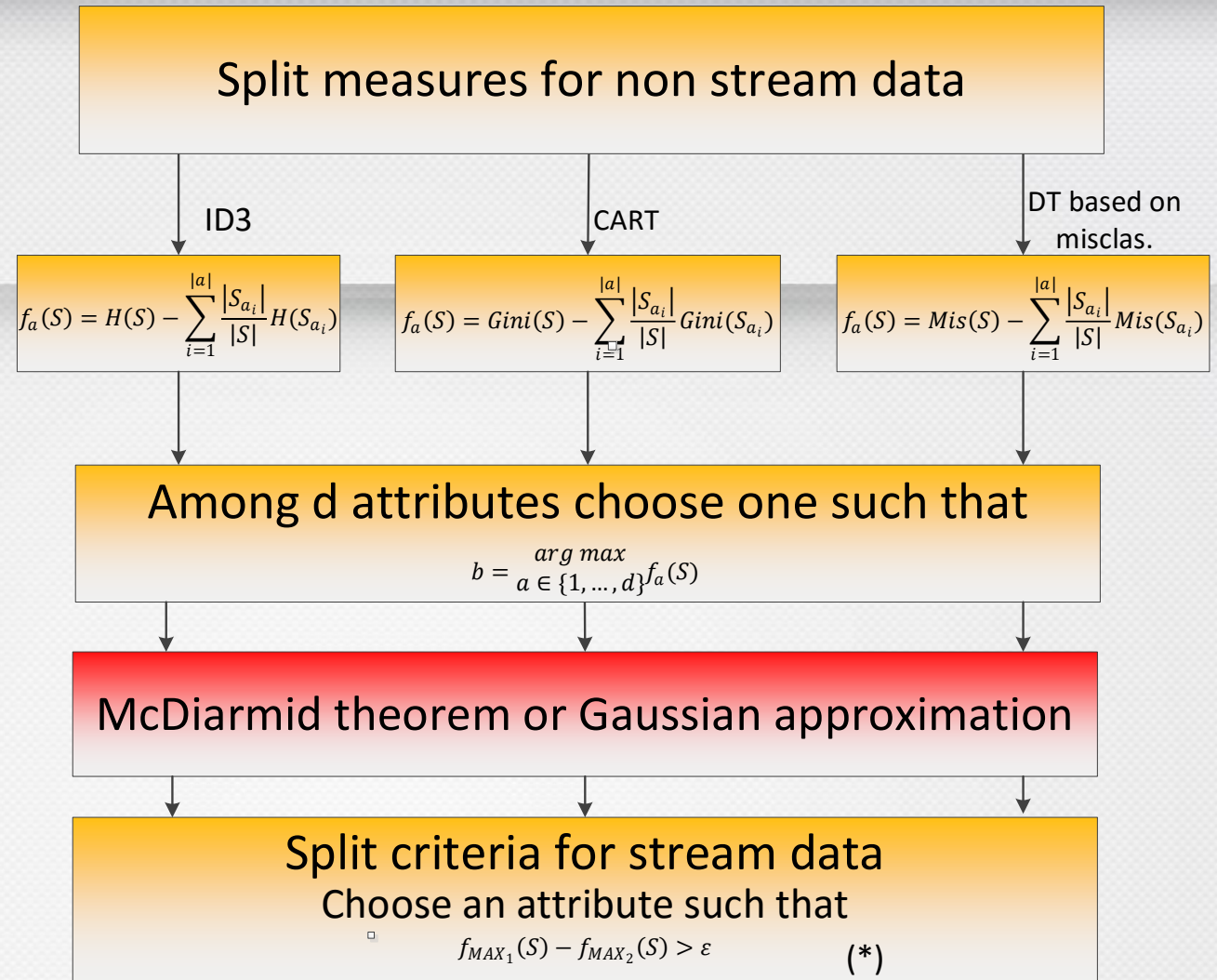
Leszek Rutkowski, Lena Pietruczuk, Piotr Duda, Maciej Jaworski, *Decision trees for mining data streams based on the McDiarmid's bound*, IEEE Transactions on Knowledge and Data Engineering, vol. 25, no. 6, pp. 1272–1279, 2013.

Leszek Rutkowski, Maciej Jaworski, Lena Pietruczuk, Piotr Duda, *Decision trees for mining data streams based on the Gaussian approximation*, IEEE Transactions on Knowledge and Data Engineering, vol. 26, no. 1, pp. 108-119, Jan. 2014.

Leszek Rutkowski, Maciej Jaworski, Lena Pietruczuk, Piotr Duda, *The CART decision tree for mining data streams*, Information Sciences, vol. 266, pp. 1 – 15, 2014.

Leszek Rutkowski, Maciej Jaworski, Lena Pietruczuk, Piotr Duda, *A new method for data stream mining based on the misclassification error*, IEEE Transaction on Neural Networks and Learning Systems, vol. 26, PP 1048-1059, no. 5, 2015.

How to deal with stream data?



Under condition (*)

$$E[f_{a_{MAX_1}}] > E[f_{a_{MAX_2}}]$$

with probability $(1 - \delta)^{d-1}$.

HOW TO DETERMINE THE VALUE OF ϵ ???

McDiarmid's inequality

Let $S = \{X_1, \dots, X_N\}$ be the set of i.i.d. random variables, $X_i \in U_i$

Suppose that the (measurable) function

$\tilde{f}: \prod U_i \rightarrow \mathbb{R}$ satisfies

$$\sup_{X_1, \dots, X_N, \widehat{X}_i} |\tilde{f}(X_1, \dots, X_i, \dots, X_N) - \tilde{f}(X_1, \dots, \widehat{X}_i, \dots, X_N)| \leq c_i,$$

for some constants $c_i, i = 1, \dots, N$. Then

$$\Pr(\tilde{f}(S) - \mathbb{E}[\tilde{f}(S)] \geq \varepsilon) \leq \exp\left(\frac{-2\varepsilon^2}{\sum_{i=1}^N (c_i)^2}\right).$$

Application of the McDiarmid's inequality to stream data mining

The main result of our research is the following theorem stating that if the difference between the information gain estimates obtained for two attributes is greater than a specific value $\varepsilon(N, \delta)$, then with a fixed probability $1 - \delta$ there is, roughly speaking, a statistical difference between the expected values of information gain.

McDiarmid's inequality - information gain

Theorem 1: Let $S = \{X_1, \dots, X_N\}$ be the set of independent random variables, with each of them taking values in the set $A_1 \times \dots \times A_d \times Y$. Then, for any fixed δ and any pair of attributes a and b , where $H(S|a) - H(S|b) > 0$, if

$$\varepsilon = C_{Gain}(K, N) \sqrt{\frac{\ln 1/\delta}{2N}}$$

where

$$C_{Gain}(K, N) = 6(K \log_2 e N + \log_2 2N) + 2 \log_2 K$$

then

$$E[H(S|a)] > E[H(S|b)], \quad \text{with prob. } 1 - \delta.$$

L. Rutkowski, L. Pietruczuk, P. Duda, M. Jaworski, *Decision trees for mining data streams based on the McDiarmid's bound*, IEEE Transactions on Knowledge and Data Engineering, vol.25, no.6, pp.1272-1279, June 2013

McDiarmid's inequality - Gini gain

Theorem 2: Let $S = \{X_1, \dots, X_N\}$ be the set of independent random variables, with each of them taking values in the set $A_1 \times \dots \times A_d \times Y$. Then, for any fixed δ and any pair of attributes a and b , where $Gini(S|a) - Gini(S|b) > 0$, if

$$\varepsilon = \sqrt{\frac{8 \ln 1/\delta}{2N}}$$

then

$$E[Gini(S|a)] > E[Gini(S|b)], \text{ with prob. } 1 - \delta.$$

L. Rutkowski, L. Pietruczuk, P. Duda, M. Jaworski, *Decision trees for mining data streams based on the McDiarmid's bound*, IEEE Transactions on Knowledge and Data Engineering, vol.25, no.6, pp.1272-1279, June 2013

McDiarmid's inequality - Gini gain

Conclusion: If the number of data elements N satisfies the condition

$$N > 64 \frac{\ln 1/\delta}{2[Gini(S|a) - Gini(S|b)]}$$

then the number of data elements is sufficient enough to say that attribute a is „better” (with probability $1 - \delta$ to make a split than attribute b).

Misclassification error

$$Mis(S) = 1 - \max_{i=1\dots K} \{p_i\}$$

or equivalently

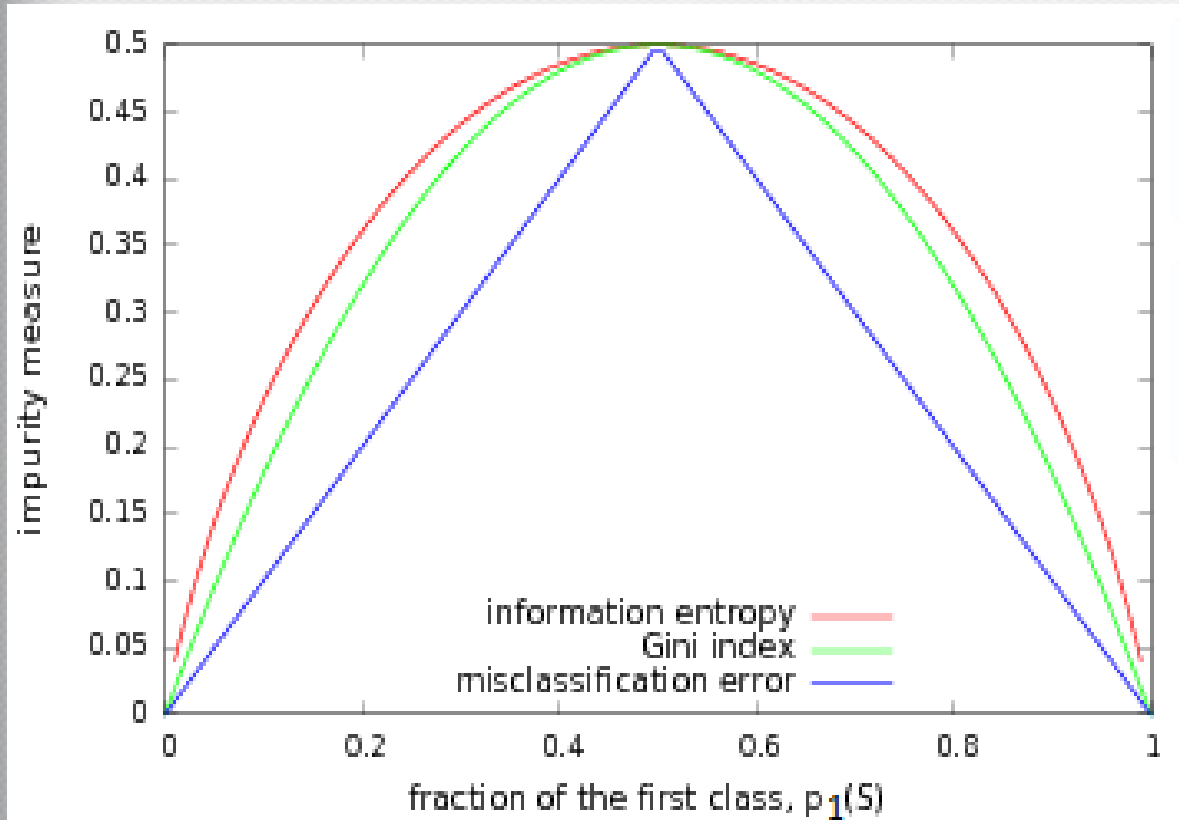
$$Mis(S) = 1 - \frac{\max_{i=1\dots K} \{N^i\}}{N}$$

where

p_i – probability that element belongs to the i -th class
(proportion of elements in S belonging to the i -th class)

N^i - the number of data elements in S from the i -th class

Misclassification error



$$H(S) = - \sum_i^K p_i \log_2 p_i$$

$$Gini(S) = 1 - \sum_i^K (p_i)^2$$

$$Mis(S) = 1 - \max_{i=1 \dots K} \{p_i\}$$

Comparison for two-class problem

Split measure function based on the misclassification error

$$\mathbf{Mis}(S|a) = \mathbf{Mis}(S) - \sum_{i=1}^{|a|} \frac{|S_i|}{|S|} \mathbf{Mis}(S_i)$$

Analogously to the information gain:

$$\mathbf{H}(S|a) = \mathbf{H}(S) - \sum_{i=1}^{|a|} \frac{|S_i|}{|S|} \mathbf{H}(S_i),$$

And Gini gain:

$$\mathbf{Gini}(S|a) = \mathbf{Gini}(S) - \sum_{i=1}^{|a|} \frac{|S_i|}{|S|} \mathbf{Gini}(S_i).$$

Splitting criterion for misclassification-based split measure function – Gaussian approximation

Theorem 3: Let us consider two attributes a and b , for which the values of split measure function based on the misclassification error were calculated for set S . If the condition $Mis(S|a) - Mis(S|b) > \varepsilon$ is satisfied, where

$$\varepsilon = z_{(1-\delta)} \frac{1}{\sqrt{2N}},$$

$z_{(1-\delta)}$ is the $(1 - \delta)$ –th quantile of the standard normal distribution $\mathcal{N}(0,1)$, then $E[Mis(S|a)]$ is greater than $E[Mis(S|b)]$ with probability $1 - \delta$.

L. Rutkowski, M. Jaworski, L. Pietruczuk, P. Duda, A new method for data stream mining based on the misclassification error, IEEE Transaction on Neural Networks and Learning Systems, vol. 26, pp. 1048-1059, no. 5, 2015

Splitting criterion for misclassification-based split measure function – Hoeffding's inequality

Theorem 4: Let us consider two attributes a and b , for which we the values of split measure function based on the misclassification error were calculated for set S . If the condition $Mis(S|a) - Mis(S|b) > \varepsilon$ is satisfied, where

$$\varepsilon = \sqrt{\frac{2 \ln 1/\delta}{N}},$$

then $E[Mis(S|a)]$ is greater than $E[Mis(S|b)]$ with probability $1 - \delta$.

Remark. Theorem 4 is based on the Hoeffding's inequality!!!

Comparison of the newest results

Impurity measure Method	Information entropy	Gini index	Misclassification error
Hoeffding bound	$\varepsilon = \sqrt{\frac{R^2 \ln 1/\delta}{2n}}$ <p>Incorrectly obtained Domingos and Hulten, ACM 2000</p>	$\varepsilon = \sqrt{\frac{R^2 \ln 1/\delta}{2n}}$ <p>Incorrectly obtained Domingos and Hulten, ACM 2000</p>	$\varepsilon = \sqrt{\frac{2 \ln 1/\delta}{n}}$
McDiarmid bound	$\varepsilon = C_{\text{Gain}}(K, n) \sqrt{\frac{\ln 1/\delta}{2n}}$ <p>$C_{\text{Gain}}(K, n) = 6(K \log_2 en + \log_2 2n) + 2 \log_2 K$ Rutkowski et. al., IEEE Trans. on Knowledge and Data Engineering, 2013</p>	$\varepsilon = 8 \sqrt{\frac{\ln 1/\delta}{2n}}$ <p>Rutkowski et. al., IEEE Trans. on Knowledge and Data Engineering, 2013</p>	$\varepsilon = \sqrt{\frac{2 \ln 1/\delta}{n}}$
Gaussian approximation			$\varepsilon = z_{(1-\delta)} \sqrt{\frac{1}{2n}}$ <p>Rutkowski et. al., IEEE Trans. on Neural Networks and Learning Systems, 2015</p>

If $f_a(S) - f_b(S) > \varepsilon$,

then with probability $1 - \delta$ $E[f_a(S)] > E[f_b(S)]$

Hybrid Splitting Criteria

,Single' Splitting Criteria

New result
2016

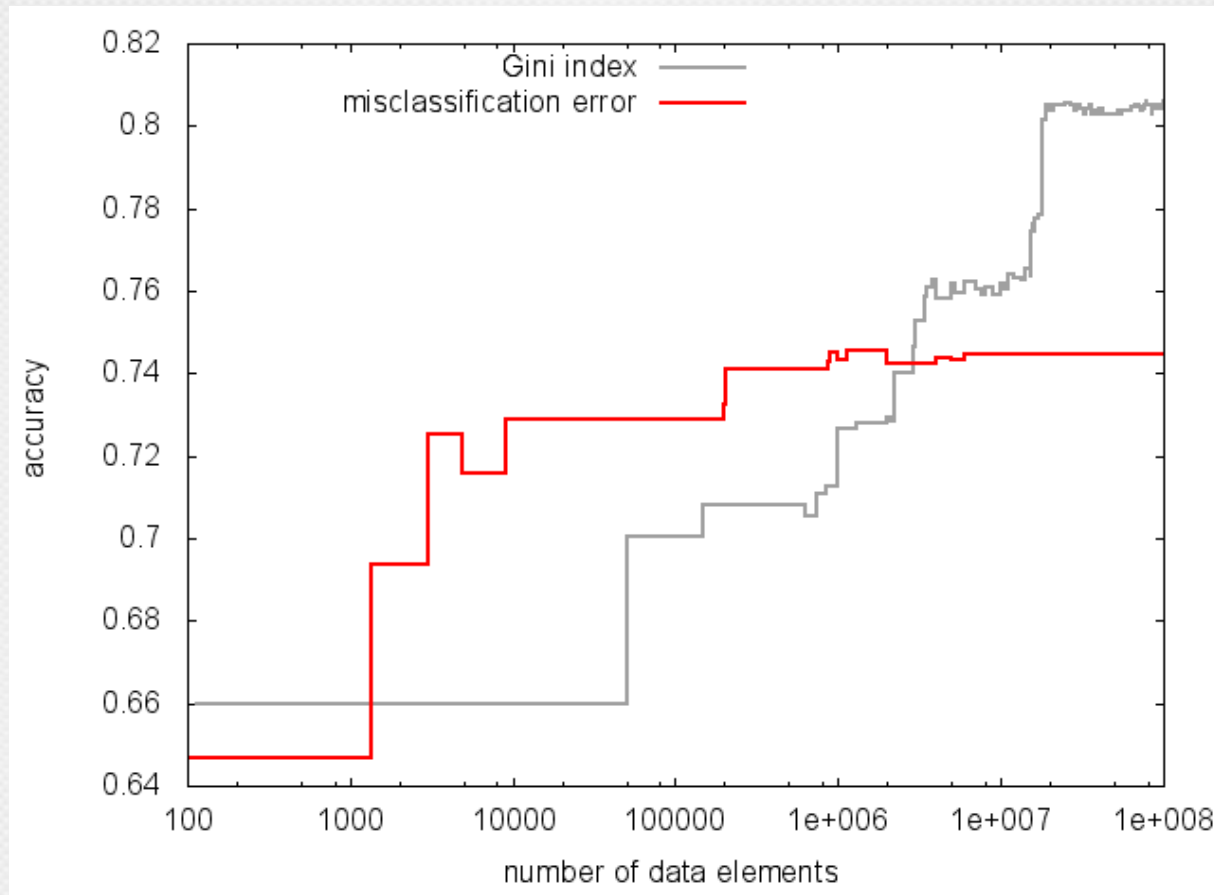
- Splitting criterion for Gini index:

$$f_{a_{max}}(\mathbf{Z}) - f_{a_{max2}}(\mathbf{Z}) > \sqrt{\frac{8 \ln\left(\frac{1}{\delta}\right)}{n(\mathbf{Z})}},$$

- Splitting criterion for misclassification error:

$$f_{a_{max}}(\mathbf{Z}) - f_{a_{max2}}(\mathbf{Z}) > z_{(1-\delta)} \sqrt{\frac{1}{2n(\mathbf{Z})}},$$

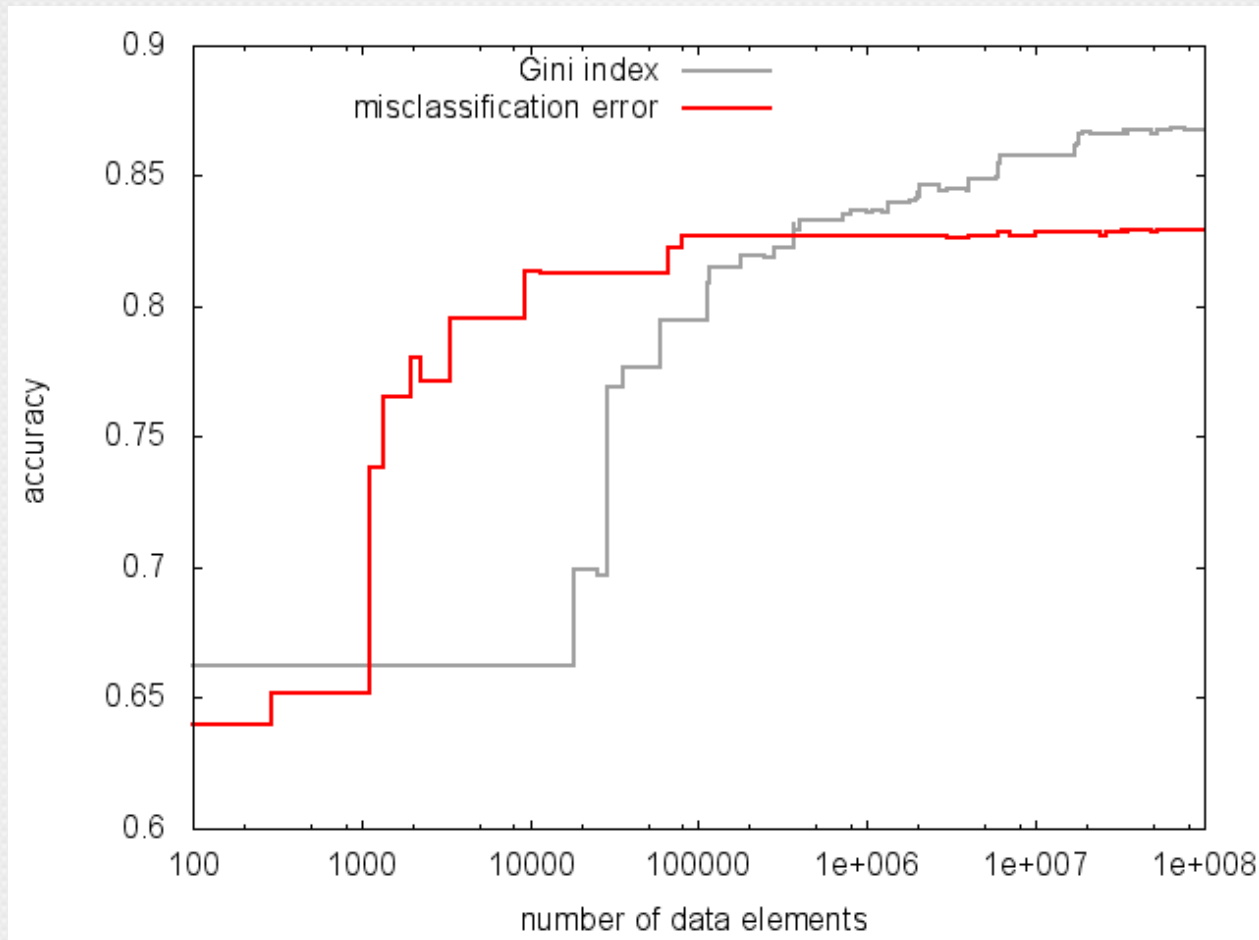
Accuracy vs number of data elements: dataset no. 1



Experimental Results

New result
2016

Accuracy vs number of data elements: dataset no. 2



Hybrid Splitting Criterion

New result
2016

$f^G_i(Z)$ - split measure for Gini index

$f^M_i(Z)$ - split measure for misclassification error

Hybrid criterion:

If:

Part 1

$$f^G_{a_{G,max}}(Z) - f^G_{a_{G,max2}}(Z) > \sqrt{\frac{8 \ln\left(\frac{1}{\delta}\right)}{n(Z)}},$$

Choose the attribute with index $a_{G,max}$ to split the node.

Else, if:

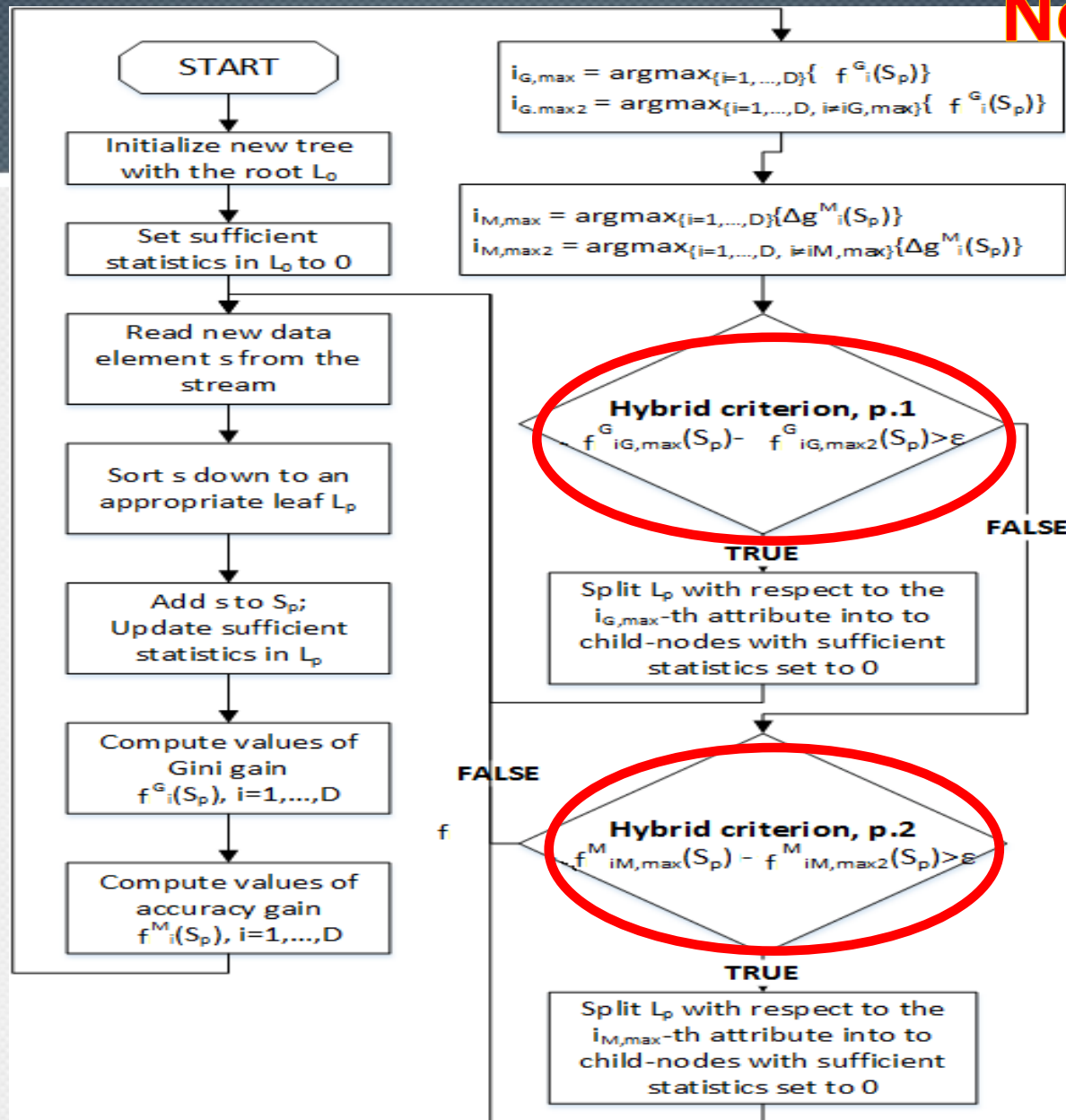
Part 2

$$f^M_{a_{M,max}}(Z) - f^M_{a_{M,max2}}(Z) > z_{(1-\delta)} \sqrt{\frac{1}{2n(Z)}},$$

Choose the attribute with index $a_{M,max}$ to split the node.

Online Decision Tree with hybrid splitting criterion

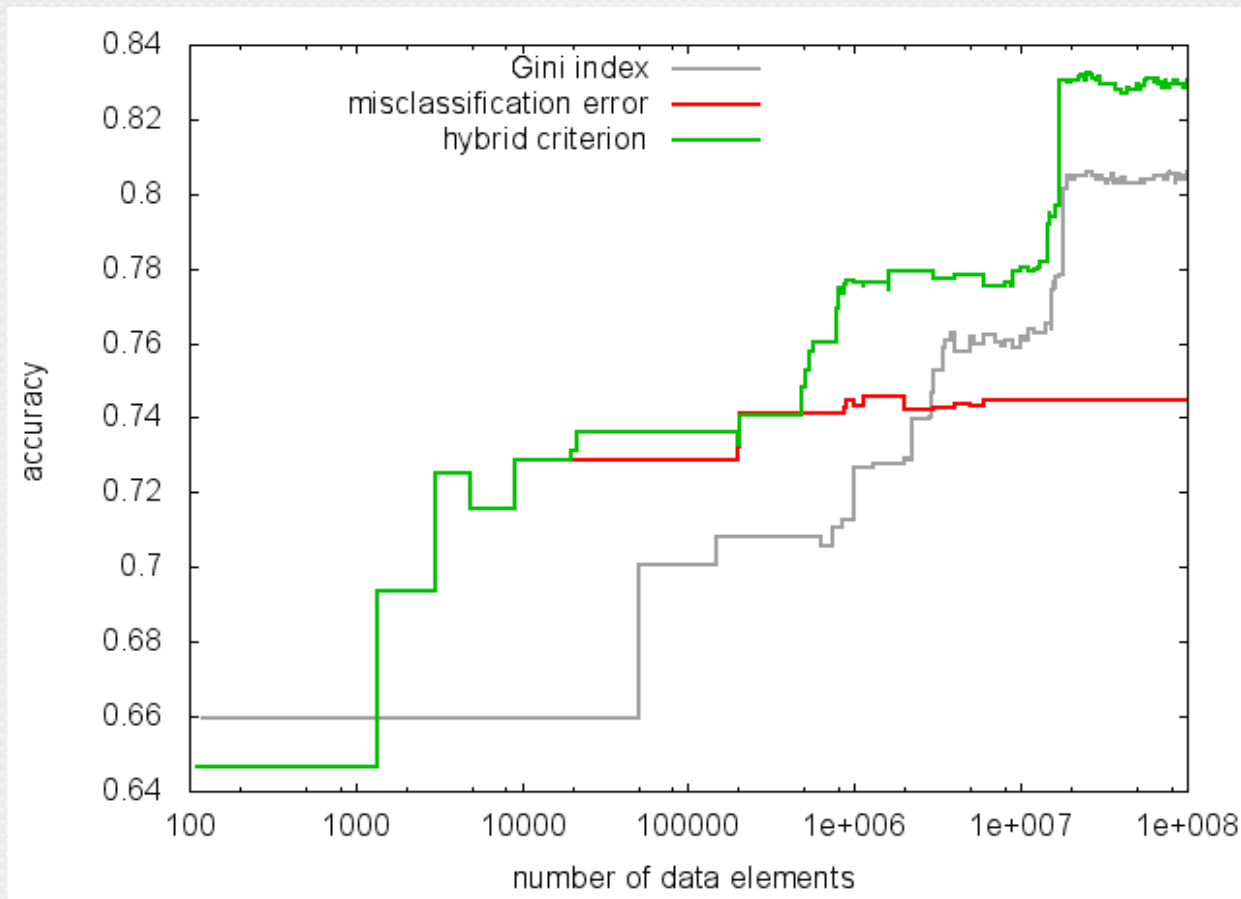
New result 2016



Experimental Results

New result
2016

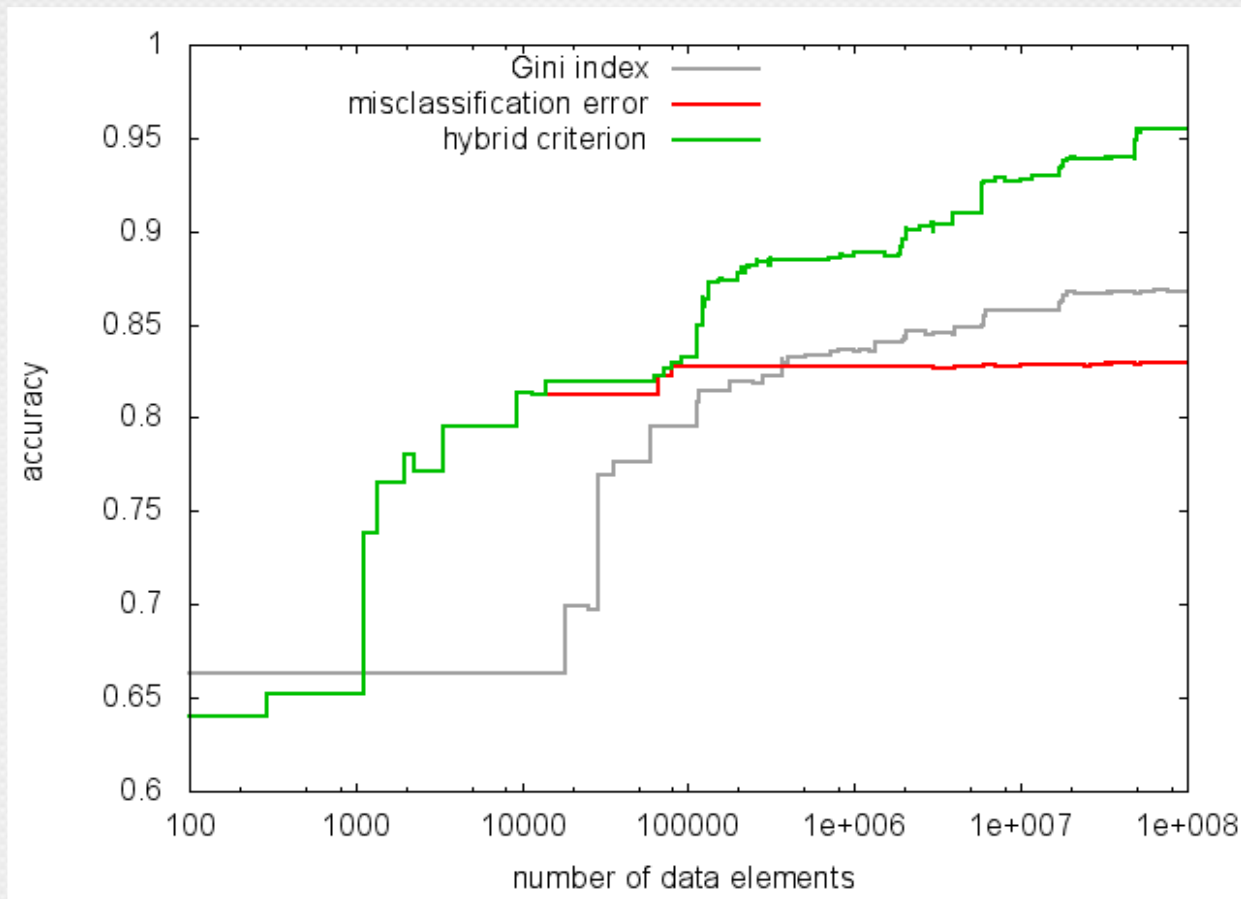
Accuracy vs number of data elements: dataset no. 1



Experimental Results

New result
2016

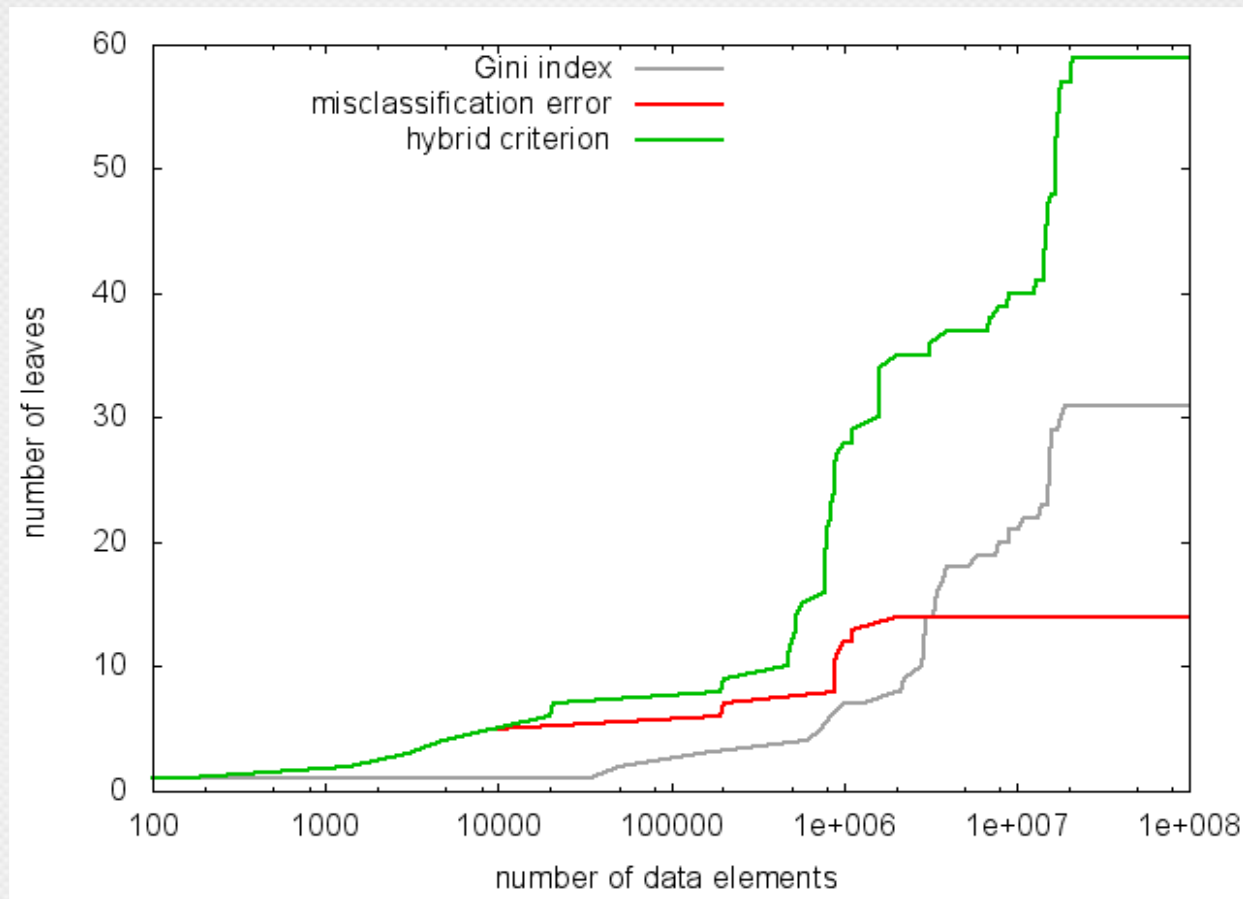
Accuracy vs number of data elements: dataset no. 2



Experimental Results

New result
2016

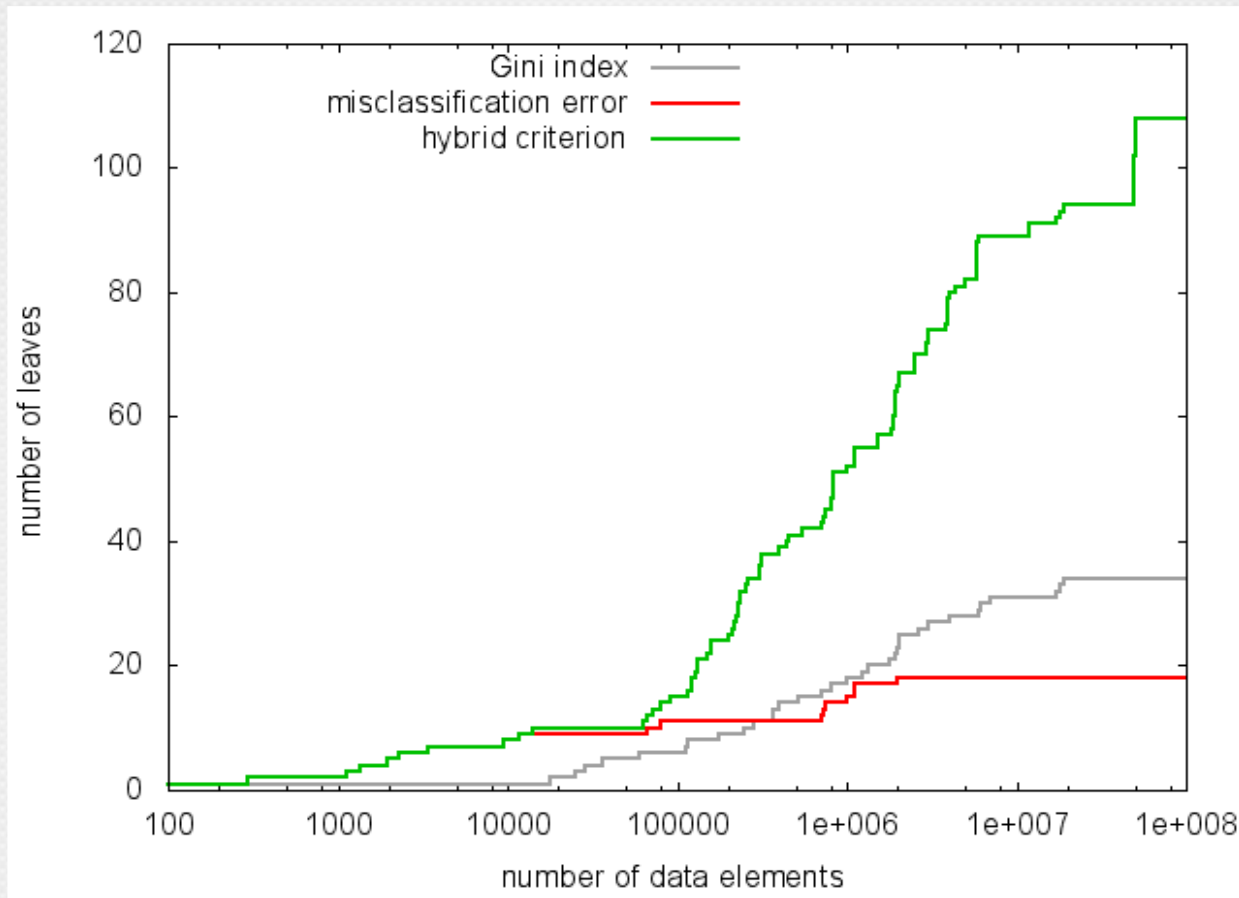
Number of leaves vs number of data elements: dataset no. 1



Experimental Results

New result
2016

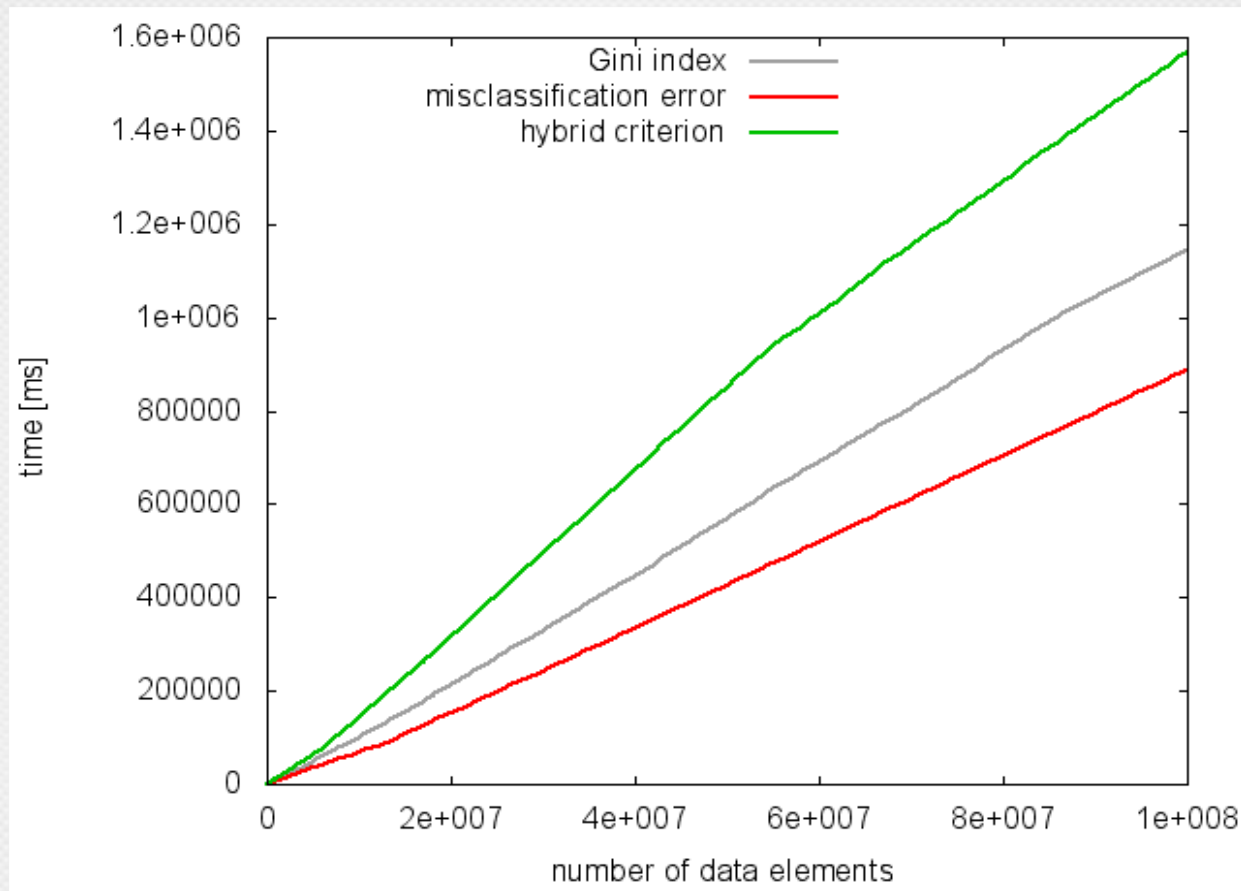
Number of leaves vs number of data elements: dataset no. 2



Experimental Results

New result
2016

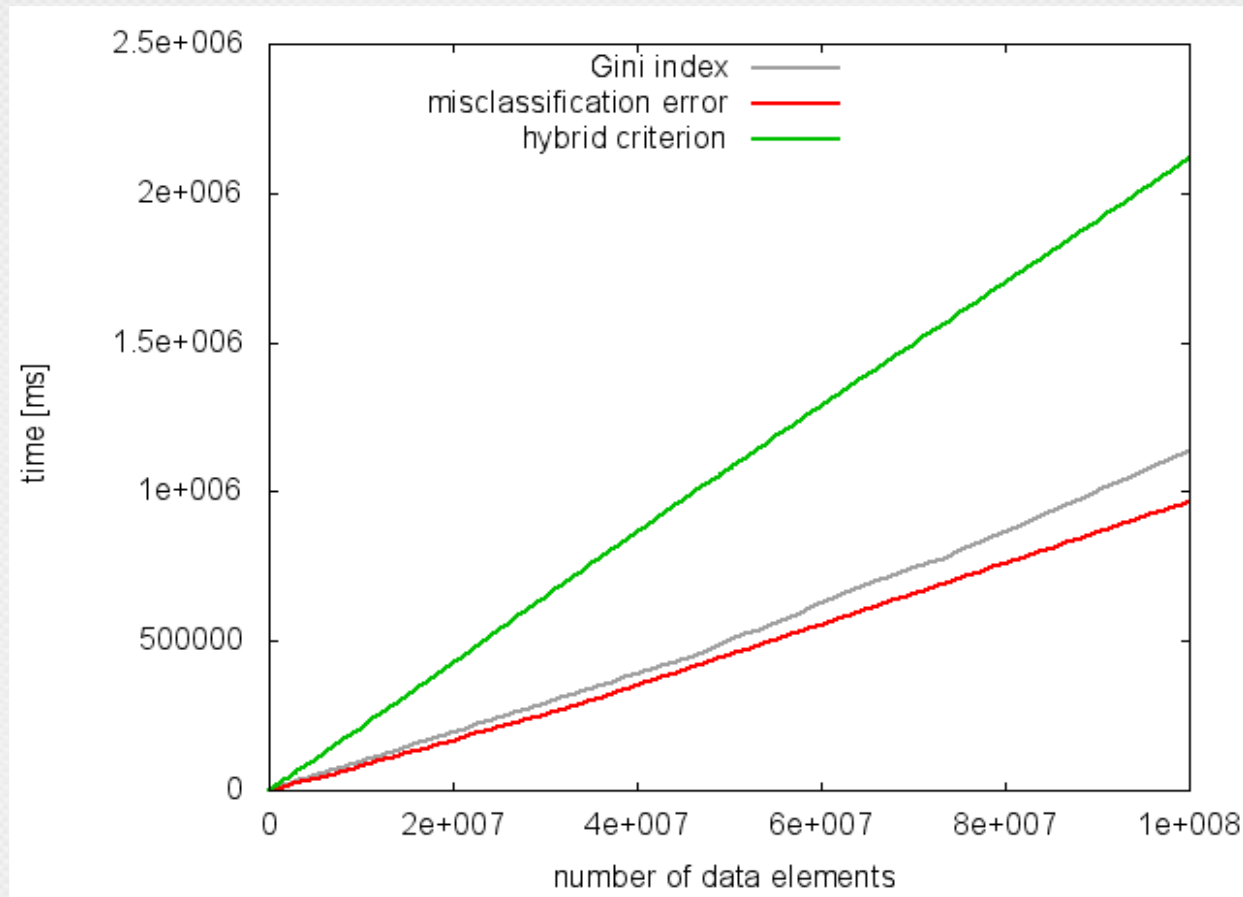
Processing time vs number of data elements: dataset no. 2



Experimental Results

New result
2016

Processing time vs number of data elements: dataset no. 2



New result

More ,single' splitting criteria 2016

- Misclassification error + Hoeffding's inequality:

$$f_{a_{\max}}(Z) - f_{a_{\max 2}}(Z) > \sqrt{\frac{2 \ln\left(\frac{1}{\delta}\right)}{n(Z)}}$$

- Gini index + McDiarmid's bound + bias ([1])

$$f_{a_{\max}}(Z) - f_{a_{\max 2}}(Z) > \sqrt{\frac{8 \ln\left(\frac{1}{\delta}\right)}{n(Z)}} + \frac{8}{\sqrt{n(Z)}}$$

[1] R. De Rosa and N. Cesa-Bianchi Splitting with confidence in decision trees with application to stream data, *Proceedings of the International Joint Conference on Neural Networks*, 2015.

Online decision trees algorithms notations

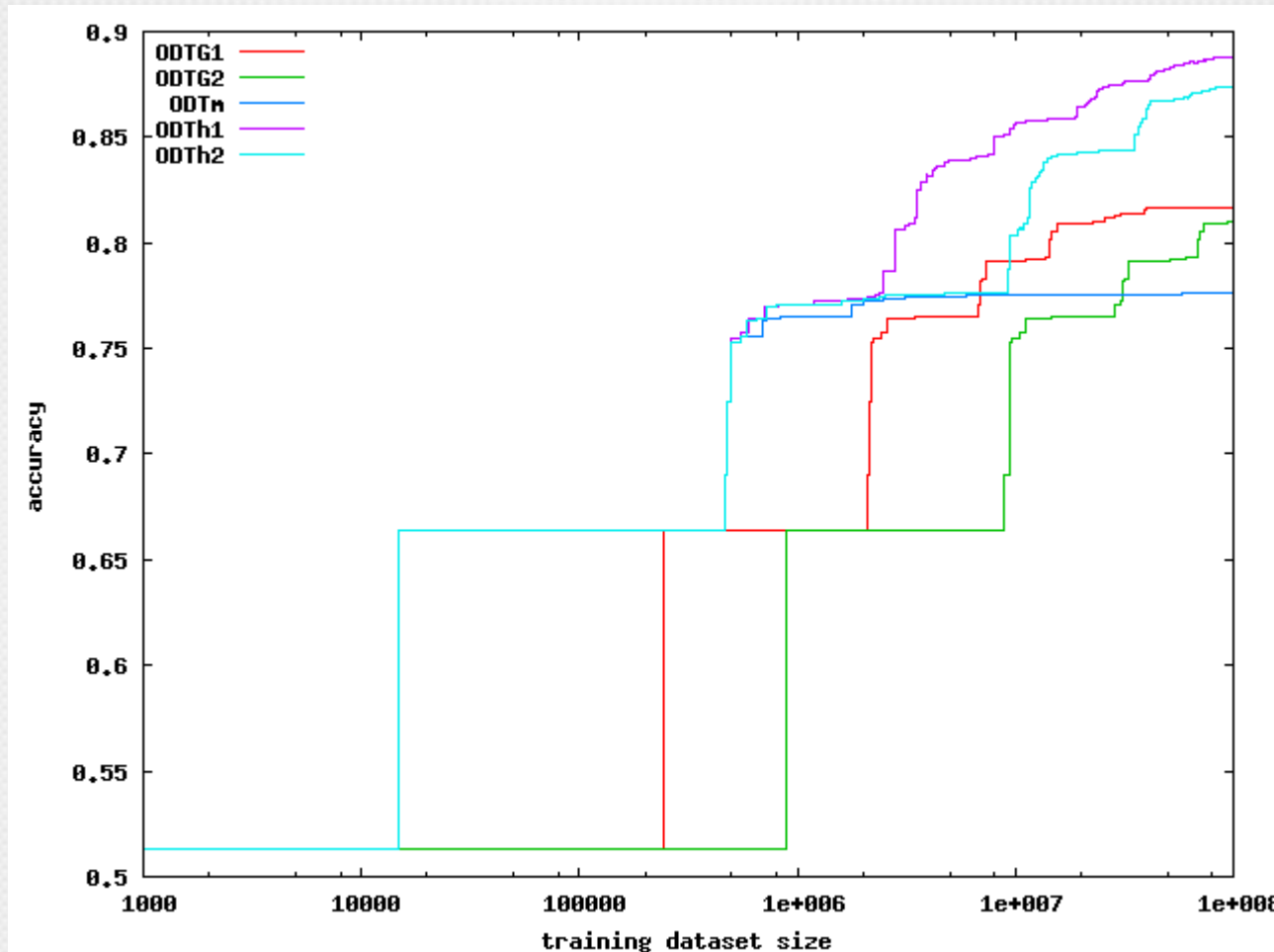
New result
2016

No.	Notation	Splitting criterion
1	ODTm	Misclassification error + Hoeffding's inequality
2	OTDG1	Gini index + McDiarmid's inequality
3	ODTG2	Gini index + McDiarmid's inequality + bias
4	ODTh1	Hybrid: 1 + 2
5	ODTh2	Hybrid: 1 + 3

Experimental Results

New result
2016

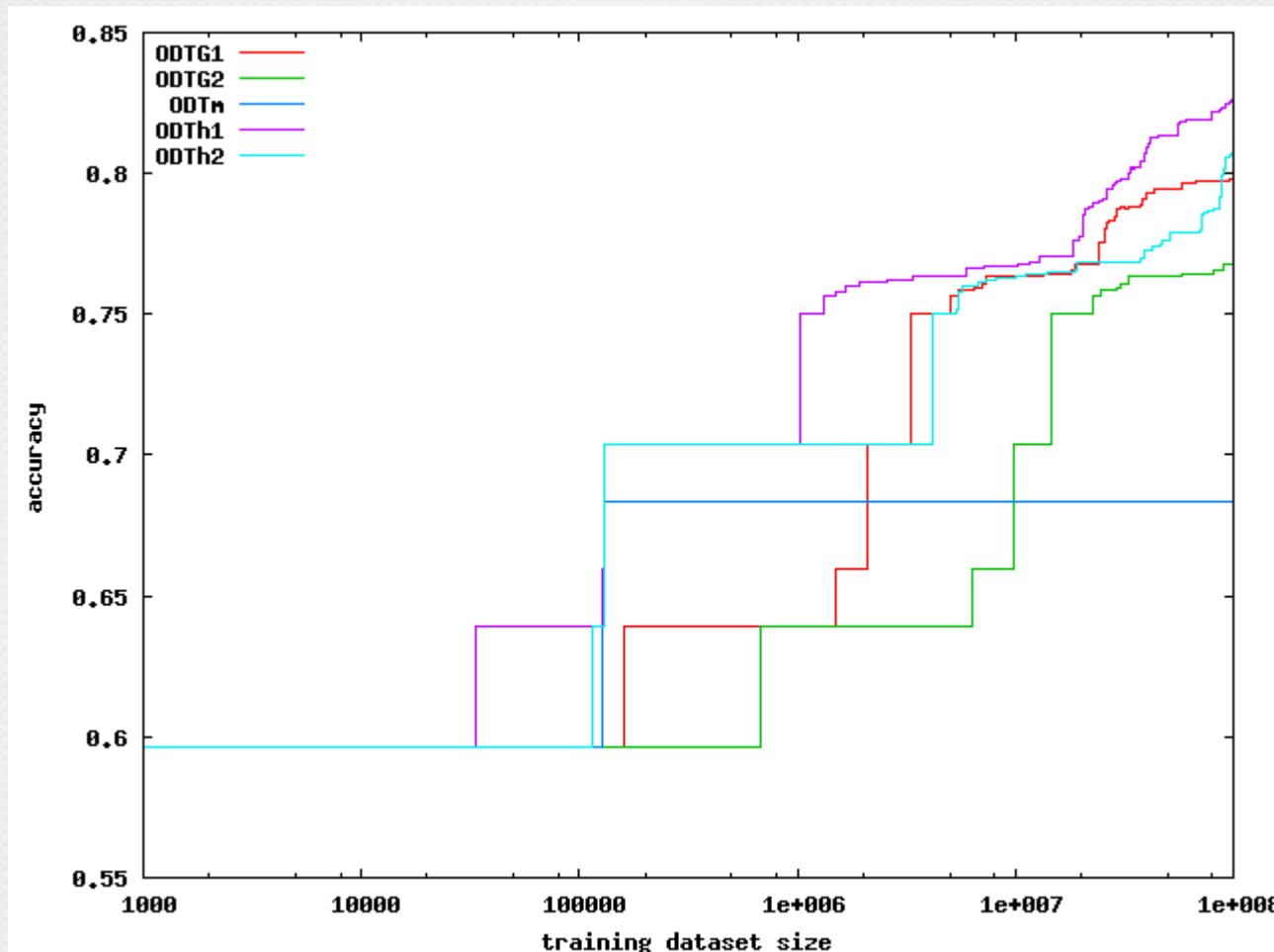
Accuracy vs number of data elements: dataset no. 1



Experimental Results

New result
2016

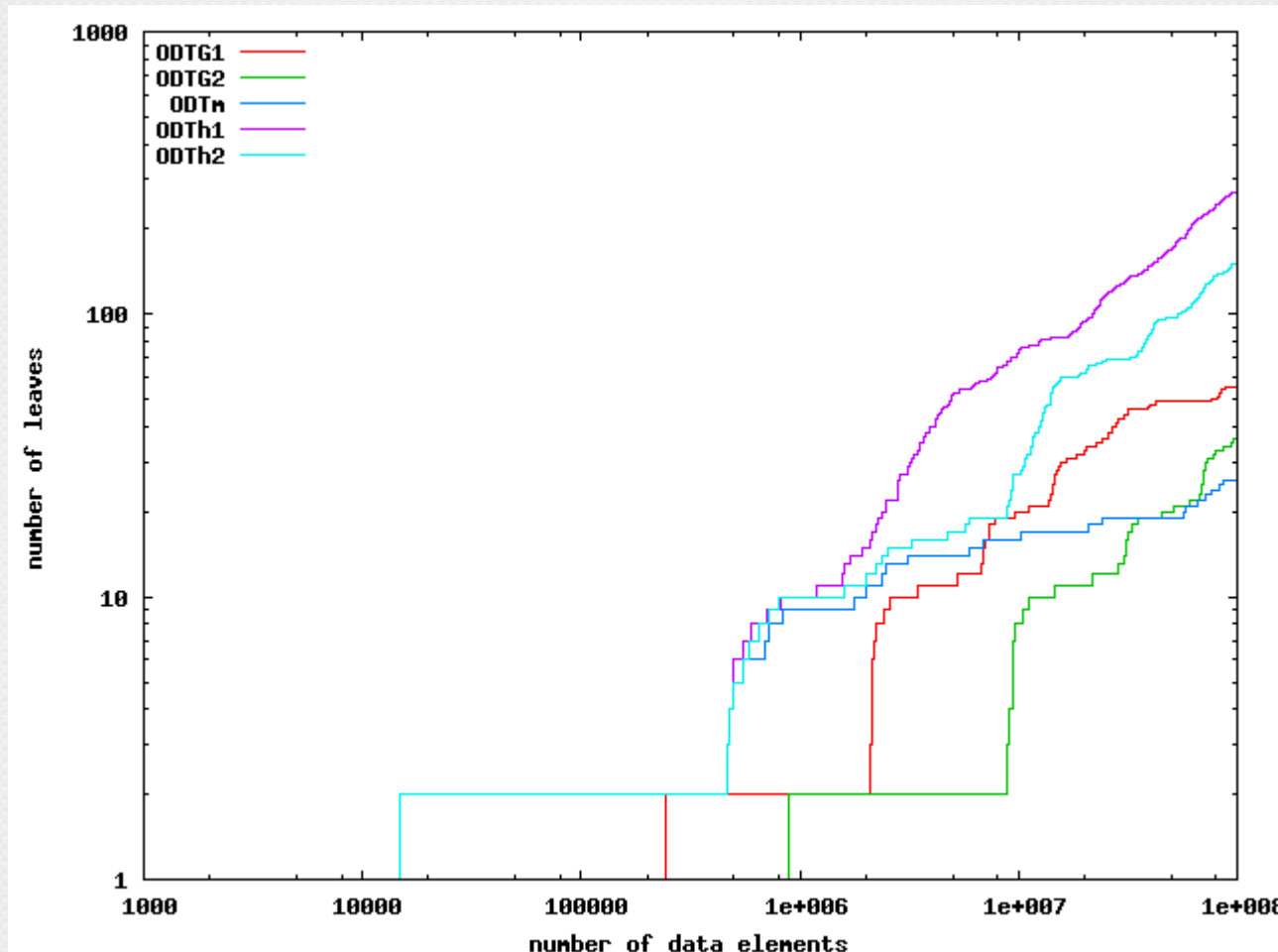
Accuracy vs number of data elements: dataset no. 2



Experimental Results

New result
2016

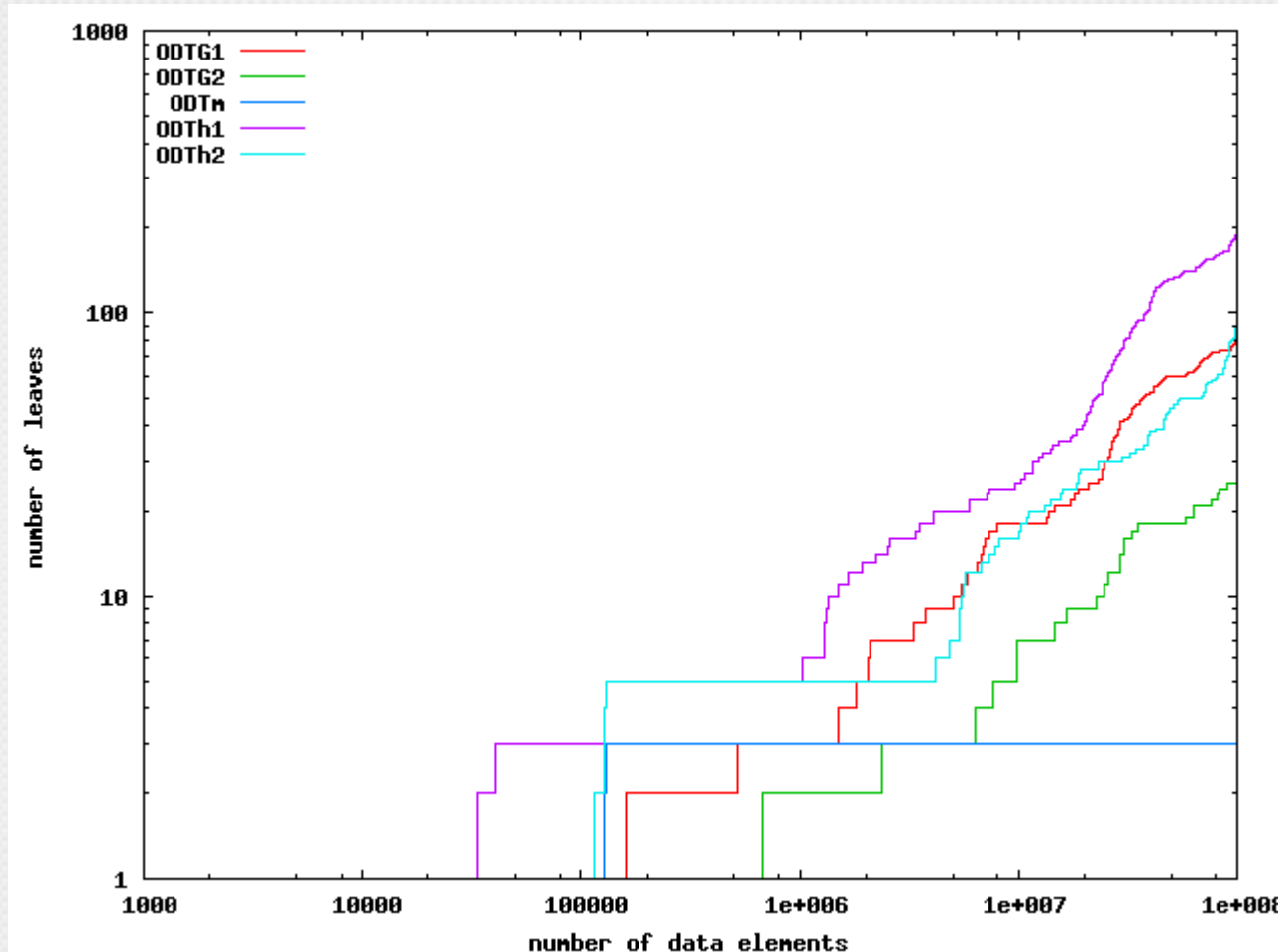
Number of leaves vs number data elements: dataset no. 1



Experimental Results

New result
2016

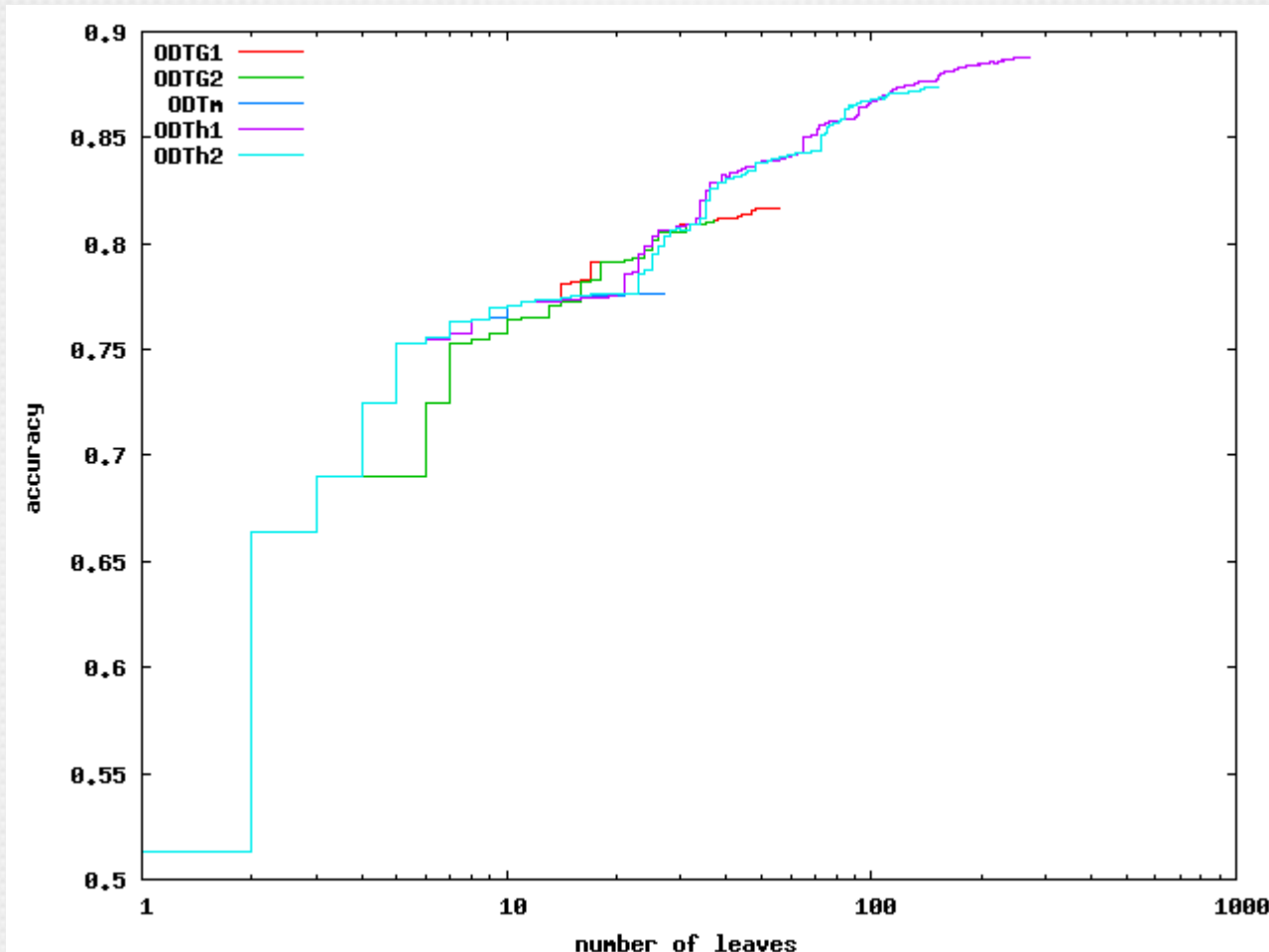
Number of leaves vs number data elements: dataset no. 2



Experimental Results

New result
2016

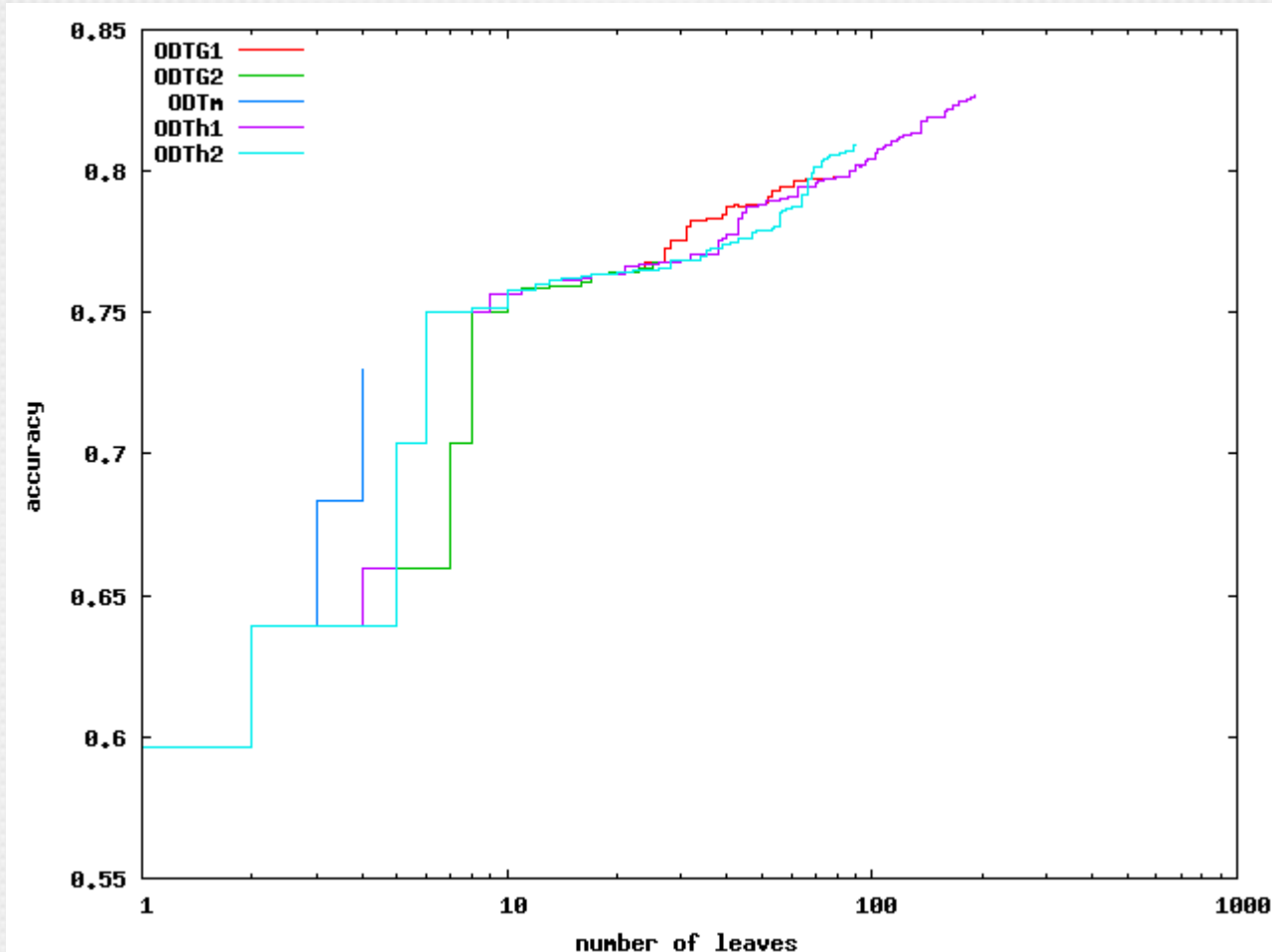
Accuracy vs number of leaves: dataset no. 1



Experimental Results

New result
2016

Accuracy vs number of leaves: dataset no. 2



Comparison with 'Hoeffding Trees'

New result
2016

- HDT:

$$f_{a_{max}}(Z) - f_{a_{max2}}(Z) > \sqrt{\frac{R \ln\left(\frac{1}{\delta}\right)}{2n(Z)}}$$

(for Gini index $R = 1$)

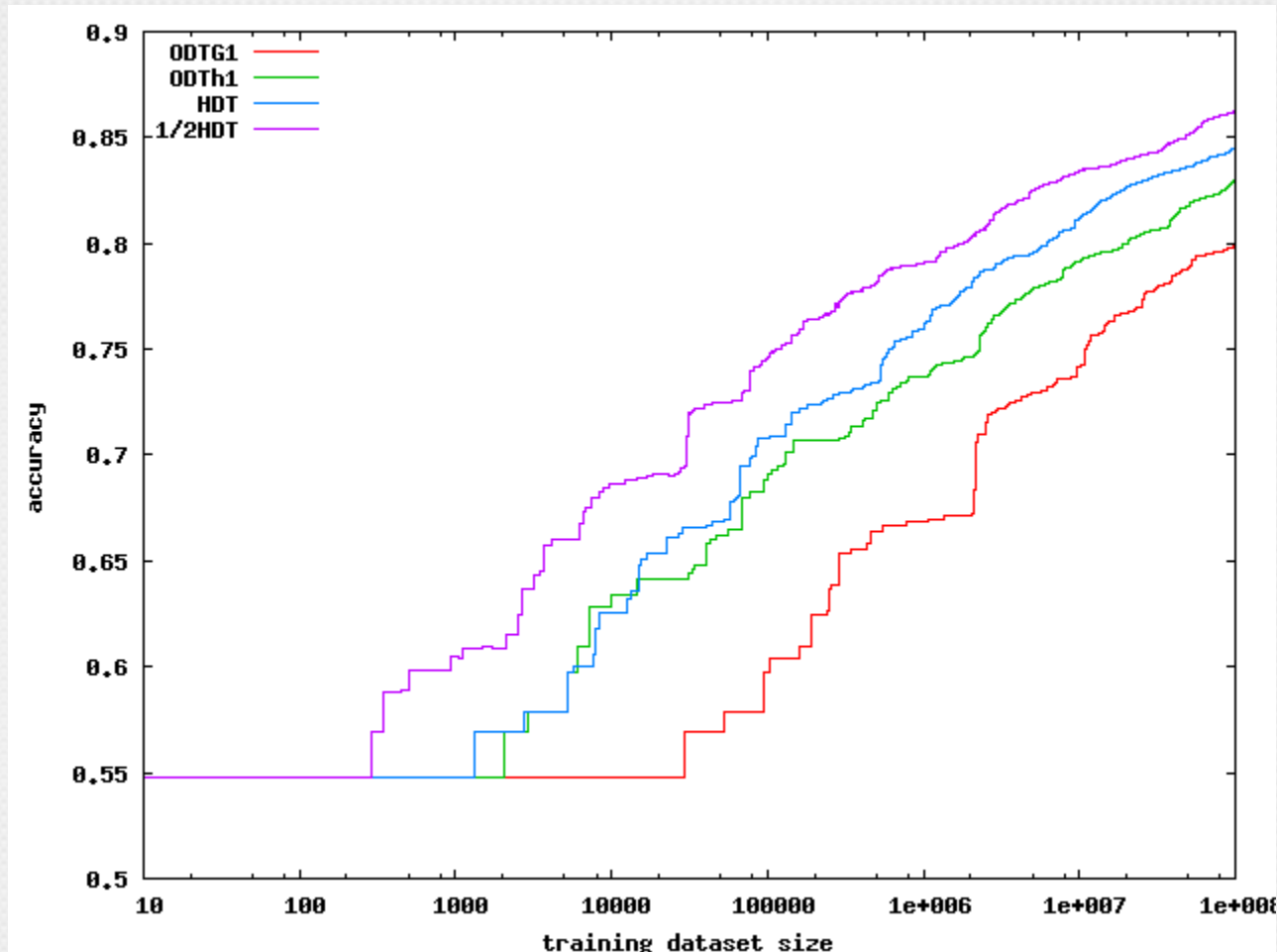
- 1/2HDT:

$$f_{a_{max}}(Z) - f_{a_{max2}}(Z) > 0.5 \sqrt{\frac{R \ln\left(\frac{1}{\delta}\right)}{2n(Z)}}$$

Experimental Results

New result
2016

Accuracy vs number of data elements



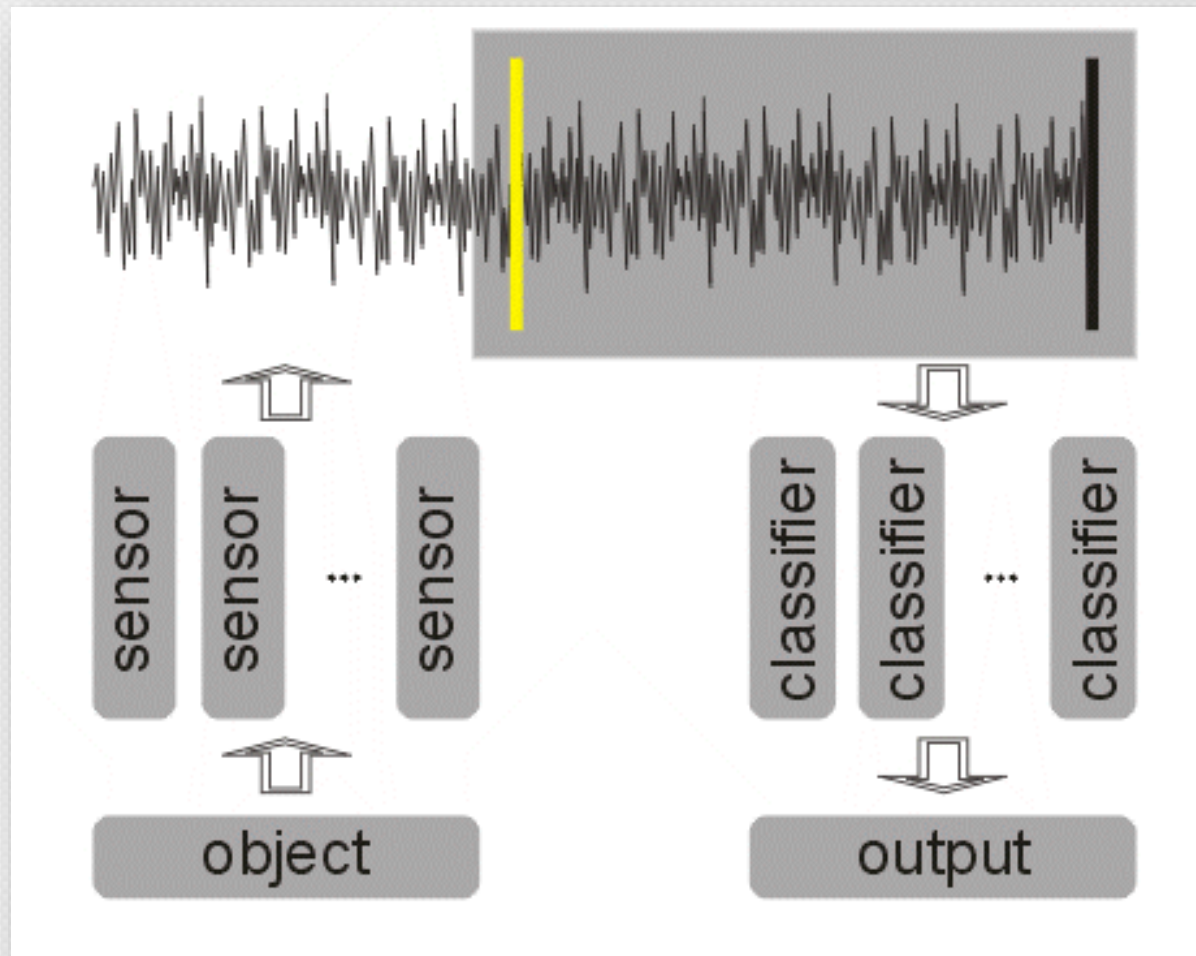
Data stream mining

- content

- Data streams – introduction to the topic
- Concept drift
- Various strategies of learning
- How to deal with concept drift?
- Data stream classification methods – short overview
- Decision trees for data streams (including new results 2016)
- **Ensemble methods for data streams (including new results 2016)**
- Probabilistic neural networks for stream data mining (including new results 2016)
- Final remarks and challenging problems
- References

Stream Data

- Ensemble methods

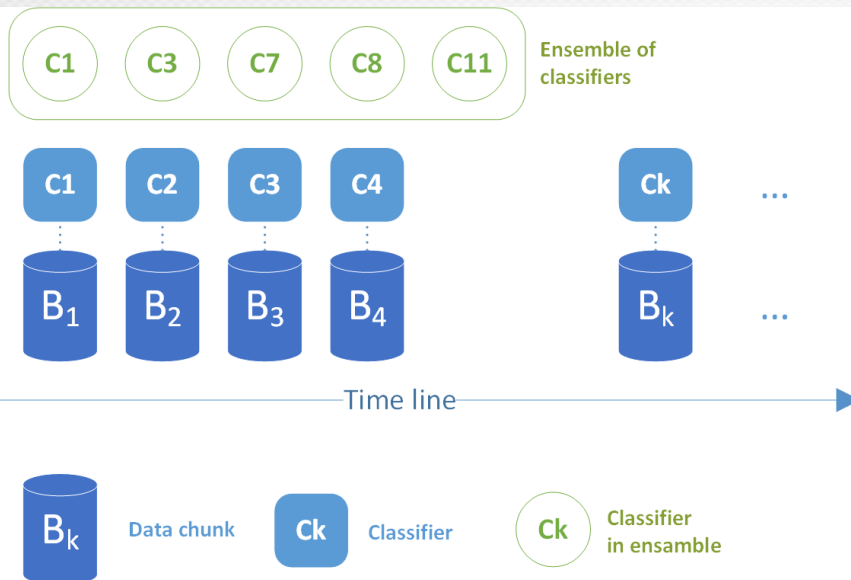


Stream Data

- Ensemble methods

Main steps of ensemble methods:

- Divide the data stream into equal sized chunks
- Train a new classifier based on current data chunk
- Keep the best L classifiers in the ensemble, e.g. $L=5$



Stream Data

- Ensemble methods

Main questions with designing ensemble system for data stream mining:

- Which components to use?
- How to train the components?
- How big the ensemble should be?
- How to establish the usefulness of component to the ensemble?
- How to assign weights to the components?
- How to reduce the ensemble size?

Stream Data

- Ensemble methods (the Adaboost)

Algorithm The Adaboost algorithm for multiclass problem

Require: S - set of N data elements X_1, X_2, \dots, X_N , $|\Gamma|$ - number of components in the ensemble

```
1: set the initial values of weights  $w(X_j) = \frac{1}{N}$ 
2: for each  $i \in \{1, \dots, |\Gamma|\}$  do
3:   create a new training dataset  $S^{(i)}$ 
4:   draw with replacement  $N$  data elements from  $S$  into  $S^{(i)}$  with probability equal
   to the value of it's weights
5:   create and train new classifier on  $S^{(i)}$ 
6:   calculate the classification error  $Err_i = \sum_{j=1}^N \mathbb{1}_{\{C(X_j) \neq \tau_i(X_j)\}} w(X_j)$ 
7:   if  $Err_i < \frac{1}{N}$  then
8:     add this classifier to the ensemble  $\Gamma$ 
9:     update weights  $w(X_j) = \begin{cases} w(X_j) \frac{Err_i}{1-Err_i} & \text{if } C(X_j) = \tau_i(X_j) \\ w(X_j) & \text{if } C(X_j) \neq \tau_i(X_j) \end{cases}$ 
10:    for each  $j \in 1, \dots, N$  do
11:       $w(X_j) = \frac{w(X_j)}{\sum_{l=1}^N w(X_l)}$ 
12:    end for each
13:  end if
14: end for each
15: for each unclassified data element  $\hat{X}$  do
16:   classify  $\hat{X}$  according to the equation
   
$$\Gamma_{cl}(\hat{X}) = \operatorname{argmax}_{C \in \xi} \sum_{j=1}^N \mathbb{1}_{\{\tau_i(\hat{X}) = C(X)\}} \log \left( \frac{1-Err_j}{Err_j} \right)$$

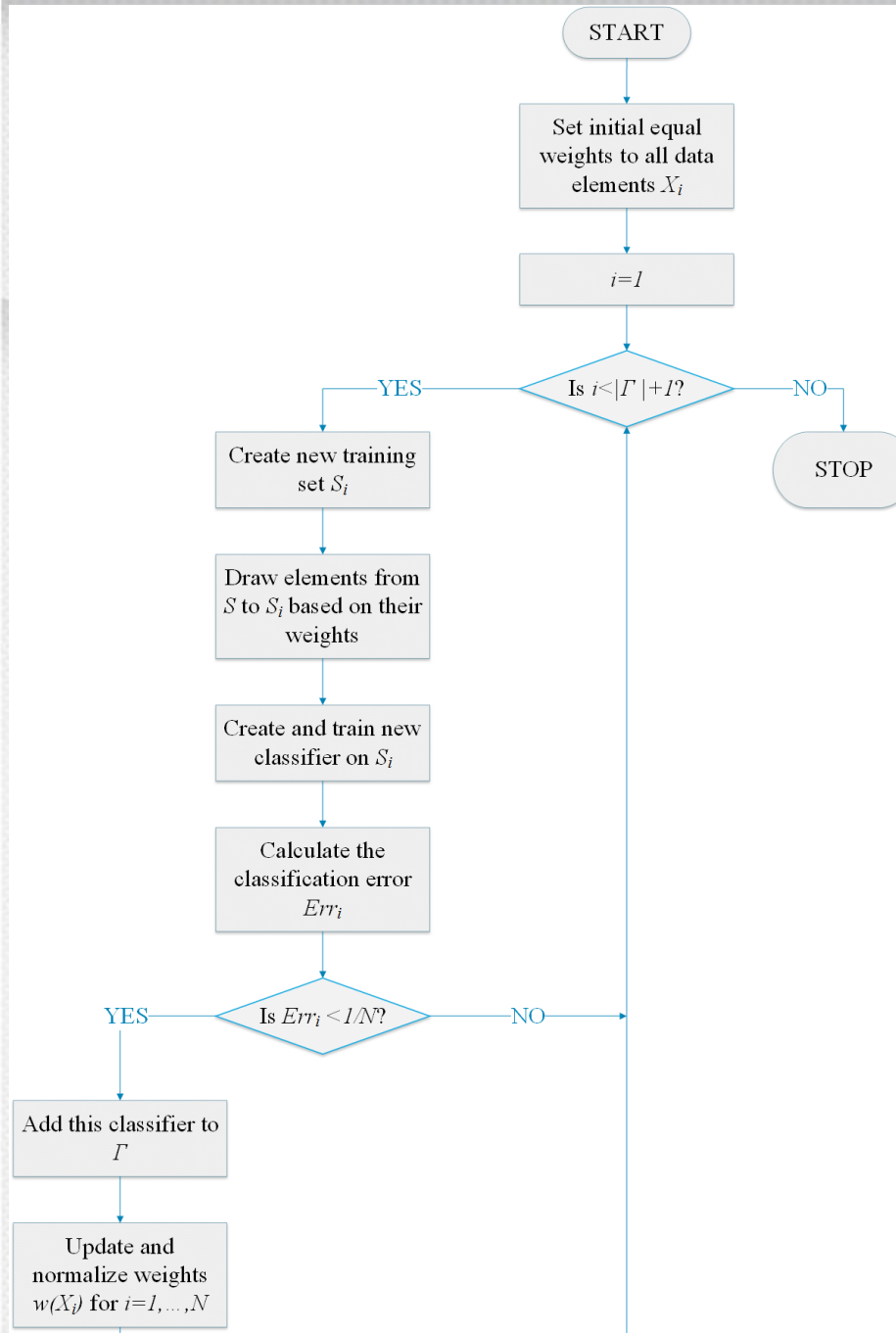
17: end for
```

* Yoav Freund and Robert E. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting, Journal of Computer and System Sciences, vol. 55, no. 1, pp. 119 –139, 1997.

Stream Data

- Ensemble methods

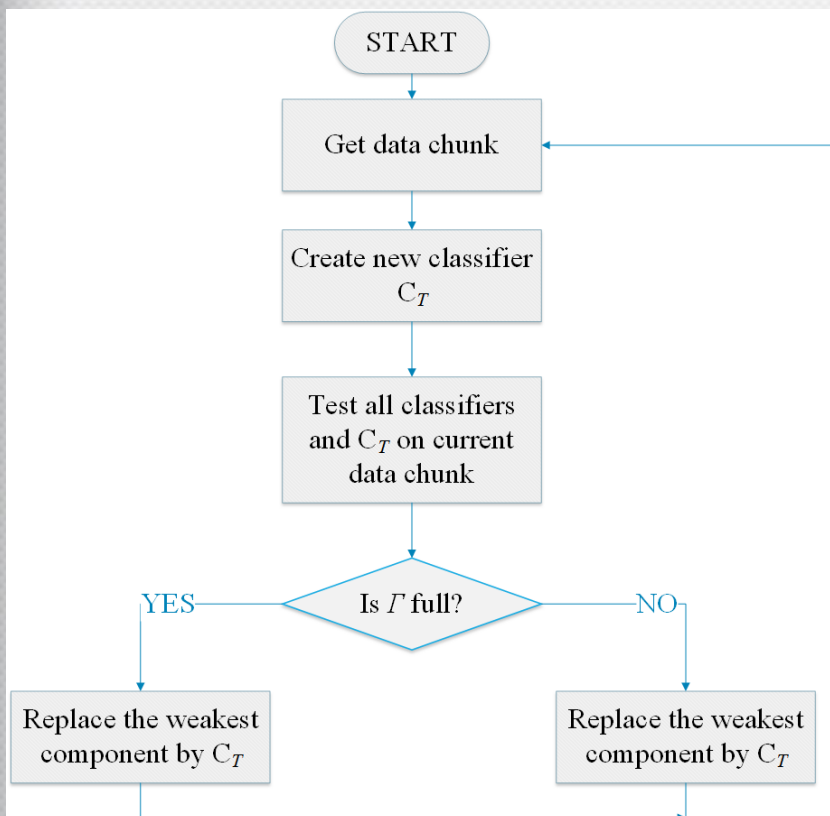
the Adabost



* Yoav Freund and Robert E. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting, Journal of Computer and System Sciences, vol. 55, no. 1, pp. 119 –139, 1997.

Stream Data

- Ensemble methods (the SEA)



P_1 denotes the percentage of votes for the most often occurring class,
 P_2 denotes the percentage of votes obtained by the second most frequently occurring class,
 P_C denotes the percentage obtained by the class of the correct class
 P_T denotes the percentage of the prediction of the new classifier.

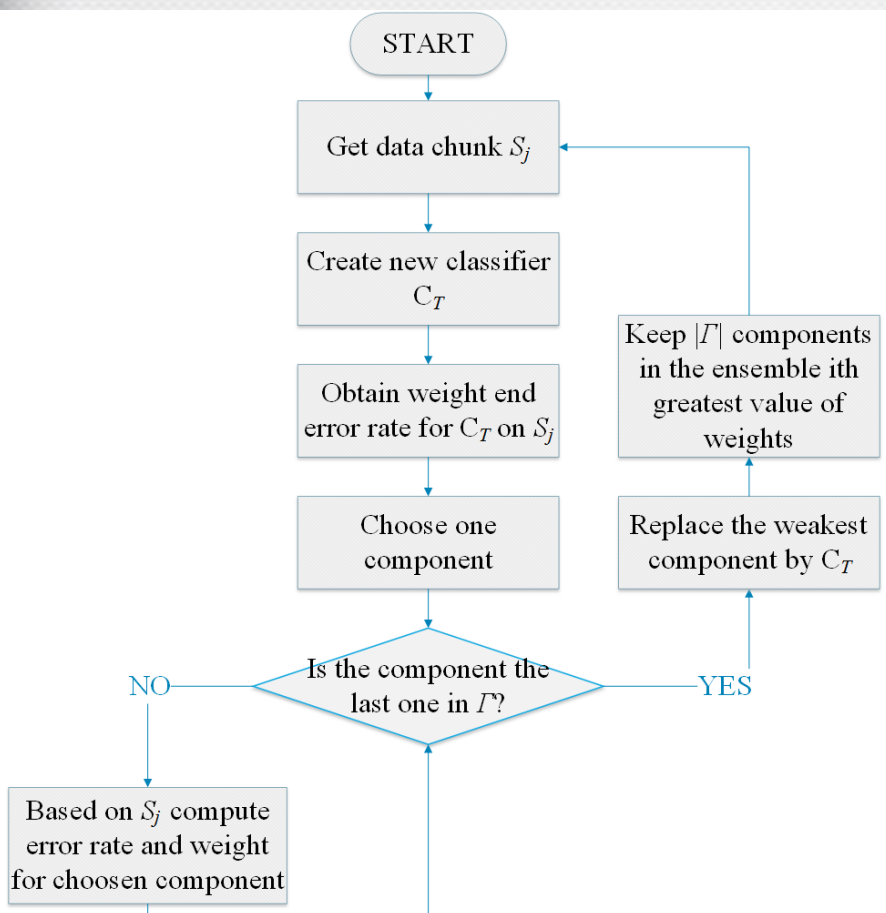
- the ensemble and the new tree give correct answers: the quality measure is increased by $1 - |P_1 - P_2|$
- the new tree gives correct answer but the ensemble is wrong: the quality measure is increased by $1 - |P_1 - P_C|$
- the tree gives wrong answer: the quality measure is decreased by $1 - |P_C - P_T|$

The obtained values of weights are used to determine if the newly created classifier should replace the least efficiently performing algorithm in the ensemble.

* W. Nick Street and Yong Seog Kim, A streaming ensemble algorithm (sea) for large-scale classification, In Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '01, pp. 377–382, New York, NY, USA, 2001, ACM.

Stream Data

- Ensemble methods (the AWE)



Accuracy Weighted Ensemble (AWE)

to each component in the ensemble the weight is assigned $\mathcal{W}(\tau_j) = MSE_r - MSE_j$ where

$$MSE_j = \frac{1}{|S^{(j)}|} \sum_{x_k \in S^{(j)}} (1 - \omega(\tau_j, C))^2,$$

$$MSE_r = \sum_C \omega(C) (1 - \omega(C))^2,$$

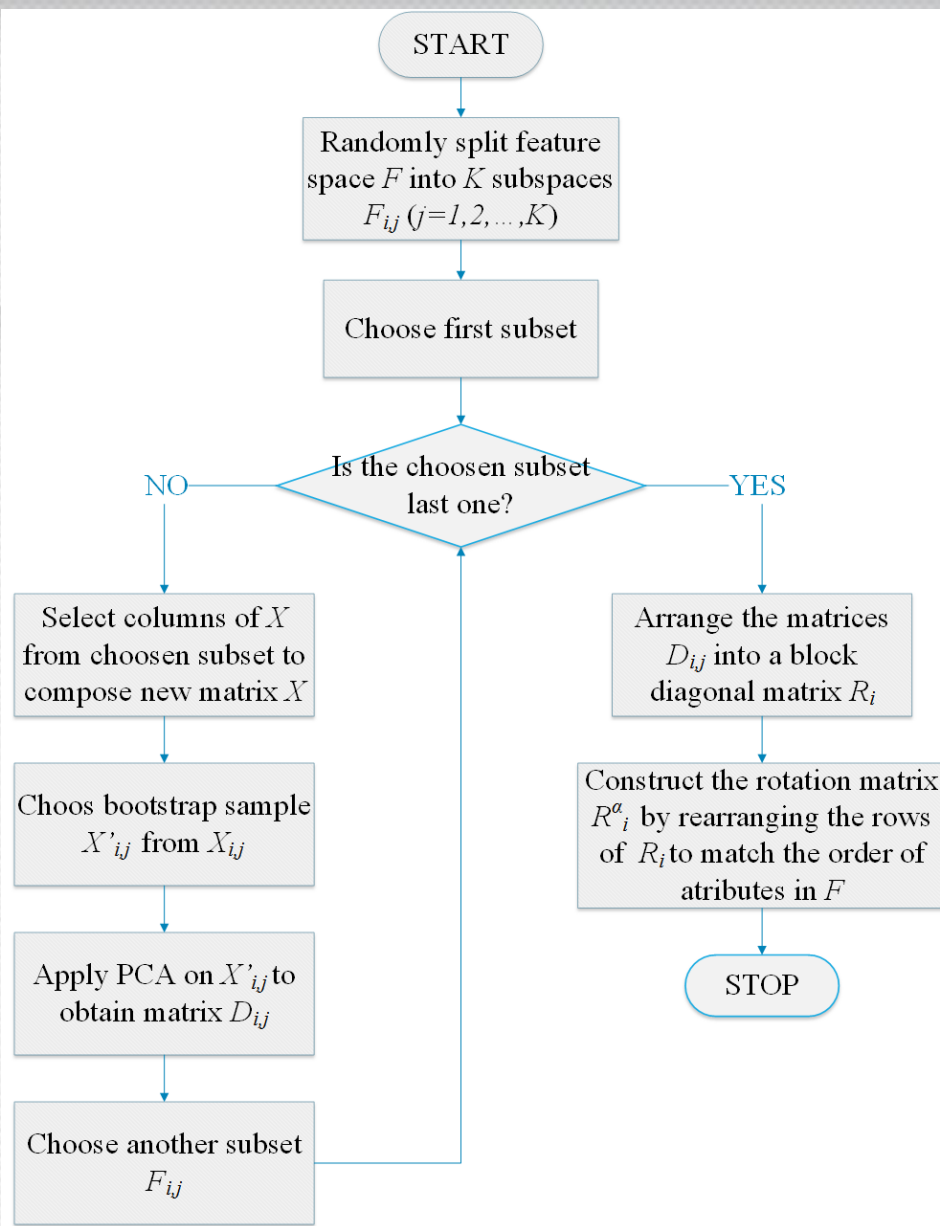
where $|S^{(j)}|$ denotes the number of data in the j -th chunk, $\omega(\tau_j, C)$ denotes the probability that classifier i will assign a correct class to data element and $\omega(C)$ denotes the probability of class C . In this case, MSE_r is the mean square error of randomly predicting classifier which defined the boundary for assessment of the minimal accuracy that the components in the ensemble should have. The members with lower accuracy are proven to decrease the performance of the whole ensemble method.

* Haixun Wang, Wei Fan, Philip S. Yu, and Jiawei Han. Mining concept-drifting data streams using ensemble classifiers. In Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '03, pp. 226–235, New York, NY, USA, 2003. ACM.

Stream Data

- Ensemble methods

the Rotation Forest



The main idea is to obtain different classifiers based on the sets of data that are transformed by the PCA algorithm. For construction of each decision tree a different mixture of features is taken into account and the principal component analysis (PCA) is applied on those sets. A new classifier is then trained based on data elements transformed linearly into new feature space. The diversity is obtained through the use of various extracted features which are the result of choosing different mixtures of feature components.

* Juan J. Rodriguez, Ludmila I. Kuncheva, and Carlos J. Alonso. Rotation forest: A new classifier ensemble method, IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 28, no. 10, pp. 1619–1630, Oct 2006.

Stream Data

- Ensemble methods

the Learn++.NSE

Algorithm The Learn++.NSE algorithm

Require: $S_\infty = \{X_1, X_2, \dots\}$ - dataset representing a new environment, S_t - sequence of data elements $X_i^t \in \{X_1^t, X_2^t, \dots, X_{N^t}^t\}$ from S_∞ , supervised learning algorithm **BaseClassifier**

- 1: **for all** $S_t, t = 1, 2, \dots$ **do**
- 2: **if** $t=1$ **then**
- 3: initialize $D^1(i) = w^t(i) = 1/m^1, \forall i$, skip to line nr 8
- 4: **end if**
- 5: compute error of existing ensemble on new data
 $E^t = \sum_{i=1}^{m^t} (1/m^t) [\Gamma^{t-1}(X_i) = C(X_i)]$
- 6: normalize error $B^t = E^t / (1 - E^t)$
- 7: update instance weights
 $w_i^t = \frac{1}{m^t} \times \begin{cases} B^t, & \Gamma^{t-1}(X_i) = C(X_i) \\ 1, & \text{otherwise} \end{cases}$
- 8: set $D^t = w^t / \sum_{i=1}^{m^t} w_i^t$ so that D^k is a distribution
- 9: call **BaseClassifier** with D_t , obtain $\tau_t : X \rightarrow C$
- 10: evaluate all existing classifiers on new dataset D_t
 $\epsilon_k^t = \sum_{i=1}^{m^t} D^t(i) [\tau_k(X_i) \neq C(X_i)], \text{ for } k = 1, \dots, t$
- 11: **if** $\epsilon_{k=t}^t > 1/2$ **then**
- 12: generate a new τ_t .
- 13: **end if**
- 14: **if** $\epsilon_{k<t}^t > 1/2$ **then**
- 15: $\epsilon_k^t = 1/2$
- 16: $\beta_k^t = \epsilon_k^t / (1 - \epsilon_k^t)$ for $k = 1, \dots, t$
- 17: **end if**
- 18: compute a weighted sum of all normalized errors for k -th classifier τ_k
- 19: $\omega_k^t = 1 / (1 + \exp -a(t - k - b)), \omega_k^t = \omega_k^t / \sum_{j=0}^{t-k} \omega_k^{t-j}$
- 20: $\hat{\beta}_k^t = \sum_{j=0}^{t-k} \omega_k^{t-j} \beta_k^{t-j}$, for $k = 1, \dots, t$
- 21: calculate classifier voting weights $\mathcal{W}_k^t = \log(1/\hat{\beta}_k^t)$, for $k = 1, \dots, t$
- 22: obtain the composite hypothesis as the current final hypothesis
 $\Gamma^t(X_i) = \operatorname{argmax}_{\{C\}} \sum_k \mathcal{W}_k^t [k_k(X_i) = C]$
- 23: **end for**

* Ryan Elwell and Robi Polikar, Incremental learning of concept drift in nonstationary environments, IEEE Transactions on Neural Networks, vol. 22, no. 10, pp. 1517–1531, Oct 2011.

Stream Data

- Ensemble methods (the Adaptive Ensemble Classifier)

Algorithm The Adaptive Ensemble Classifier

Require: $S = \{X_1, X_2, \dots, X_N\}$ - training dataset

```
1: for all data element  $X_i \in S_\infty$  do
2:   initialize the weight,  $w_i \rightarrow X_i \in S$  using Instance_Weighting procedure
3: end for
4: for  $N \in \{1, 2, 3\}$  do
5:   if  $N=1$  then
6:     generate a new dataset,  $S_{NEW}$ , form  $S$  using a selection and replacement technique
7:   else
8:     generate a new dataset,  $S_{NEW}$ , form  $S$  with  $X_i$  of higher weights
9:   end if
10:  build a decision tree,  $\tau$ , from  $S_{new}$ 
11:  for all leaf node in  $T$  do
12:    cluster the data instances of  $S_{new}$  using Similarity_Based_Clustering
13:    calculate the threshold value based on the ratio of percentage of instances in this leaf node and instances in  $S_{new}$ 
14:  end for
15:  initialize the weight,  $w_i \rightarrow \tau$ , based on its classification accuracy rate for the categorization of the instances,  $X_i \in S$ 
16:  update the weight,  $w_i$ , of each  $X_i \in S$ 
17: end for
```

* Dewan Md. Farid, Li Zhang, Alamgir Hossain, Chowdhury Mofizur Rahman, Rebecca Strachan, Graham Sexton, and Keshav Dahal, An adaptive ensemble classifier for mining concept drifting data streams, Expert Systems with Applications, vol. 40, no. 15, pp. 5895 – 5906, 2013.

Stream Data

- Ensemble methods (the AUE)

Algorithm The AUE algorithm

Require: S_∞ - stream of data elements X_1, X_2, \dots , size of chunks of data, Γ - ensemble of algorithms, $\hat{\Gamma}$ - set of the top best classifiers in the ensemble, $|\Gamma_{MAX}|$ - maximal number of components in the ensemble,

```
1:  $\Gamma \leftarrow \emptyset$ 
2: for all chunks of data  $S^{(j)}$  do
3:   create new classifier  $\tau_T$  based on  $S^{(j)}$ 
4:   compute error  $MSE$  of  $\tau_T$  via cross validation on  $S^{(j)}$ 
5:   obtain weight for  $\tau_T$  using (1)
6:   for all classifiers  $\tau_i$  in the ensemble  $\Gamma$  do
7:     test  $\tau_i$  on  $S^{(j)}$  to obtain  $MSE_i$  value
8:     compute weight of classifier  $\tau_i$  based on (1)
9:   end for
10:  assigned top weighted classifiers in  $\Gamma \cup \{\tau_T\}$  to  $\hat{\Gamma}$ 
11:   $\Gamma \leftarrow \Gamma \cup \{\tau_T\}$ 
12:  for all classifiers  $\tau_i$  in  $\hat{\Gamma}$  do
13:    if  $W_i > \frac{1}{MSE_r}$  and  $\tau_i \neq \tau_T$  then
14:      update classifier  $\tau_i$  with  $S^{(j)}$ 
15:    end if
16:  end for
17: end for
```

$$W_j = \frac{1}{MSE_j + \epsilon}$$

where ϵ is a very small constant value which is used to avoid the problem of division by zero.

* Dariusz Brzezinski and Jerzy Stefanowski, Reacting to different types of concept drift: The accuracy updated ensemble algorithm, IEEE Transactions on Neural Networks and Learning Systems, vol. 25, no. 1, pp. 81–94, Jan 2014.

Algorithm The OAUE algorithm

Require: S_∞ - stream of data elements X^t where t denotes the point in time, window size, Γ - ensemble of algorithms, $|\Gamma_{MAX}|$ - maximal number of algorithms in the ensemble, m - memory limit

```
1:  $\Gamma \leftarrow \emptyset$ 
2: create new classifier  $\tau_T$ 
3: for all data examples  $x^t \in S_\infty$  do
4:   calculate the prediction error of all classifiers  $\tau_i$  on  $x^t$ 
5:   if  $t > 0$  and  $t \bmod d = 0$  then
6:     if  $|\Gamma| < |\Gamma_{MAX}|$  then
7:        $\Gamma \leftarrow \Gamma \cup \{\tau_T\}$ 
8:     else
9:       weight all classifiers  $\tau_i \in \Gamma$  and  $\tau_T$  using (1)
10:      substitute least accurate classifier in  $\Gamma$  with  $\tau_T$ 
11:    end if  $\tau_T \leftarrow$  new candidate classifier
12:    if  $memory\_usage(\Gamma) > m$  then
13:      prune (decrease size of) component classifiers
14:    end if
15:  else
16:    incrementally train classifier  $\tau_T$  with  $X^t$ 
17:    weight all classifiers  $\tau_i \in \Gamma$  using (1)
18:  end if
19:  for all classifiers  $\tau_i \in \Gamma$  do
20:    incrementally train classifier  $\tau_i$  with  $X^t$ 
21:  end for
22: end for
```


Stream Data

- Ensemble methods (the OAUE)

The authors calculate and update the weight w_i^t for the i -th classifier in the ensemble in time t by using d latest examples in the following way

$$MSE_i^t = \begin{cases} MSE_i^{t-1} + \frac{e_i^t}{d} + \frac{e_i^{t-d}}{d}, & t - \tau_i > d \\ \frac{t - \tau_i - 1}{t - \tau_i} MSE_i^{t-1} + \frac{e_i^t}{t - \tau_i}, & 1 \leq t - \tau_i \leq d \\ 0, & t - \tau_i = 0 \end{cases},$$

$$e_i^t = (1 - f_{iC}^t(X^t))^2,$$

$$MSE_r^t = \begin{cases} MSE_r^{t-1} - r^{t-1}(C^t) - r^{t-1}(C^{t-d}) + r^t(C^t) + r^t(C^{t-d}), & t > d \\ \sum_C r^t(C) & t = d \end{cases},$$

$$r^t(C) = p^t(C)(1 - p^t(C))^2,$$

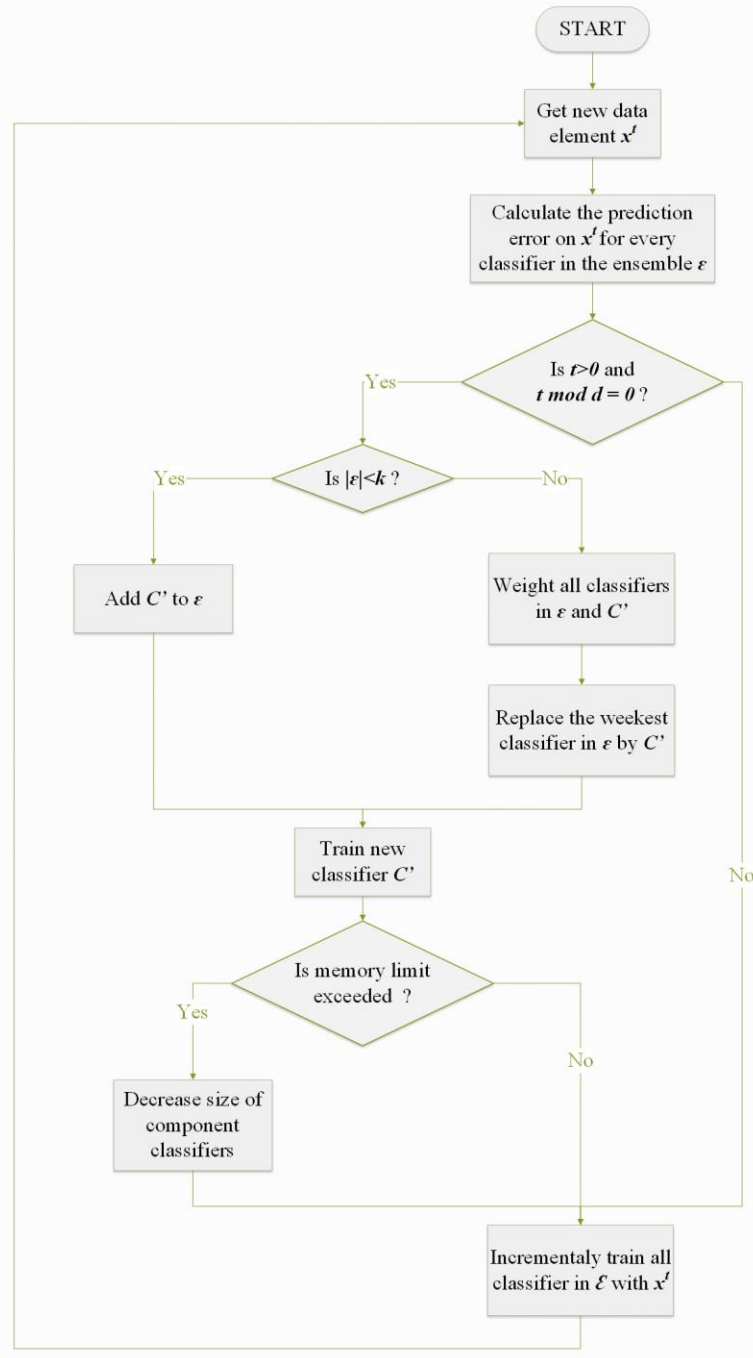
$$W_i^t = \frac{1}{MSE_r^t + MSE_i^t + \varepsilon},$$

where $f_{iC}^t(X^t)$ denotes the probability given by classifier τ_i that data element X^t collected in time t is from class C^t . To overcome the problem of division by zero, a very small value was added in the form of ϵ .

Stream Data

- Ensemble methods

the OAUE



Brzezinski D., Stefanowski J., „Combining block-based and online methods in learning ensembles from concept drifting data streams”, Information Sciences, Vol. 265, pp. 50-67, 2014

Stream Data

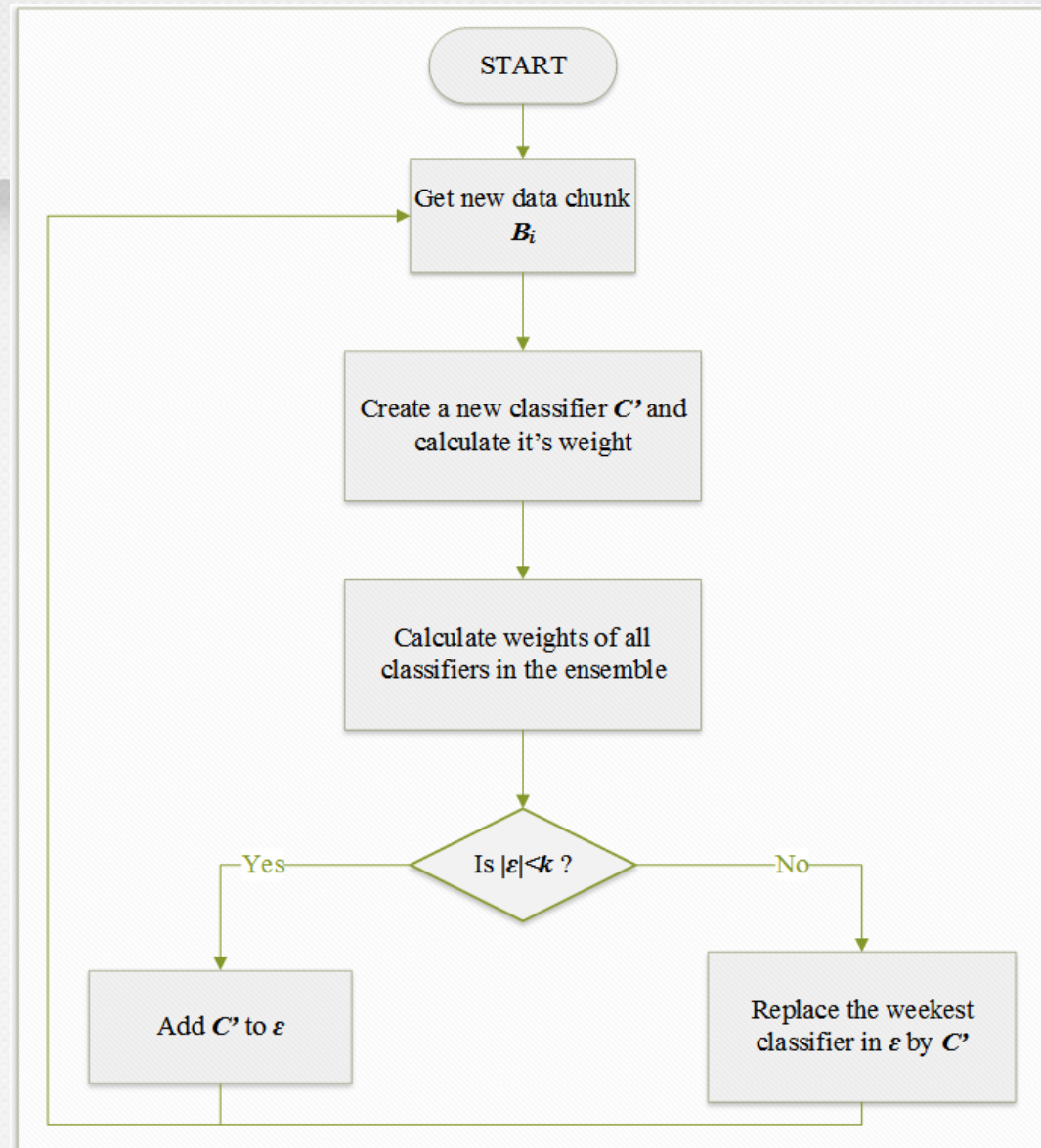
- Ensemble methods

General strategies for adaptation of algorithms for mining data streams based on ensemble methods.

Stream Data

- General Ensemble Strategies

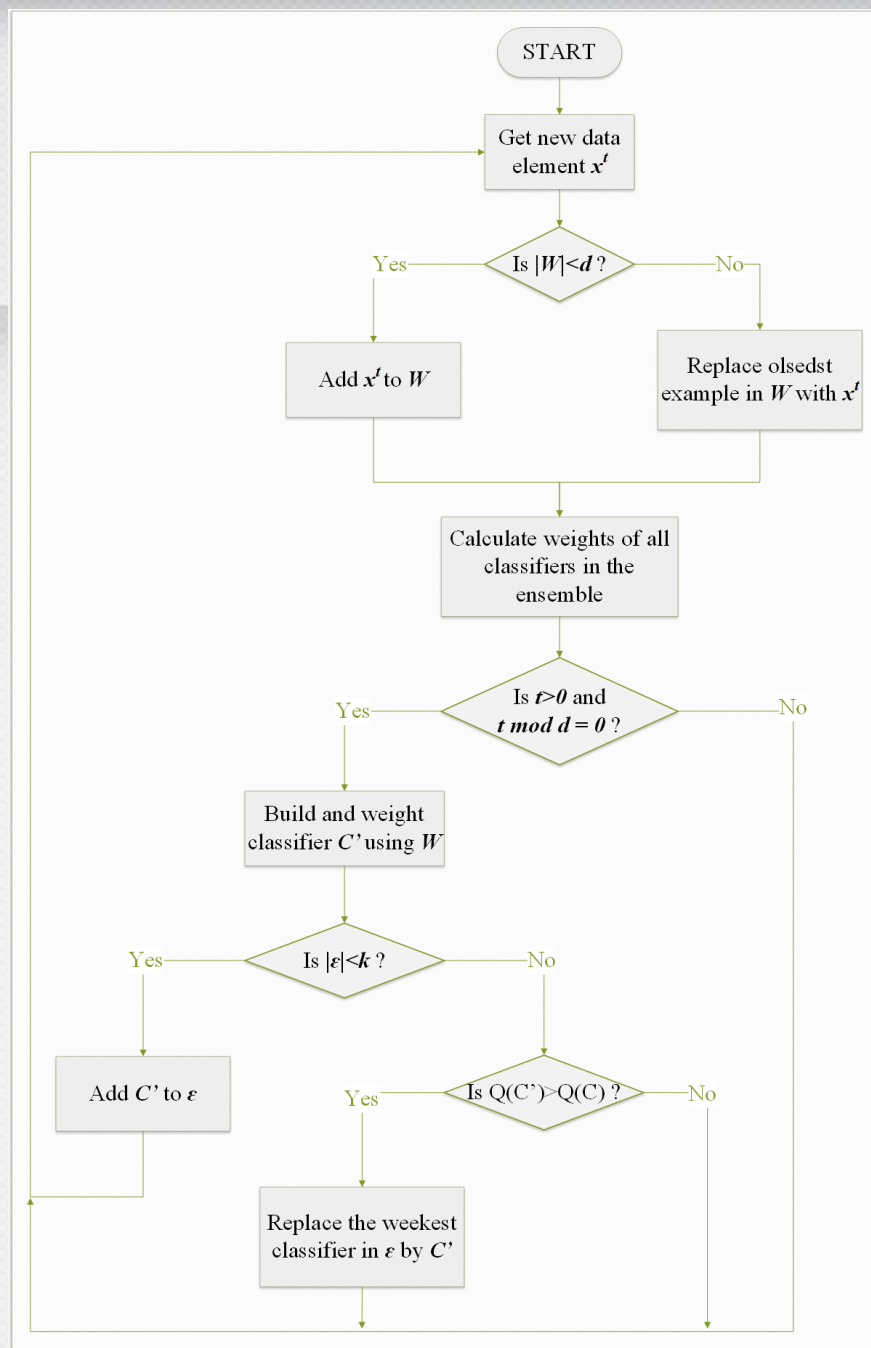
Generic block ensemble training scheme



Stream Data

- General Ensemble Strategies

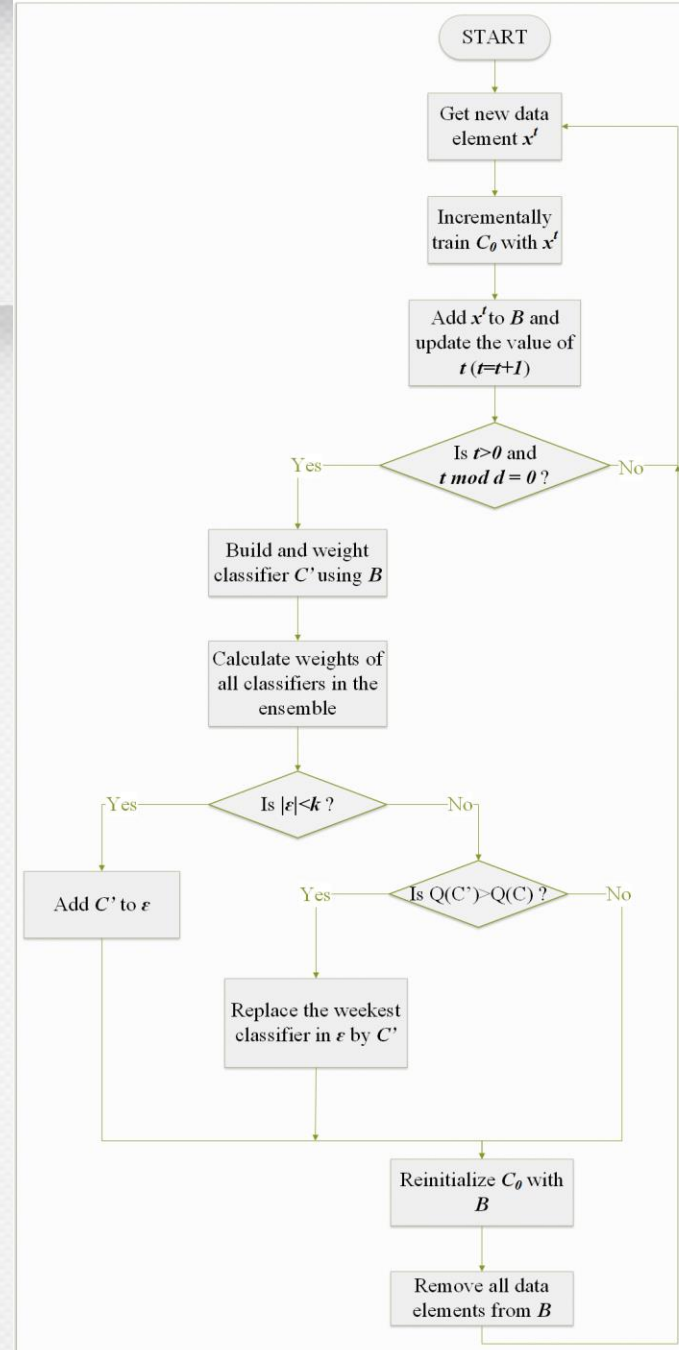
Windowing strategy



Stream Data

- General Ensemble Strategies

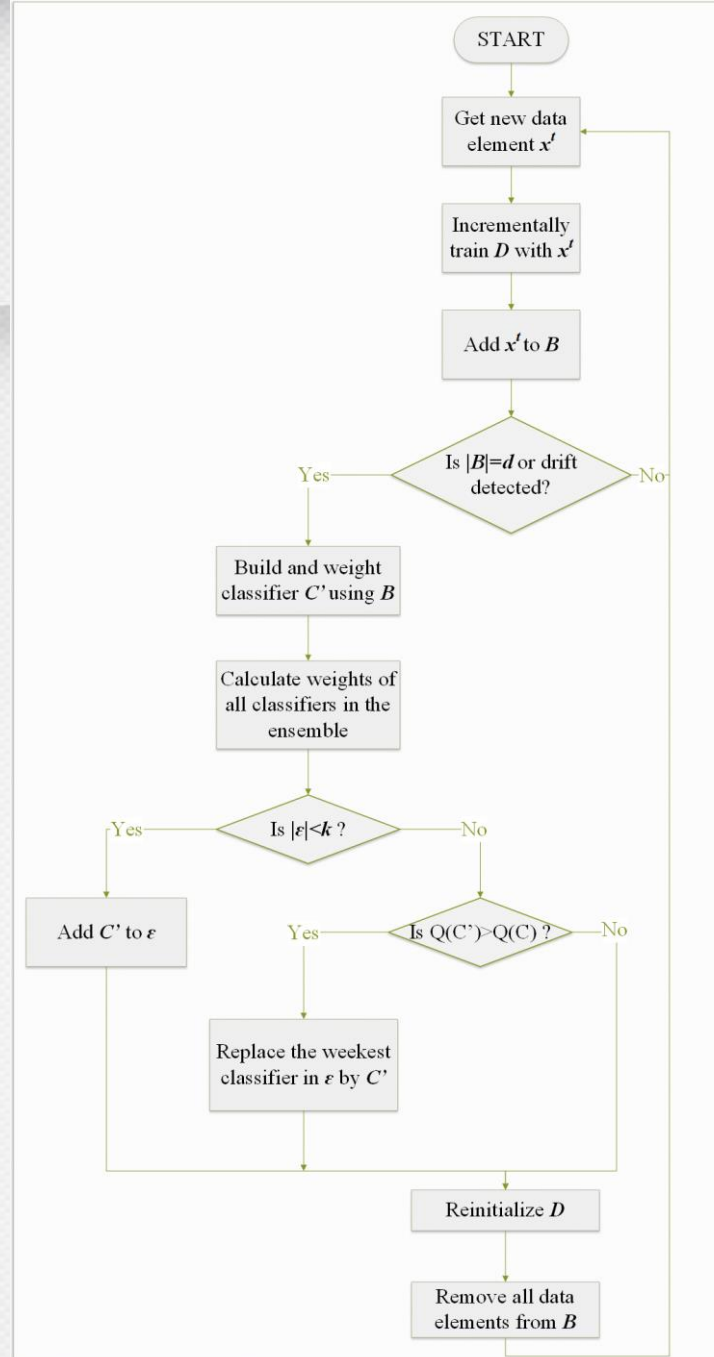
Additional incremental learner strategy



Stream Data

- General Ensemble Strategies

Drift detector strategy



Stream Data

- Ensemble methods

How to decide if a new component should be added to the ensemble?

New result 2016

Stream Data

- The Automatically Adjusting Size of Ensemble Algorithm (ASE)

Theorem 1 Let S_∞ denote the data stream and $S_n \subset S_\infty$ is a set of n independent random variables $S_n = \{X_1, X_2, \dots, X_n\}$. Moreover, let Γ and Γ^+ denote two ensembles of components where $\Gamma = \{\tau_1, \tau_2, \dots, \tau_{|\Gamma|}\}$ and $\Gamma^+ = \Gamma \cup \{\tau_{|\Gamma|+1}\}$. If the following inequality is true

$$P_{\Gamma^+}(S) - P_\Gamma(S) > z_{1-\gamma_1} \frac{1}{\sqrt{N}},$$

where $z_{1-\gamma_1}$ is the $(1 - \gamma_1)$ quantile of the standard normal distribution $\mathcal{N}(0, 1)$ and $P_{\Gamma^+}(S)$ ($P_\Gamma(S)$) denote the accuracy of ensemble Γ^+ (Γ), then with probability $1 - \gamma_1$

$$P_{\Gamma^+}(S_\infty) - P_\Gamma(S_\infty) > 0.$$

Stream Data

- Ensemble methods

**How to decide if we should to
remove a component from the
ensemble?**

New result
2016

Stream Data

- The Automatically Adjusting Size of Ensemble Algorithm (ASE)

Theorem 2 Let S_∞ denote the data stream and $S_n \subset S_\infty$ is a set of n independent random variables $S_n = \{X_1, X_2, \dots, X_n\}$. Moreover, let Γ and Γ^- denote two ensembles of algorithms where $\Gamma = \{\tau_1, \tau_2, \dots, \tau_{|\Gamma|}\}$ and $\Gamma^- = \{\tau_1, \tau_2, \dots, \tau_{t-1}, \tau_{t+1}, \dots, \tau_{|\Gamma|}\}$. If the following inequality is true

$$P_\Gamma(S) - P_{\Gamma^-}(S) > z_{1-\gamma_2} \frac{1}{\sqrt{N}},$$

where function $P_\Gamma(S)$ ($P_{\Gamma^-}(S)$) denote the accuracy of ensemble Γ (Γ^-), then with probability $1 - \gamma_2$

$$P_\Gamma(S_\infty) - P_{\Gamma^-}(S_\infty) > 0.$$

Therefore, if the above conditions are true, then the removal of classifier τ_t will decrease the accuracy of the ensemble.

New result 2016

Stream Data

- The Automatically Adjusting Size of Ensemble Algorithm (ASE)

Algorithm The Automatically Sized Ensemble Algorithm

Require: S_∞ - stream of data elements, value of parameters

$N, \gamma_1, \gamma_2, \eta$

```
1: while not the end of the stream do
2:   obtain new  $k$ -th chunk of  $N$  data elements  $S$ 
3:   create new component  $\tau_{temp}$  and
4:   calculate tree weights
5:    $\Gamma^+ = \Gamma \cup \{\tau_{temp}\}$ 
6:   for  $i = 1, \dots, N$  do
7:     obtain value of  $G_\Gamma(X_i)$  and  $G_{\Gamma^+}(X_i)$ 
8:   end for
9:   calculate  $P_\Gamma(S)$  and  $P_{\Gamma^+}(S)$  using (2)
10:  if  $P_\Gamma(S) - P_{\Gamma^+}(S) - \frac{z_{1-\gamma_1}}{\sqrt{N}} > 0$  then
11:     $\Gamma = \Gamma \cup \tau_{temp}$ 
12:  else
13:    remove  $\tau_{temp}$ 
14:  end if
15:  if  $k \bmod \eta = 0$  then
16:    for  $i = 1, \dots, |\Gamma|$  do
17:       $\Gamma^- = \Gamma \setminus \{\tau_i\}$ 
18:      calculate  $P_\Gamma(S)$  and  $P_{\Gamma^-}(S)$ 
19:      if  $P_\Gamma(S) - P_{\Gamma^-}(S) \leq \frac{z_{1-\gamma_2}}{\sqrt{N}}$  then
20:        remove  $\tau_i$  from  $\Gamma$ 
21:      end if
22:    end for
23:  end if
24: end while
```

**New result
2016**

Stream Data

- The Automatically Adjusting Size of Ensemble Algorithm (ASE)

	Bayes	AUE	AWE	Bagging	Bagging Adwin	Boosting	WMA	Hoefding Adaptive Tree	Hoeffding Opt.Tree	Hoeffding Tree	AASE $\eta=20$	AASE $\eta=40$
GEN1	80.86%	95.55%	96.86%	96.19%	96.19%	95.90%	95.23%	94.13%	95.62%	92.17%	90.68%	90.23%
GEN2	50.00%	87.18%	86.33%	89.69%	88.71%	87.80%	89.35%	86.91%	89.16%	89.35%	92.42%	92.30%
GEN3	86.85%	70.94%	69.27%	79.72%	70.02%	70.26%	80.19%	71.93%	80.08%	80.08%	91.31%	91.06%
GEN4	74.88%	73.61%	73.87%	74.45%	73.63%	73.91%	73.73%	73.12%	73.81%	73.73%	90.46%	90.30%

New result 2016

Stream Data

- The Automatically Adjusting Size of Ensemble Algorithm (ASE)

	GEN1		GEN2		GEN3		GEN4	
γ_1	max tree number	end tree number	max tree number	end tree number	max tree number	end tree number	max tree number	end tree number
0.01	7	5	5	2	6	3	9	2
0.05	9	6	7	2	7	5	14	5
0.10	14	9	10	4	11	7	21	7
0.15	17	11	14	4	13	9	19	9
0.20	22	13	18	6	18	10	29	12
Accuracy								
0.01	91.89%		92.45%		91.91%		87.84%	
0.05	92.37%		93.19%		92.62%		88.97%	
0.10	93.08%		93.96%		93.16%		89.93%	
0.15	93.24%		94.19%		93.41%		90.27%	
0.20	93.43%		94.38%		93.67%		90.47%	

New result 2016

Stream Data

- The Automatically Adjusting Size of Ensemble Algorithm (ASE)

	GEN1		GEN2		GEN3		GEN4	
γ_2	max tree number	end tree number	max tree number	end tree number	max tree number	end tree number	max tree number	end tree number
0.01	21	13	18	6	18	10	26	11
0.05	22	13	18	6	18	10	29	12
0.10	24	15	17	6	17	11	31	11
0.15	23	15	17	5	17	10	27	11
0.20	23	15	17	7	17	12	26	11
Accuracy								
0.01	93.43%		94.33%		93.68%		90.48%	
0.05	93.43%		94.38%		93.67%		90.47%	
0.10	93.30%		94.33%		93.65%		90.46%	
0.15	93.28%		94.26%		93.53%		90.56%	
0.20	93.27%		94.29%		93.64%		90.47%	

**New result
2016**

Stream Data

- The Automatically Adjusting Size of Ensemble Algorithm (ASE)

	GEN1	GEN2	GEN3	GEN4
0.01	539	458	484	613
0.05	539	432	472	620
0.10	517	410	475	615
0.15	521	399	454	605
0.20	509	397	438	569

difference	30	61	46	44
-------------------	-----------	-----------	-----------	-----------

New result
2016

Stream Data

- The Dynamically Expanded Ensemble Algorithm (DEEA)

Algorithm The Dynamically Expanded Ensemble Algorithm

Require: $S_\infty, \Psi \geq 0$

```
1: get first data chunk and create (learn)  $\tau$ 
2: get next data chunk and create (learn)  $\tau_{temp}$ 
3: add  $\tau$  to  $\Gamma$  and  $\Gamma^+$ , and  $\tau_{temp}$  to  $\Gamma^+$ 
4: while not the end of stream do
5:   get new data chunk  $S = X_1, X_2, \dots, X_n$ 
6:   for  $i = 1, 2, \dots, n$  do
7:     obtain value of  $G_\Gamma(X_i)$  and  $G_{\Gamma^+}(X_i)$ 
8:   end for
9:   calculate  $P_\Gamma(S)$  and  $P_{\Gamma^+}(S)$ 
10:  if  $P_{\Gamma^+}(S) - P_\Gamma(S) - \frac{z_{1-\gamma_2}}{\sqrt{n}} > \Psi$  then
11:     $\Gamma = \Gamma \cup \tau_{temp}$ 
12:  else
13:    remove  $\tau_{temp}$ 
14:  end if
15:  calculate weights (learn) for all components in  $\Gamma$ 
16:  create (learn) new component  $\tau_{temp}$ 
17:  add  $\tau_{temp}$  to  $\Gamma^+$ 
18: end while
```

New result 2016

Stream Data

- The Dynamically Expanded Ensemble Algorithm (DEEA)

It is the procedure of making a decision whether a particular classifier should be added to the ensemble.

First, an ensemble Γ^+ is created by combining classifiers from Γ and temporary classifier τ_{temp} created on the last data chunk. Next, a new chunk of data elements is collected. Based on the new data chunk the values of P_Γ and P_{Γ^+} are calculated using:

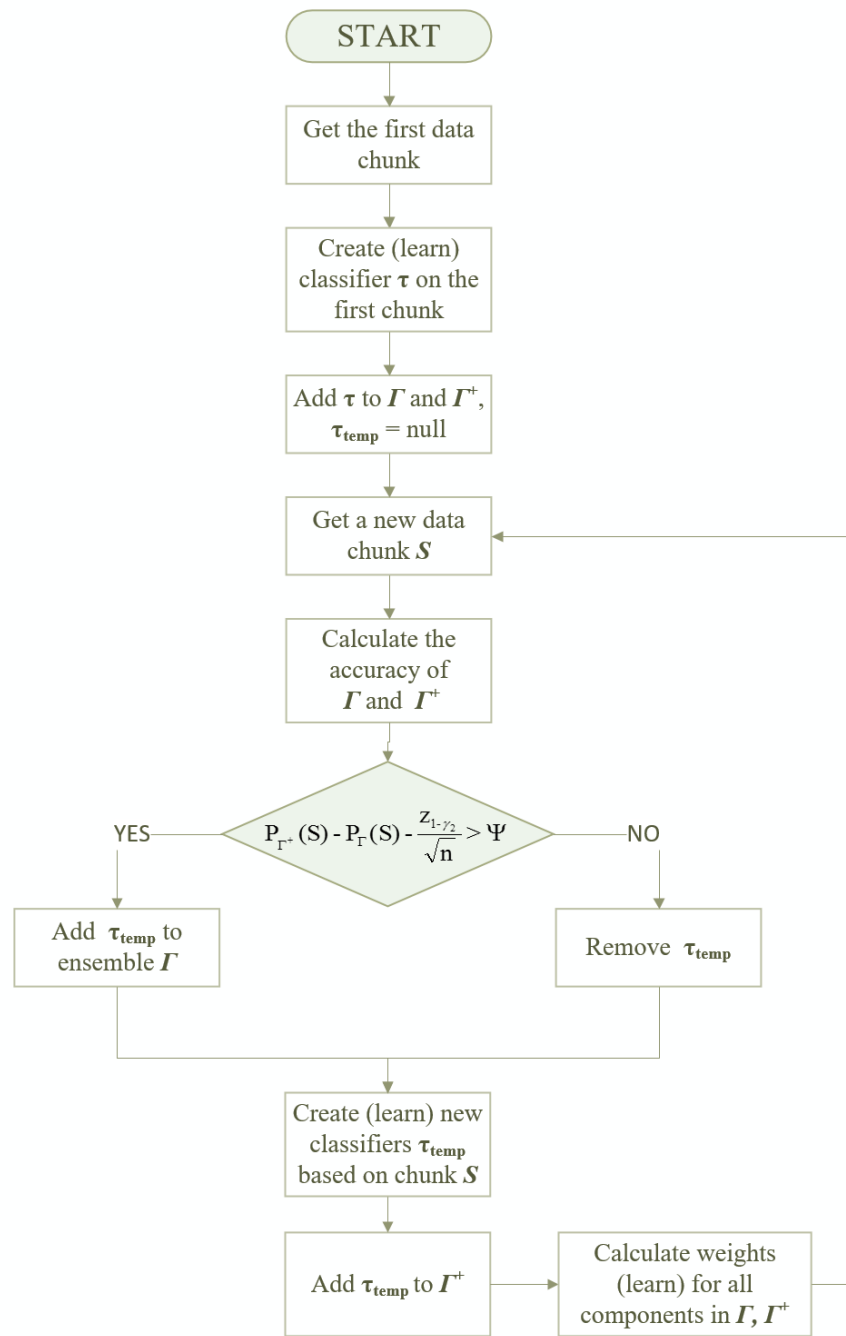
$$P_\Gamma(S_n) = \frac{\sum_{i=1}^n G_\Gamma(X_i)}{n}.$$

Then, if the condition

$$P_{\Gamma^+}(S) - P_\Gamma(S) - \frac{z_{1-\gamma_2}}{\sqrt{n}} > \Psi$$

is satisfied, the classifier τ_{temp} is added to the ensemble. However, if this condition is not fulfilled, then the investigated classifier is discarded. Next a new classifier is created based on the investigated data chunk and is labeled by τ_{temp} .

The Dynamically Expanded Ensemble Algorithm (DEEA)



Data stream mining

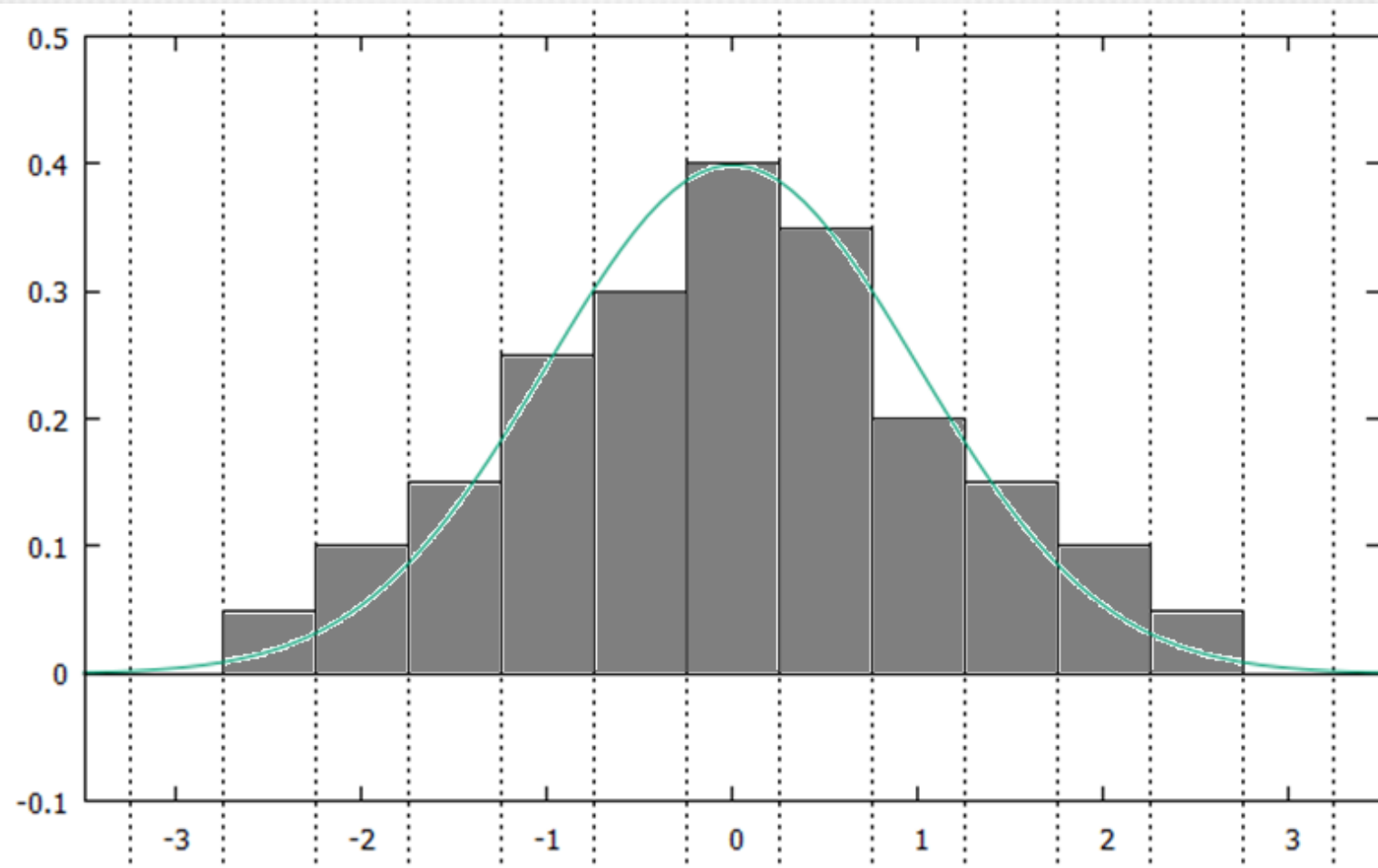
- content

- Data streams – introduction to the topic
- Concept drift
- Various strategies of learning
- How to deal with concept drift?
- Data stream classification methods – short overview
- Decision trees for data streams (including new results 2016)
- Ensemble methods for data streams (including new results 2016)
- **Probabilistic neural networks for stream data mining
(including new results 2016)**
- Final remarks and challenging problems
- References

Probabilistic Neural Networks (PNN) for stream data mining

PNN for stream data mining

Histograms



PNN for stream data mining

Non-recursive procedures

Let X_1, \dots, X_n be a sequence of independent, identically distributed random variables taking values in $A \subset \mathbb{R}^p$ and having a probability density function f . The general estimator of the probability density function f is given by the following formula

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n K_n(x, X_i)$$

where

$$K_n(x, X_i) = \frac{1}{h_n} K(x, X_i)$$

PNN for stream data mining

Non-recursive procedures

Parzen (1962):

$$\hat{f}_n(x) = \frac{1}{nh_n} \sum_{i=1}^n K\left(\frac{x - X_i}{h_n}\right)$$

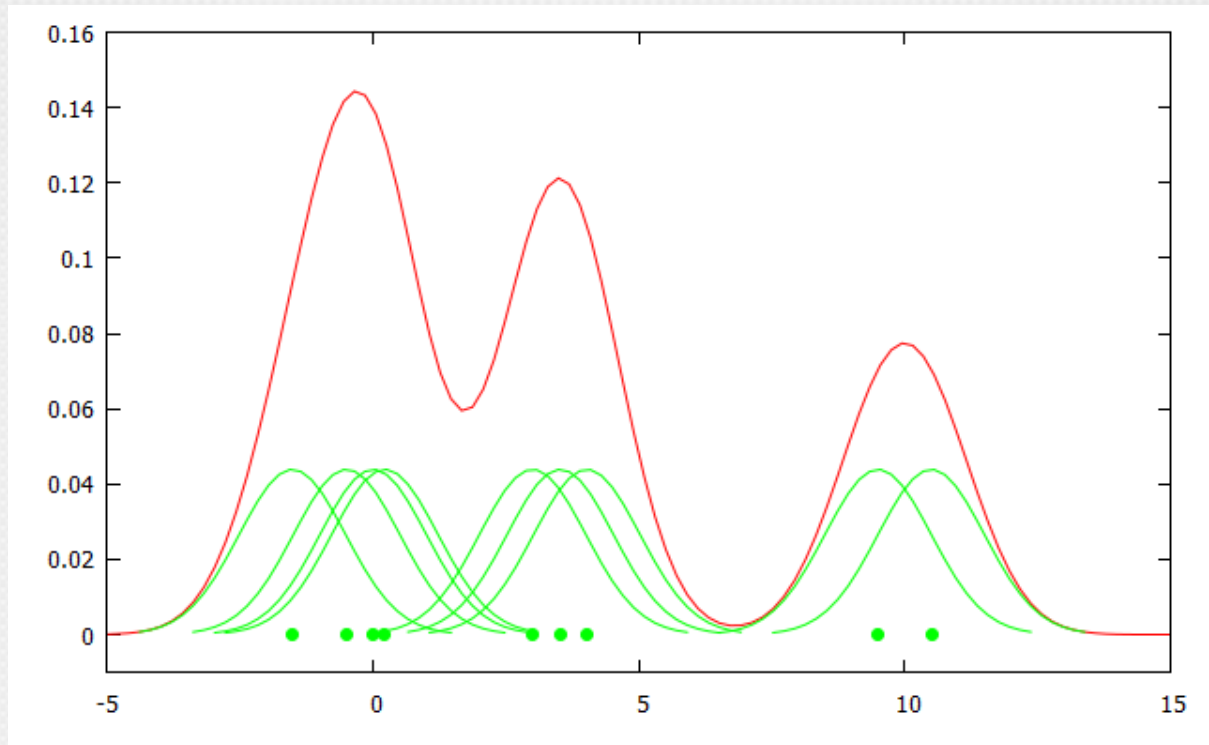
Cacullous (1965):

$$\hat{f}_n(x) = \frac{1}{nh_n^p} \sum_{i=1}^n K\left(\frac{x - X_i}{h_n}\right)$$

$$h(n) > 0, \quad h_n \xrightarrow{n} 0, \quad nh_n^p \xrightarrow{n} \infty$$

PNN for stream data mining

Non-recursive procedures

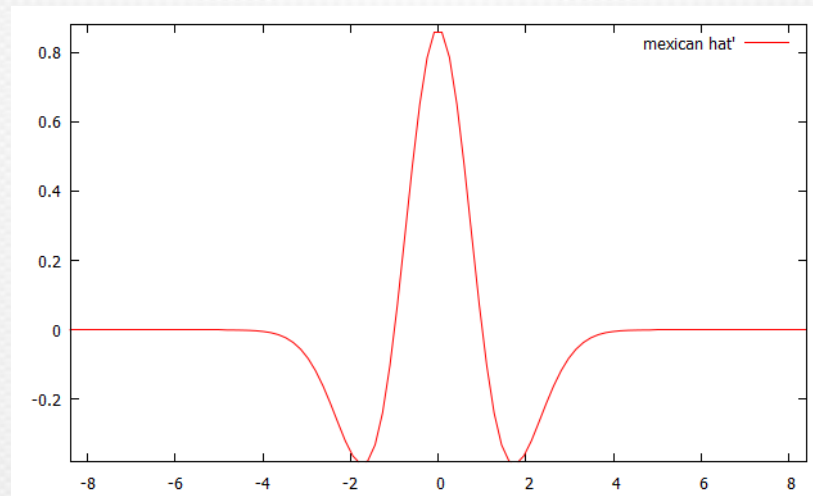
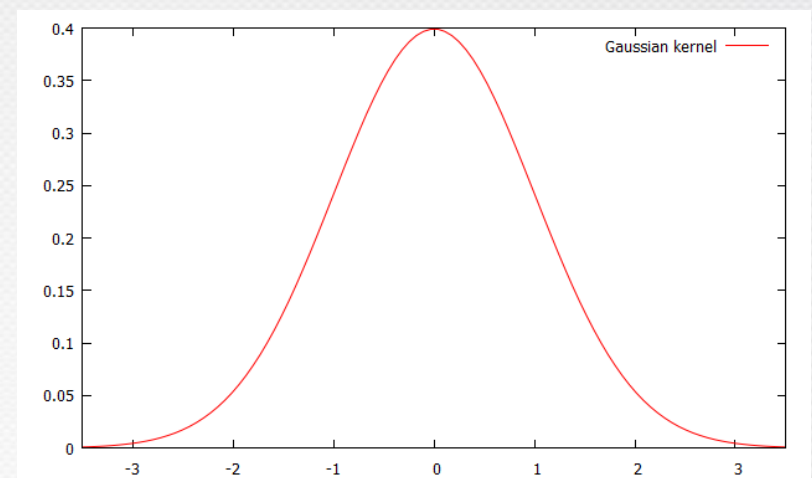
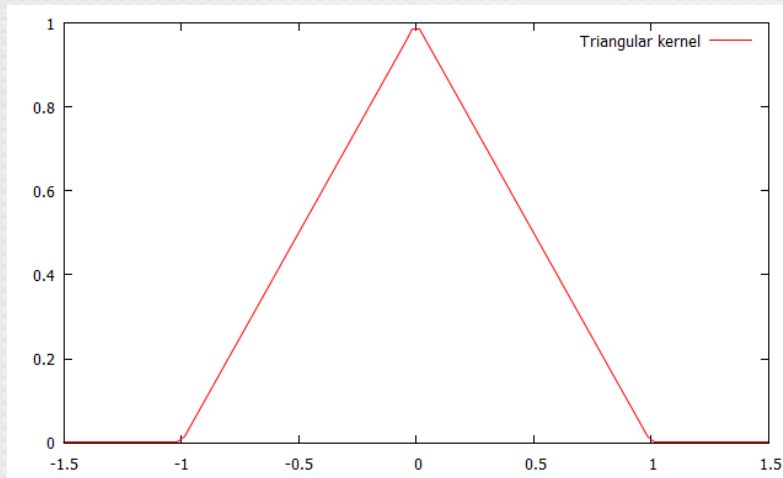


Parzen (1962):

$$\hat{f}_n(x) = \frac{1}{nh_n} \sum_{i=1}^n K\left(\frac{x - X_i}{h_n}\right)$$

PNN for stream data mining

Example of kernels



PNN for stream data mining

Non-recursive procedures

Example (Gaussian kernel)

$$K_n(x, u) = h_n^{-p} K\left(\frac{x - u}{h_n}\right)$$

$$K(x) = \prod_{j=1}^p H(x^{(j)})$$

$$K_n(x, u) = h_n^{-p} \prod_{j=1}^p H\left(\frac{x^{(j)} - u^{(j)}}{h_n}\right)$$

PNN for stream data mining

Non-recursive procedures

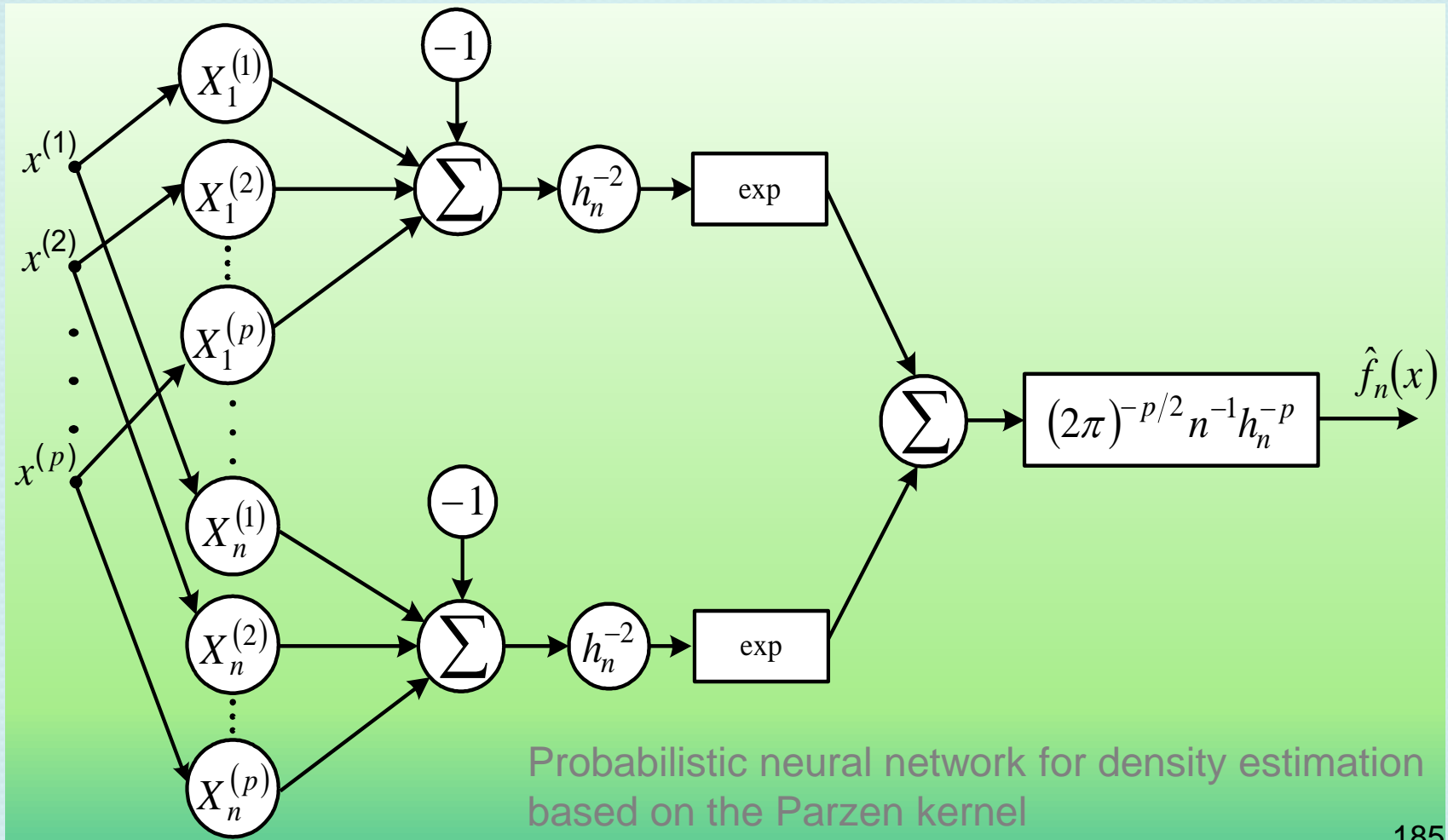
$$H(v) = (2\pi)^{-\frac{1}{2}} e^{-\frac{1}{2}v^2}$$

$$K_n(x, u) = h_n^{-p} (2\pi)^{-\frac{1}{2}} \prod_{j=1}^p \exp\left(-\frac{1}{2} \left(\frac{x^{(j)} - u^{(j)}}{h_n}\right)^2\right)$$

$$\hat{f}_n(x) = \frac{1}{(2\pi)^{\frac{1}{2}} n h_n^p} \sum_{i=1}^n \prod_{j=1}^p \exp\left(-\frac{1}{2} \left(\frac{x^{(j)} - X_i^{(j)}}{h_n}\right)^2\right)$$

PNN for stream data mining

Recursive procedures



PNN for stream data mining

Non-recursive procedures

Let us modify estimator as follows

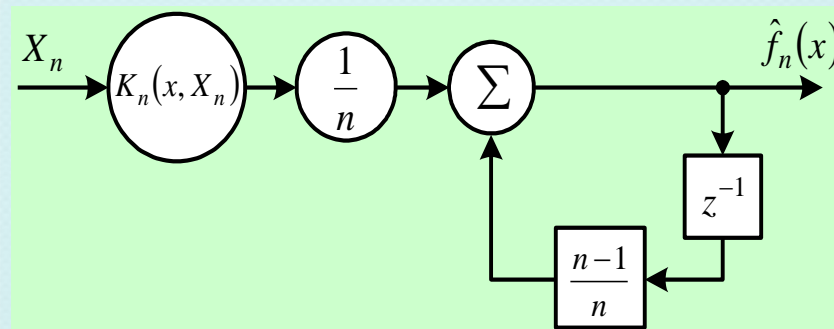
$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n K_n(x, X_i) \rightarrow \hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n K_i(x, X_i)$$

PNN for stream data mining

Recursive procedures

Observe that estimator is computationally equivalent to the recursive procedure

$$\hat{f}_{n+1}(x) = \hat{f}_n(x) + \frac{1}{n+1} [K_{n+1}(x, X_{n+1}) - \hat{f}_n(x)]$$
$$\hat{f}_0(x) = 0$$

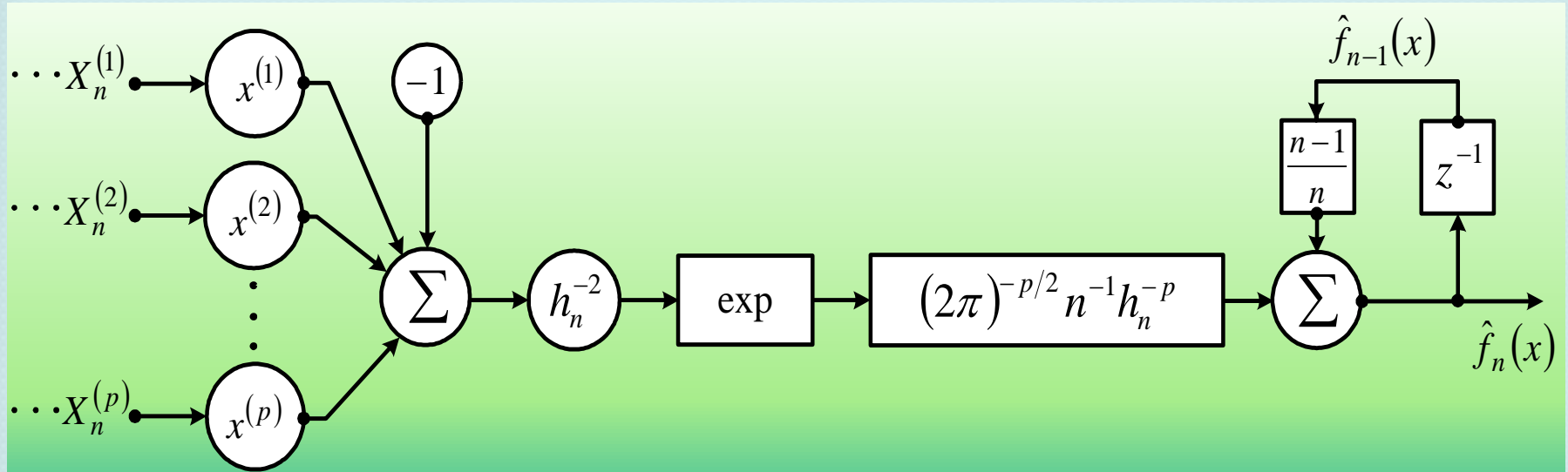


Recursive probabilistic neural network for density estimation

187

PNN for stream data mining

Recursive procedures



Recursive probabilistic neural network for estimation based on the Parzen kernel

PNN for stream data mining

for pattern classification

Let $(X, Y), (X_1, Y_1), \dots, (X_n, Y_n)$ be a sequence of i.i.d. pairs of random variables, Y takes values in the set of classes $S = \{1, \dots, M\}$, whereas x takes values in $A \subset \mathbb{R}^p$. The problem is to estimate Y from x and w_n , where $w_n = (X_1, Y_1), \dots, (X_n, Y_n)$ is a learning sequence. Suppose that p_m and f_m , $m = 1, \dots, M$ are the prior class probabilities and class conditional densities, respectively. We define a discriminant function of class j :

$$d_j(x) = p_j f_j(x)$$

PNN for stream data mining

for pattern classification

Let $L(i, j)$ be the loss incurred in taking action $i \in S$ when the class is j . We assume 0 – 1 loss function. For a decision function $\varphi: A \rightarrow S$ the expected loss is

$$R(\varphi) = \sum_{j=1}^M p_j \int_A L(\varphi(x), j) f_j(x) dx$$

PNN for stream data mining

for pattern classification

A decision function φ^* which classifies every x as coming from any class m for which

$$p_m f_m(x) = \max_j p_j f_j(x) = \max_j d_j(x)$$

is a Bayes decision function and

$$R^* = R(\varphi^*) = \sum_{j=1}^M p_j \int_A L(\varphi^*(x), j) f_j(x) dx$$

is the minimal Bayes risk. The function $d_m(x)$ is called the Bayes discriminant function.

PNN for stream data mining

for pattern classification

A decision function φ^* which classifies every $x \in A$ as coming from any class m for which

$$p_m f_m(x) = \max_j p_j f_j(x) = \max_j d_j(x)$$

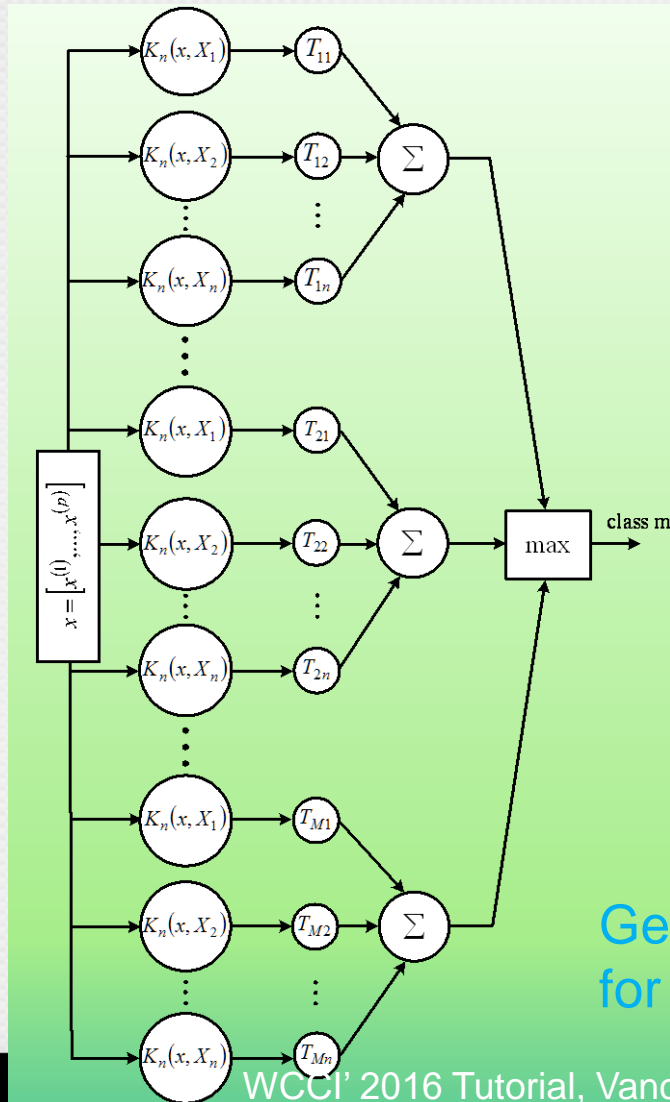
is a Bayes decision function and

$$R^* = R(\varphi^*) = \sum_{j=1}^M p_j \int_A L(\varphi^*(x), j) f_j(x) dx$$

is the minimal Bayes risk. The function $d_m(x)$ is called the Bayes discriminant function.

PNN for stream data mining

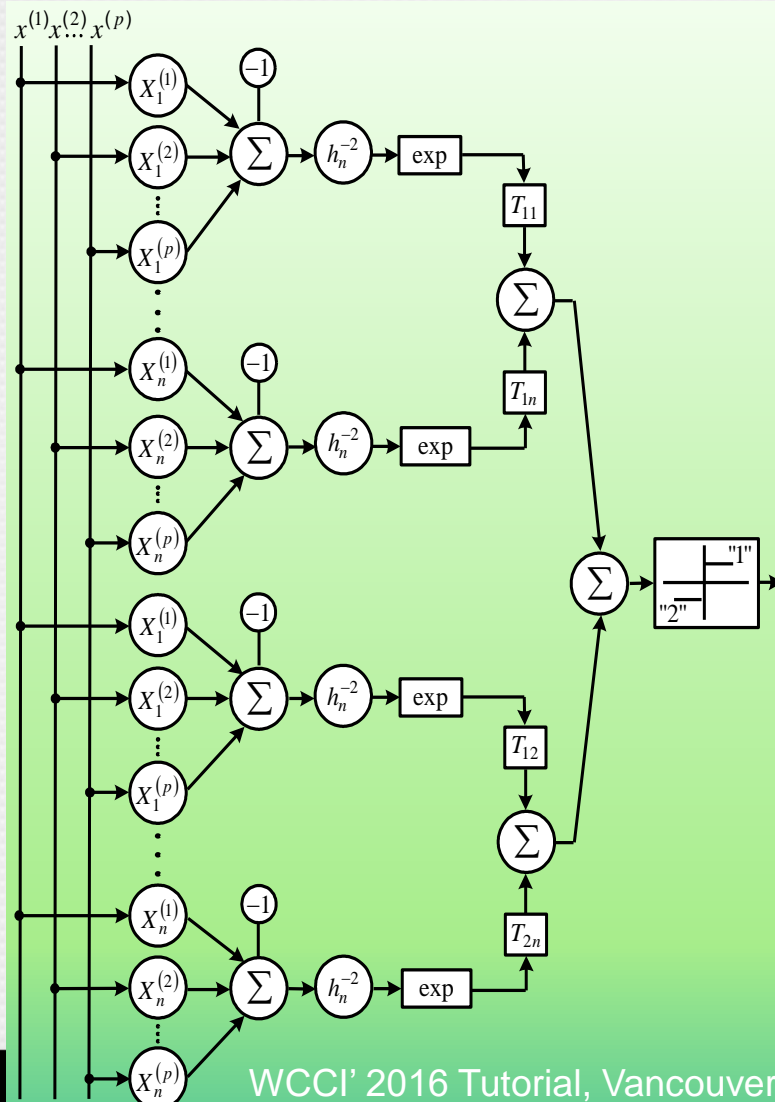
for pattern classification



Generalized regression neural network
for pattern classification

PNN for stream data mining

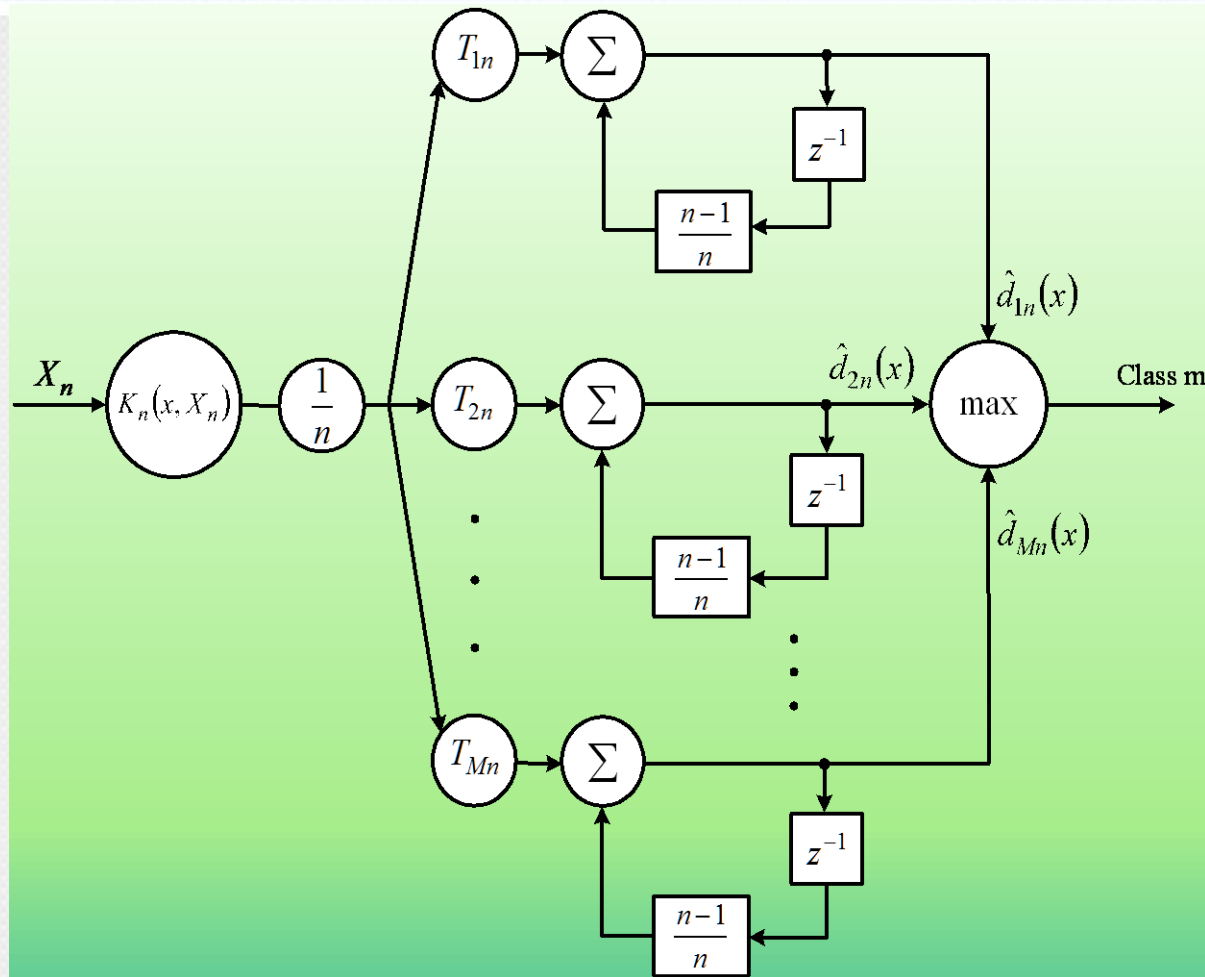
for pattern classification



Generalized regression neural network based on the Parzen kernel for pattern classification ($M=2$)

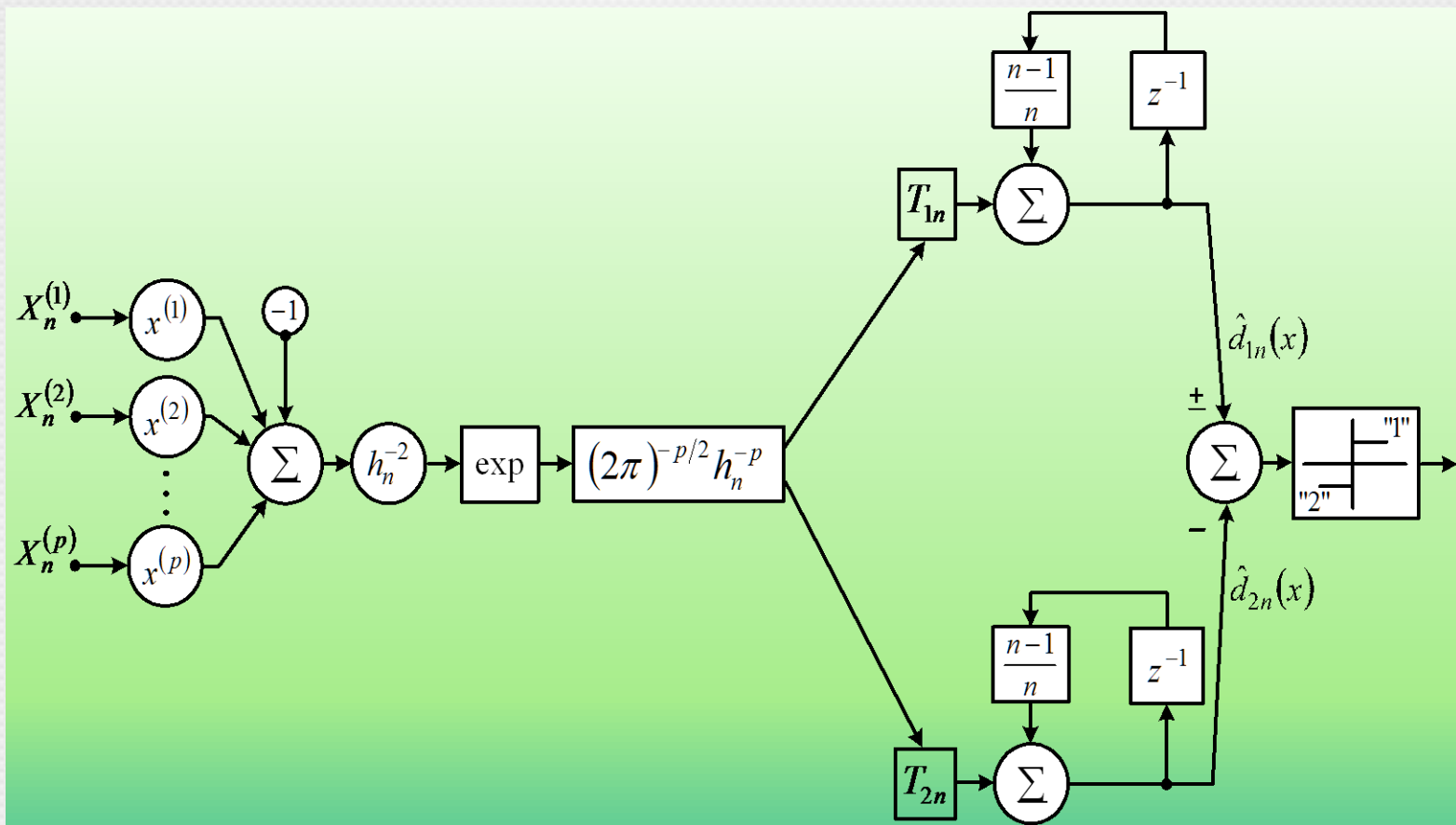
PNN for stream data mining

for pattern classification



PNN for stream data mining

for pattern classification



Recursive generalized regression neural network based on the Parzen kernel for pattern classification ($M=2$)

PNN for stream data mining

for pattern classification in a time-varying environment

K – number of classes

Concept Drift:

$f_{mn}(x)$ – time-varying probability density of
class m ($m = 1, \dots, K$)
at the instant n ($n = 1, 2, \dots$)

p_{mn} – time-varying a priori probabilities

PNN for stream data mining

for pattern classification in a time-varying environment

We want to estimate

$$d_{mn}(x) = p_{mn} f_{mn}(x), \quad m = 1, \dots, K, \quad n = 1, 2, \dots,$$

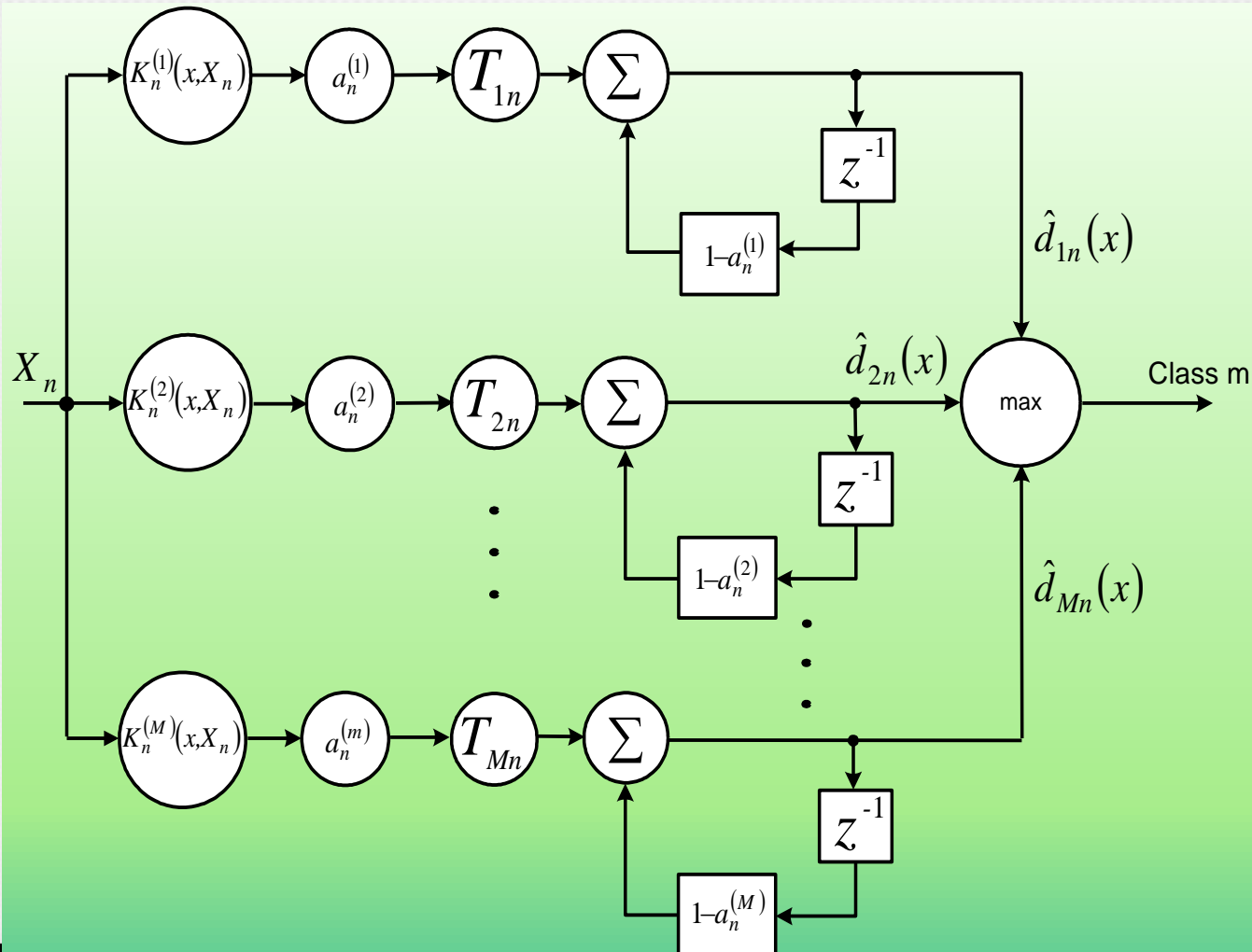
Find estimate \hat{d}_{mn} of d_{mn} such that

$$\left| \hat{d}_{mn}(x) - d_{mn}(x) \right| \xrightarrow{n} 0$$

in probability
(with probability one)

PNN for stream data mining

for pattern classification in a time-varying environment



PNN for stream data mining

for pattern classification in a time-varying environment

Example

Consider a two-category classification problem with $p_{1n} = p_{2n} = \frac{1}{2}$ and

$$f_{1n}(x) = f_1(x - n^t), f_{2n}(x) = f_2(x - n^t)$$

where

$$f_1(x) = \mathcal{N}(0,1), f_2(x) = \mathcal{N}(2,1)$$

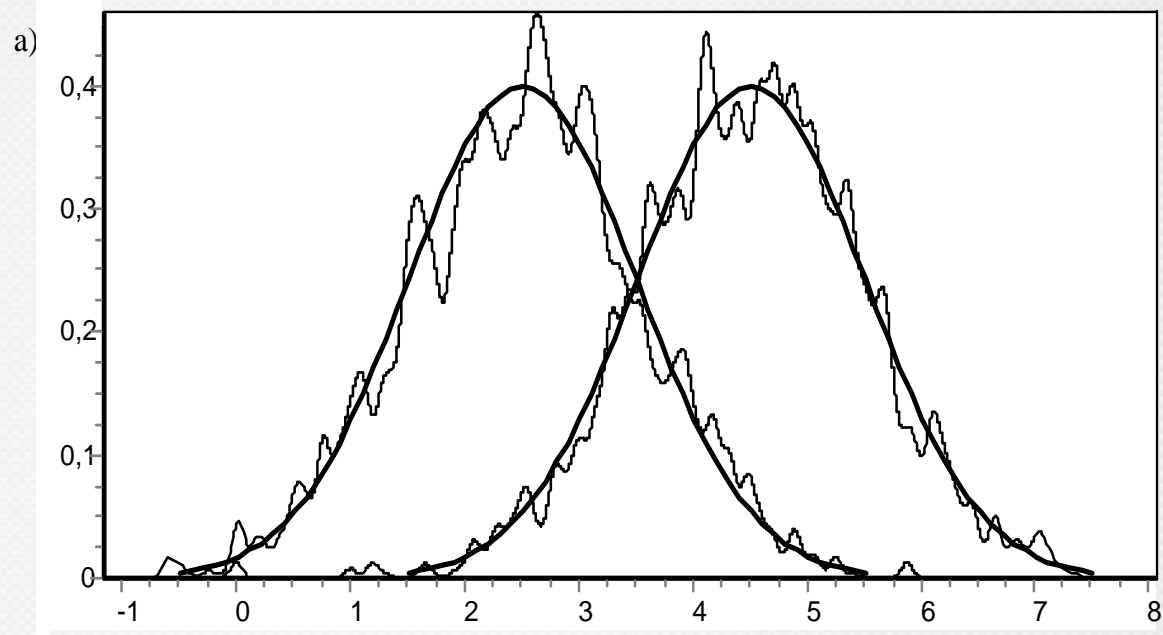
In this case the minimum probability of error is given by

$$P_e = \frac{1}{\sqrt{2\pi}} \int_1^{\infty} e^{-u^2/2} du = 0.159$$

PNN for stream data mining

for pattern classification in a time-varying environment

n	Error
1000	0.1559
2000	0.1602
3000	0.1573
4000	0.1691
5000	0.1566
6000	0.1598
7000	0.1562
8000	0.1607
9000	0.1561
10000	0.1572

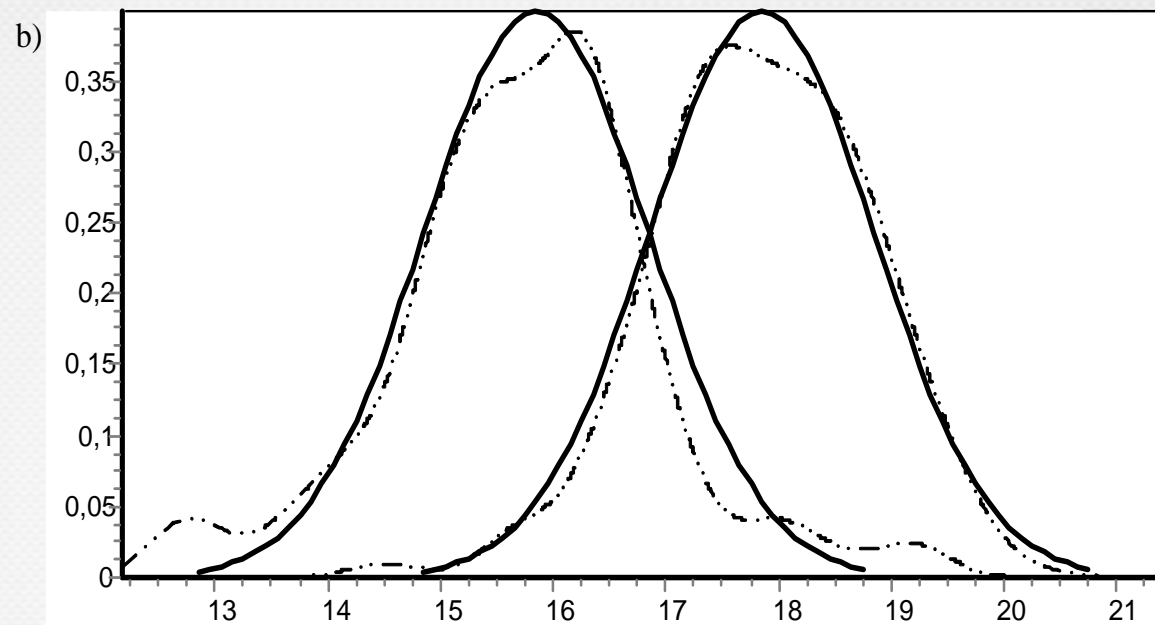


Empirical probability of misclassification
for $t=0.1$, $a=0.7$, $H=0.5$

PNN for stream data mining

for pattern classification in a time-varying environment

n	Error
1000	0.1639
2000	0.1648
3000	0.1687
4000	0.1649
5000	0.1564
6000	0.1611
7000	0.1558
8000	0.1586
9000	0.1569
10000	0.1562



Empirical probability of misclassification
for $t=0.3$, $a=0.4$, $H=0.3$

PNN for stream data mining

Generalized regression neural networks

Let (X, Y) be a pair of random variables X takes values in a Borel $A, A \subset \mathbb{R}^p$, whereas Y takes values in \mathbb{R} . Let f be the marginal Lebesgue density of X . Based on a sample $(X_1, Y_1), \dots, (X_n, Y_n)$ of independent observations of (X, Y) we wish to estimate the regression function

$$\phi(x) = E[Y|X = x]$$

PNN for stream data mining

Generalized regression neural networks

To estimate function ϕ we propose the following formula

$$\hat{\phi}_n(x) = \frac{\hat{R}_n(x)}{\hat{f}_n(x)}$$

where

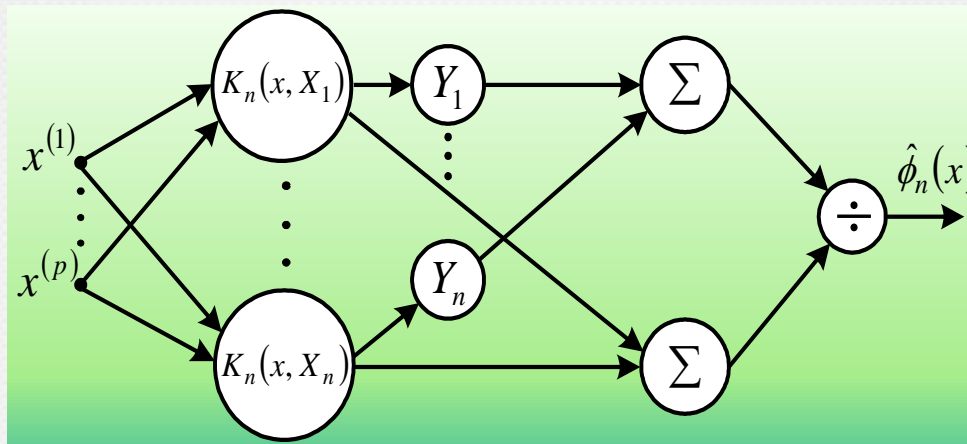
$$\hat{R}_n(x) = \frac{1}{n} \sum_{i=1}^n Y_i K_i(x, X_i)$$

and

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n K_i(x, X_i)$$

PNN for stream data mining

Generalized regression neural networks



Scheme of generalized regression neural network

PNN for stream data mining

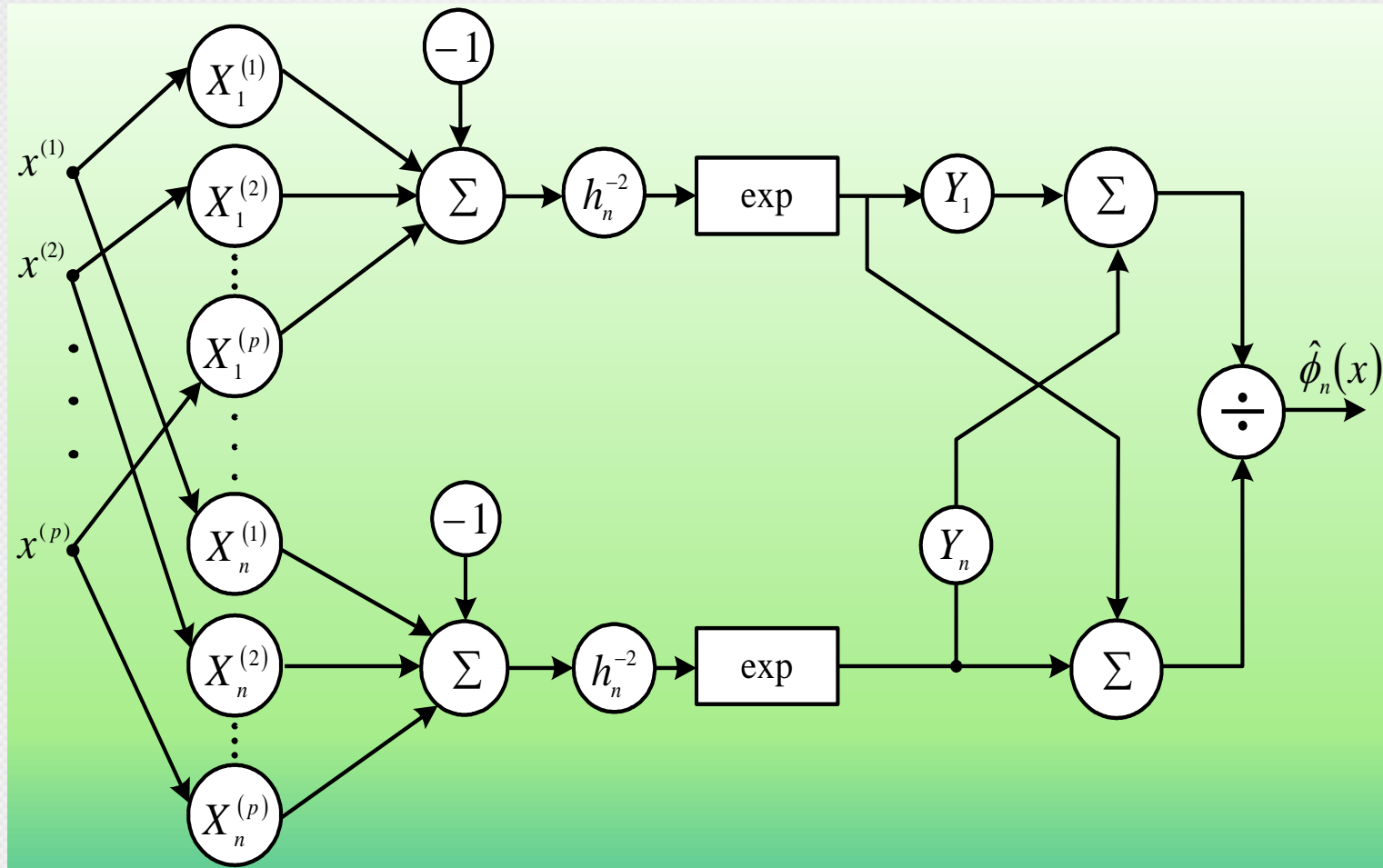
Generalized regression neural networks

Example (Nadaraya and Watson). Applying the Parzen kernel we get

$$\hat{\phi}_n(x) = \frac{\sum_{i=1}^n Y_i K\left(\frac{x - X_i}{h_n}\right)}{\sum_{i=1}^n K\left(\frac{x - X_i}{h_n}\right)}$$

PNN for stream data mining

Generalized regression neural networks



Generalized regression neural network based on the Parzen kernel

PNN for stream data mining

Generalized regression neural networks

The recursive version of procedure is given as follows

$$\hat{\phi}_n(x) = \frac{\hat{R}_n(x)}{\hat{f}_n(x)}$$

where

$$\hat{R}_n(x) = \frac{1}{n} \sum_{i=1}^n Y_i K_i(x, X_i)$$

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n K_i(x, X_i)$$

PNN for stream data mining

Generalized regression neural networks

$$\hat{R}_{n+1}(x) = \hat{R}_n(x) + \frac{1}{n+1} \left[Y_{n+1} K_{n+1}(x, X_{n+1}) - \hat{R}_n(x) \right]$$

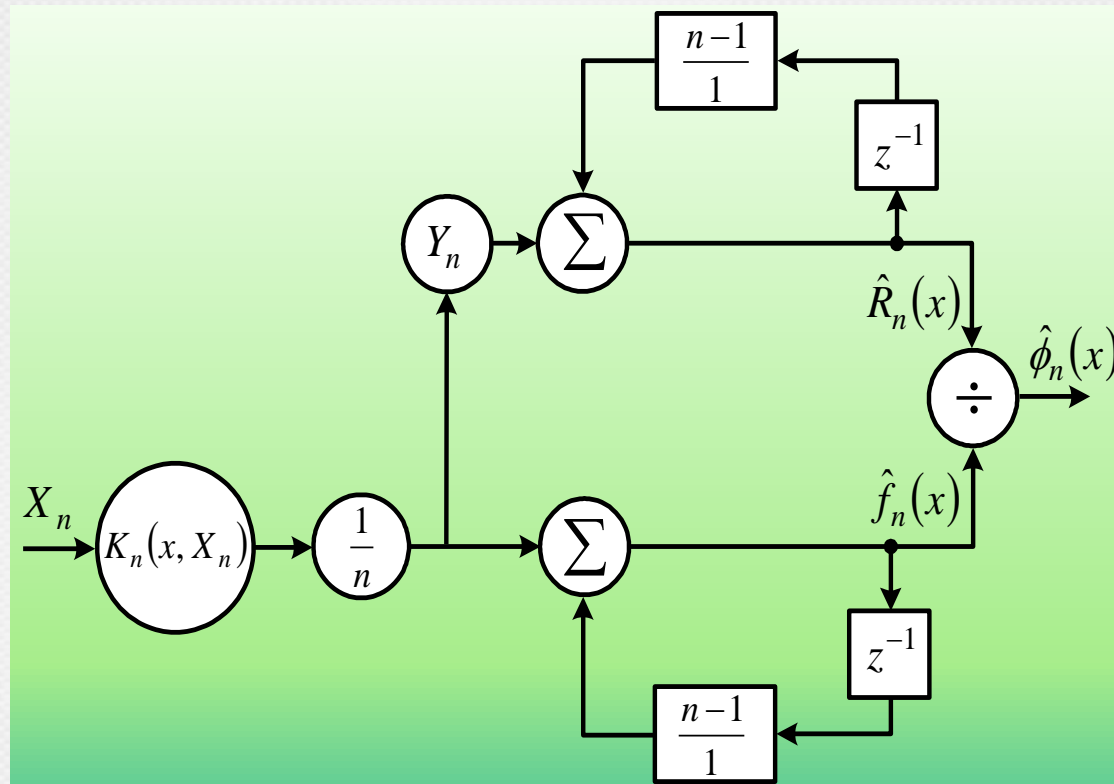
$$\hat{f}_{n+1}(x) = \hat{f}_n(x) + \frac{1}{n+1} \left[K_{n+1}(x, X_{n+1}) - \hat{f}_n(x) \right]$$

$$\hat{f}_0(x) = 0$$

where $\hat{R}_0(x) = 0$ and $\hat{f}_0(x) = 0$

PNN for stream data mining

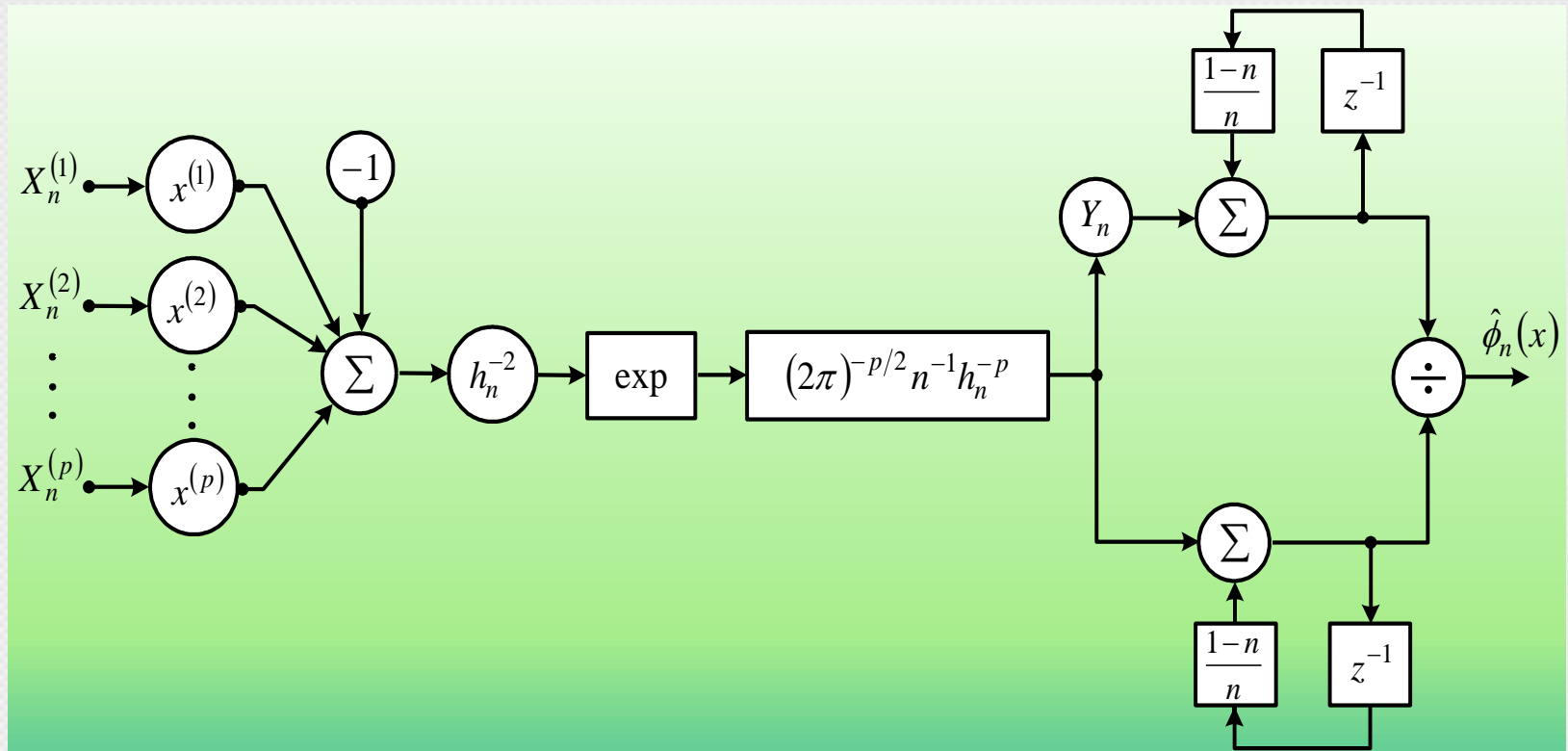
Generalized regression neural networks



Recursive generalized regression neural network

PNN for stream data mining

Generalized regression neural networks



Recursive generalized regression neural network
based on the Parzen kernel

PNN for stream data mining

Generalized regression neural networks in time-varying environment

In the non-stationary regression we consider a sequence of random variables $(X_n, Y_n), n = 1, 2, \dots$, having time-varying cumulative probability density functions $f_n(x, y)$. The problem is to find a measurable function $\phi_n: \mathbb{R}^p \rightarrow \mathbb{R}$ such that the L_2 risk

$$E[\phi_n(X_n) - Y_n]^2$$

attains minimum. The solution is the regression function

$$\phi_n^*(x) = E[Y_n | X_n = x], \quad n = 1, 2, \dots,$$

changing with time.

PNN for stream data mining

Generalized regression neural networks in time-varying environment

For the illustration of the capability of our GRNN we may consider an application to modelling of non-stationary plants described by

$$Y_n = \phi_n^*(X_n) + Z_n$$

where concept drift ϕ_n^* is given by:

- (i) $\phi_n^*(x) = \alpha_n \phi(x)$
- (ii) $\phi_n^*(x) = \phi(x) + \beta_n$
- (iii) $\phi_n^*(x) = \phi(x \omega_n)$
- (iv) $\phi_n^*(x) = \phi(x - \lambda_n)$

PNN for stream data mining

Generalized regression neural networks in time-varying environment

The problem of non-parametric regression boils down to finding an adaptive algorithm that could follow the changes of optimal characteristics expressed by formula. This algorithm should be constructed on the basis of a learning sequence, i.e. observations of the following random variables $(X_1, Y_1), (X_2, Y_2), \dots$

We assume that pairs of the above random variables are independent.

In points x , where $f(x) \neq 0$, the characteristics of the best model can be expressed as

$$\phi_n^*(x) = R_n(x) / f(x), \quad n = 1, 2, \dots,$$

where $R_n(x) = \phi_n^*(x)f(x)$.

PNN for stream data mining

Generalized regression neural networks in time-varying environment

The algorithm has the form

$$\hat{\phi}_n(x) = \hat{R}_n(x) / \hat{f}_n(x)$$

$$\hat{R}_{n+1}(x) = \hat{R}_n(x) + a_{n+1} [Y_{n+1} K_{n+1}(x, X_{n+1}) - \hat{R}_n(x)]$$

$$\hat{f}_{n+1}(x) = \hat{f}_n(x) + \frac{1}{n+1} (K_{n+1}(x, X_{n+1}) - \hat{f}_n(x))$$

where $\hat{R}_0^*(x)=0$ and $\hat{f}_0^*(x)=0$

PNN for stream data mining

Generalized regression neural networks in time-varying environment

Let

$$a_n > 0, \quad a_n \xrightarrow{n} 0, \quad \sum_{n=1}^{\infty} a_n = \infty$$

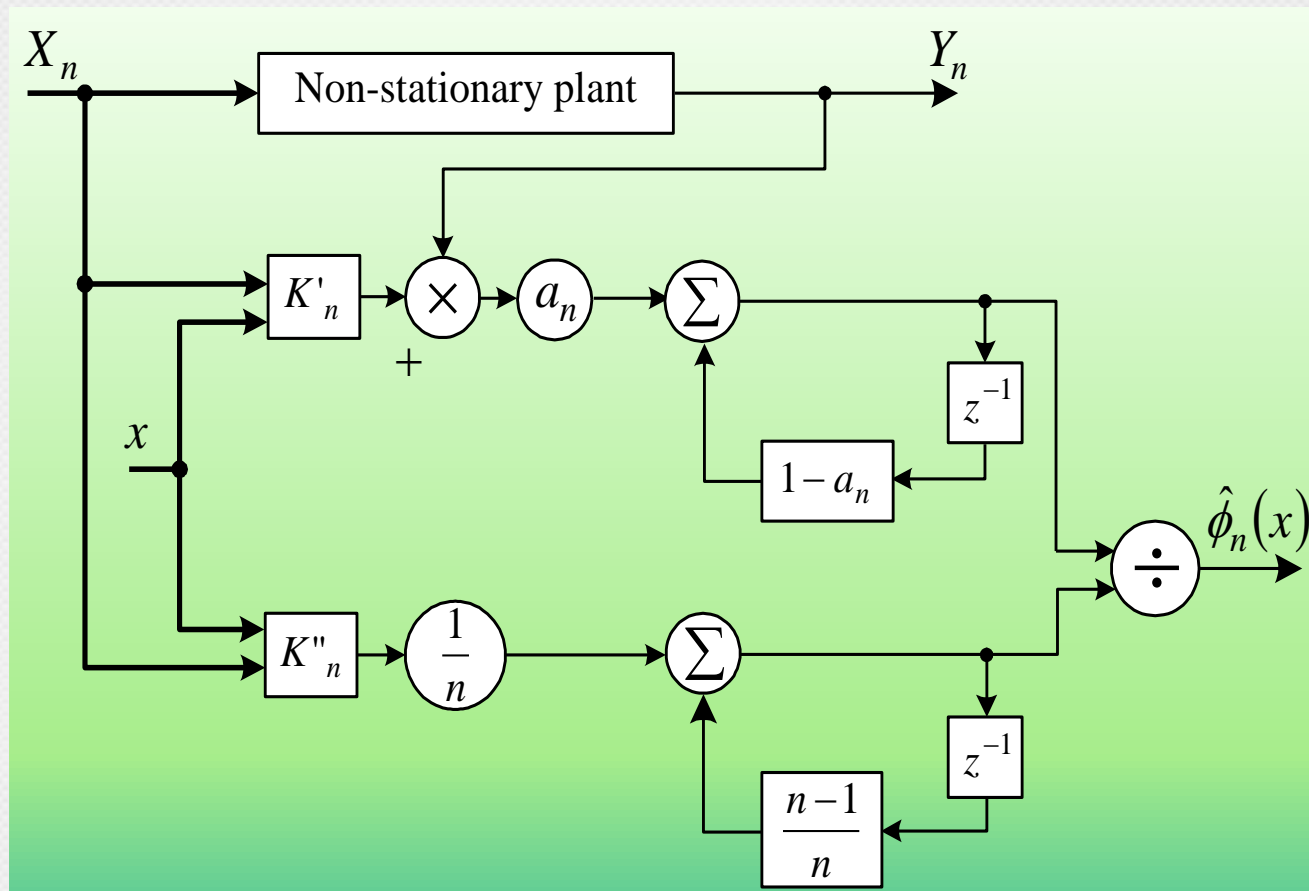
Under specific conditions imposed on parameters a_n and h_n and the rate of change of function ϕ_n the following holds

$$\left| \hat{\phi}_n(x) - \phi_n(x) \right| \xrightarrow{n} 0$$

in probability
(with probability one)

PNN for stream data mining

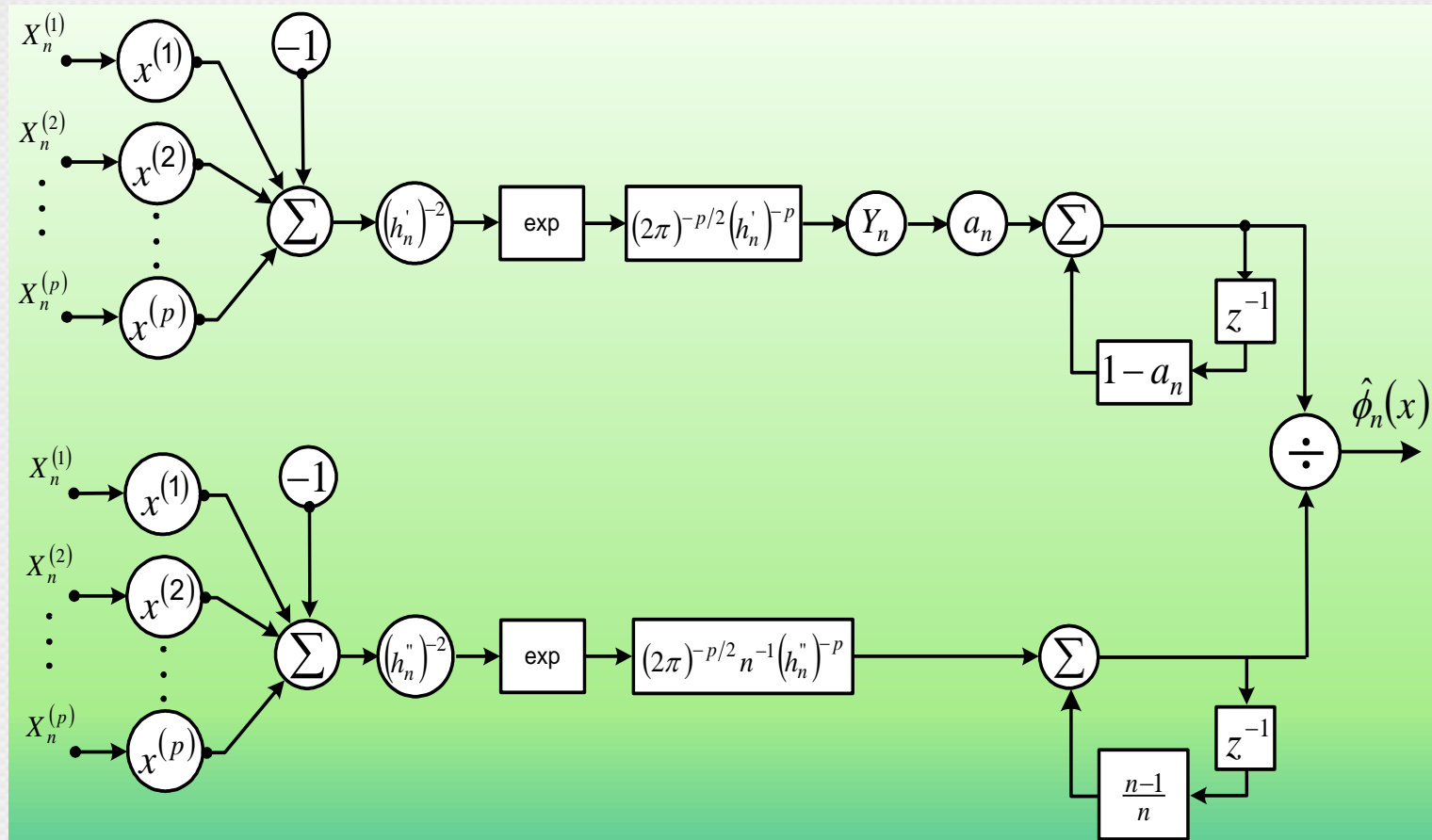
Generalized regression neural networks in time-varying environment



Block diagram of the GRNN applied to modelling of non-stationary plant 217

PNN for stream data mining

Generalized regression neural networks in time-varying environment



Block diagram of the GRNN based on the Parzen kernel

PNN for stream data mining

Generalized regression neural networks in time-varying environment

Example

In order to track changes of the system described by

$$y_n = (c_1 n^t + c_2 \log n + c_3) \phi(x_n) + z_n$$

where $t > 0$ is an unknown parameter and ϕ is an unknown function, it is possible to use algorithm if

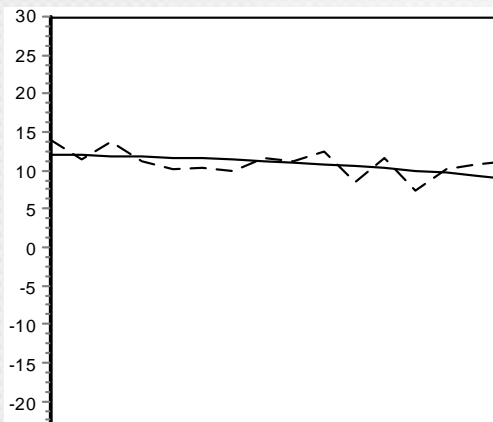
- a) $0 < t < \frac{1}{3}$ for weak convergence
- b) $0 < t < \frac{1}{6}$ for strong convergence

PNN for stream data mining

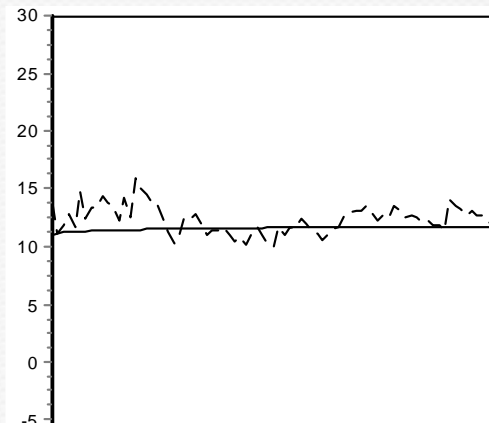
Generalized regression neural networks in time-varying environment

Example

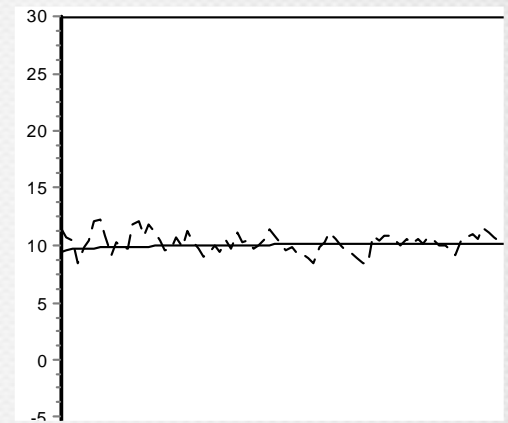
$$\phi_n^*(x_n) = 10 \cos(x_n) + n^{0.1} + z_n$$



a) $n=1000$



b) $x=0.2$



c) $x=0.6$

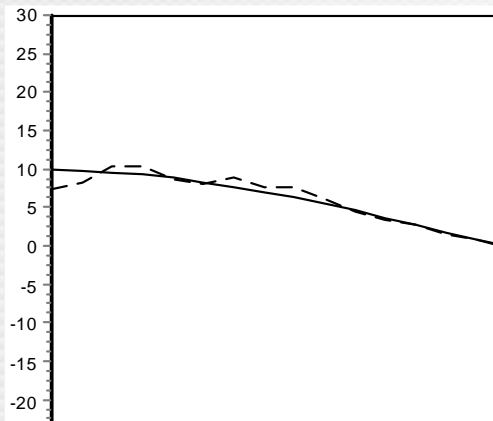
GRNN for modeling regressions with additive non-stationarity

PNN for stream data mining

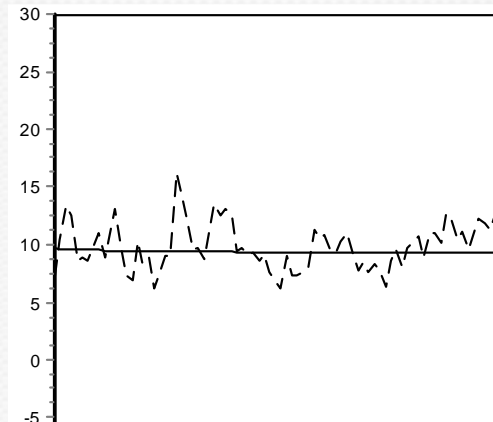
Generalized regression neural networks in time-varying environment

Example

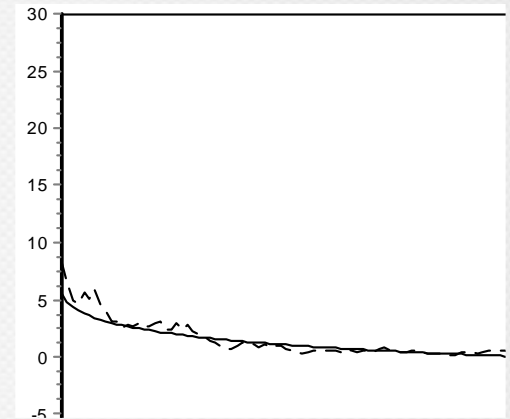
$$\phi_n^*(x_n) = 10 \cos(x_n n^{0.1}) + z_n$$



a) $n=1000$



b) $x=0.2$



c) $x=0.8$

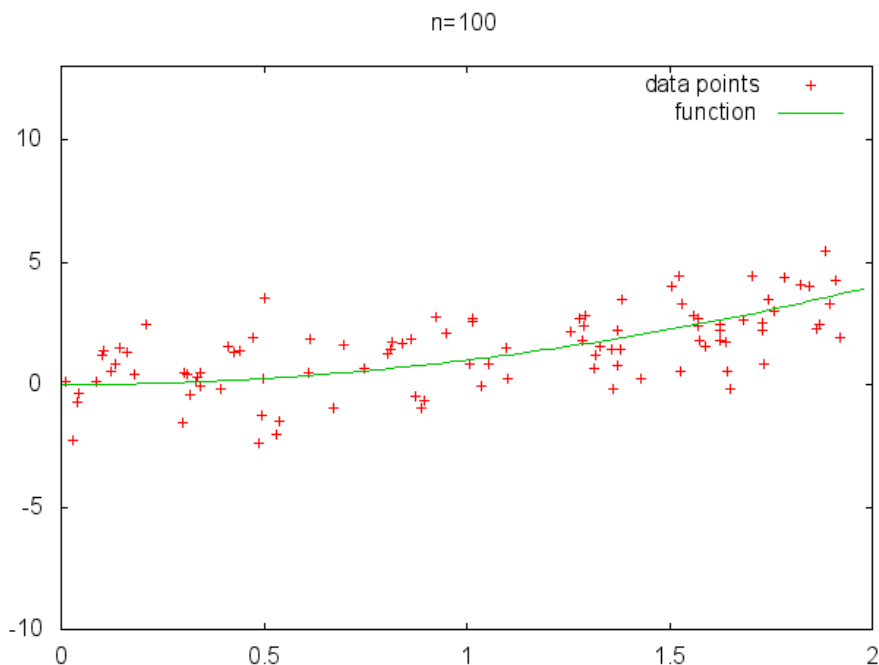
GRNN for modeling regressions with non-stationarity of the type “scale change”

New result 2016

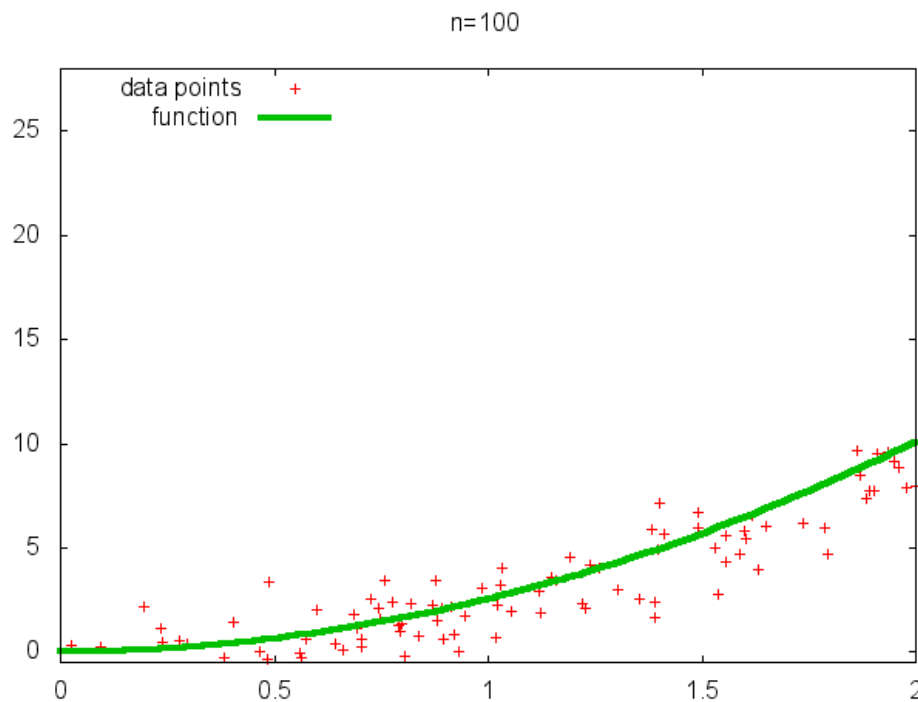
Three types of non-stationarities

$$1) Y_n = \phi(X_n) + Z_n$$

where variance of Z_n is changing over time



$$2) Y_n = \phi_n(X_n) + Z_n$$



$$3) Y_n = \phi_n(X_n) + Z_n \text{ with changing variance of } Z_n$$

New result 2016

Forgetting mechanism

$$\Phi_n^\lambda(x) = \frac{\sum_{i=1}^n \lambda^{n-i} Y_i \frac{1}{h_i^\lambda} K\left(\frac{x - X_i}{h_i^\lambda}\right)}{\sum_{i=1}^n \lambda^{n-i} \frac{1}{h_i'^\lambda} K\left(\frac{x - X_i}{h_i'^\lambda}\right)}$$

New result 2016

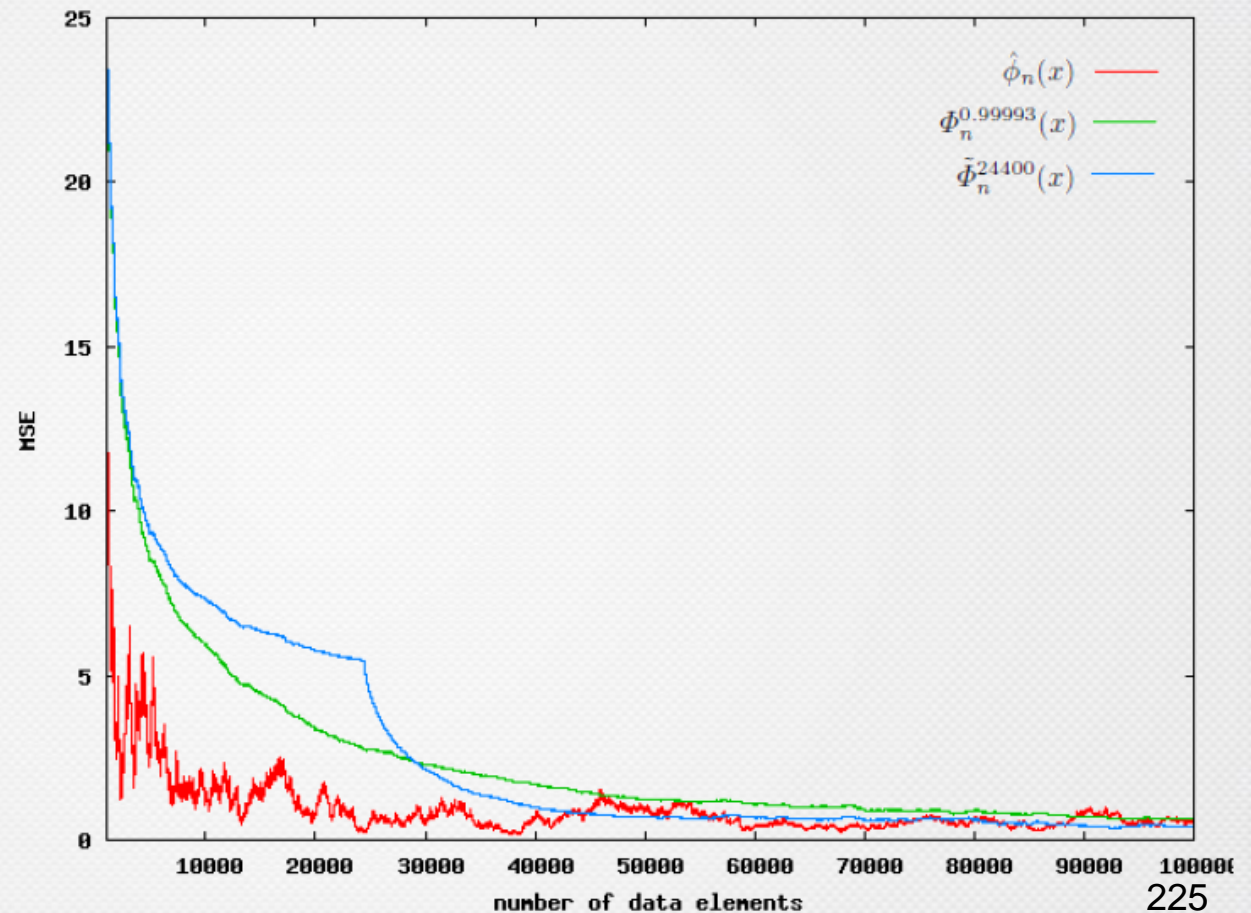
Sliding window

$$\tilde{\Phi}_n^W(x) = \frac{\sum_{i=n-W+1}^n Y_i \frac{1}{h_i^W} K\left(\frac{x-X_i}{h_i^W}\right)}{\sum_{i=n-W+1}^n \frac{1}{h_i'^W} K\left(\frac{x-X_i}{h_i'^W}\right)}$$

New result 2016

Experimental results

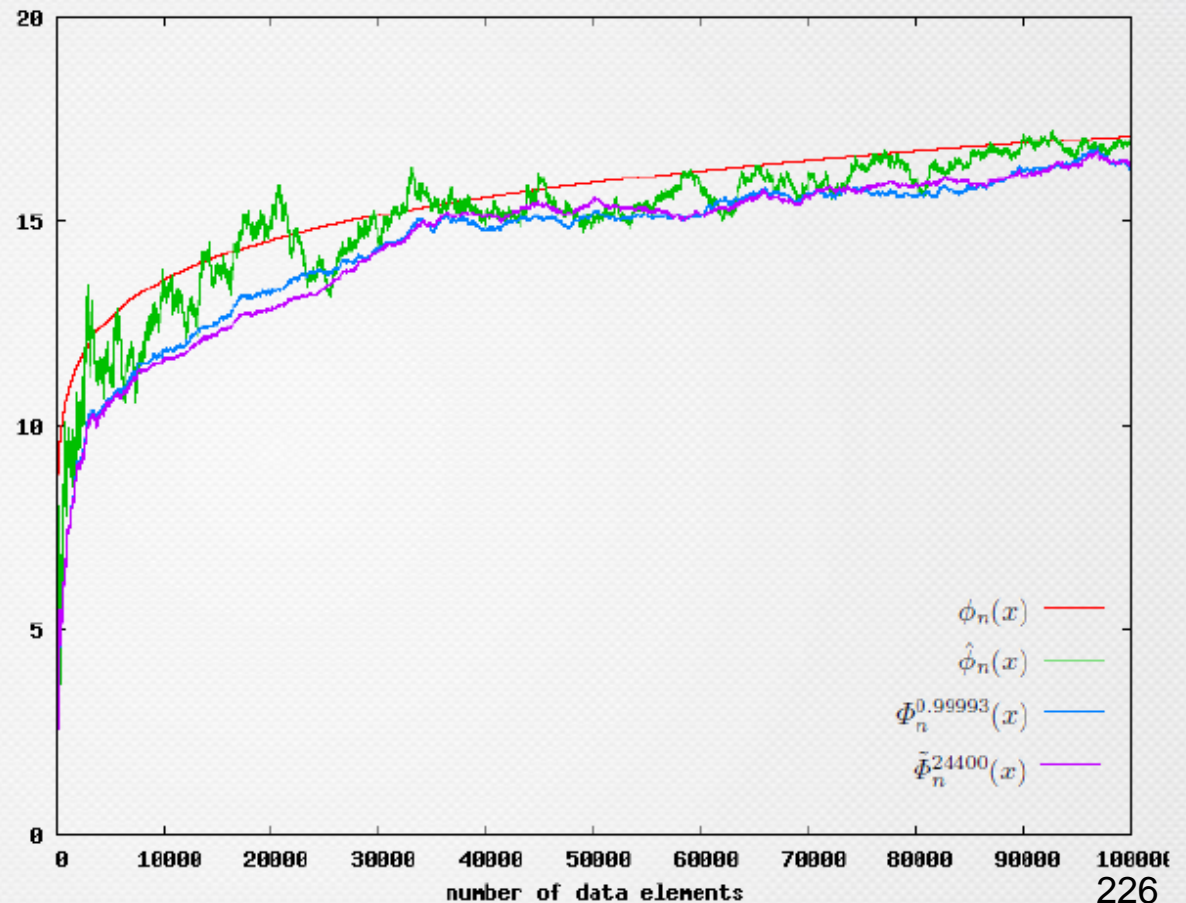
$$\phi_n(x) = 10 \cdot \cos(x) \cdot n^{0.1}$$



New result 2016

Experimental results

$$\phi_n(x) = 10 \cdot \cos(x) \cdot n^{0.1}$$



Data stream mining

- content

- Data streams – introduction to the topic
- Concept drift
- Various strategies of learning
- How to deal with concept drift?
- Data stream classification methods – short overview
- Decision trees for data streams (including new results 2016)
- Ensemble methods for data streams (including new results 2016)
- Probabilistic neural networks for stream data mining (including new results 2016)
- **Final remarks and challenging problems**
- References

Data stream mining

- Final remarks and challenging problems

To be presented at the end of the tutorial.

Data stream mining

- content

- Data streams – introduction to the topic
- Concept drift
- Various strategies of learning
- How to deal with concept drift?
- Data stream classification methods – short overview
- Decision trees for data streams (including new results 2016)
- Ensemble methods for data streams (including new results 2016)
- Probabilistic neural networks for stream data mining (including new results 2016)
- Final remarks and challenging problems
- **References**

References

Stream data mining algorithms (concept drift)

Rutkowski L., Duda P., Jaworski M., Pietruczuk L., Stream Data Mining: Algorithms and Their Probabilistic Properties, Studies in Big Data, Springer, 2017.

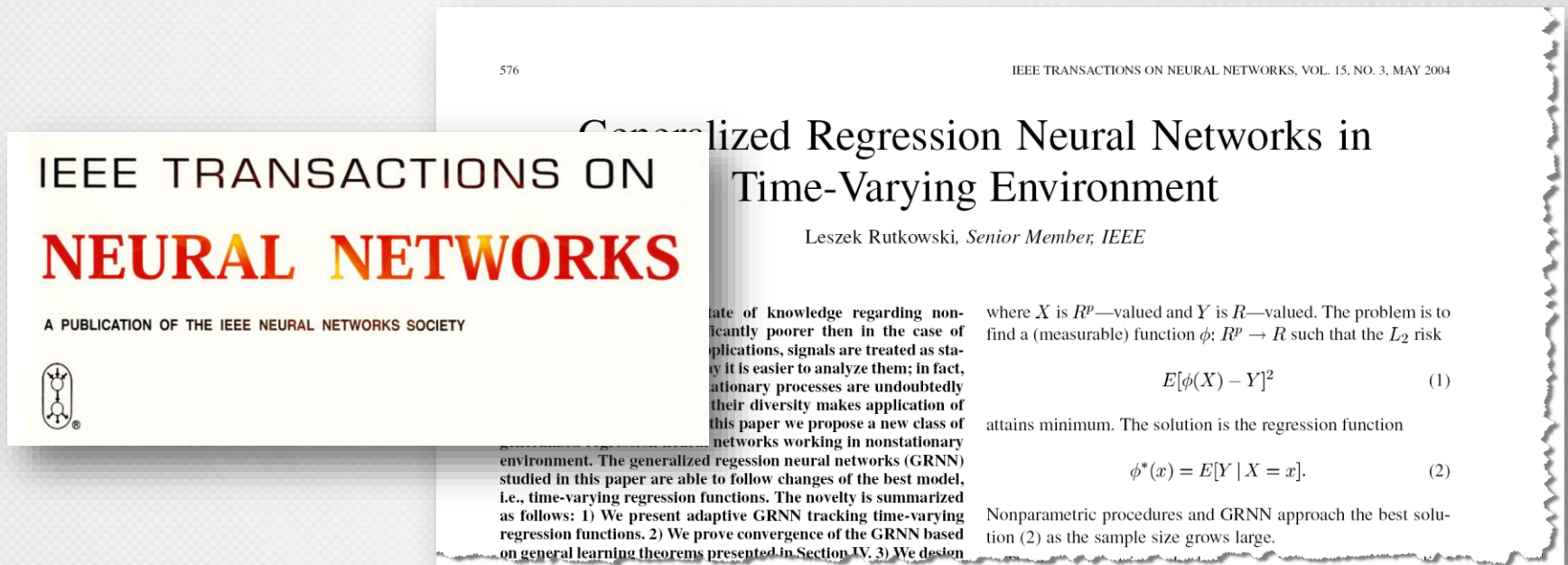
This book shows methods and algorithms which are mathematically justified.

- (1) It shows how to adopt the static decision trees, like ID.3 or CART, to deal with data streams.
- (2) A new technique, based on the McDiarmid bound, is developed.
- (3) New decision trees are designed, by a proper combination of the Gini index and misclassification error impurity measures, leading to the original concept of the hybrid decision trees.
- (4) The problem of designing ensembles and automatic choosing their sizes is described and solved.
- (5) Nonparametric techniques based on the Parzen – kernels and orthogonal series, are adopted to deal with concept drift in the problem of non-stationary regressions and classification in time-varying environment. Nonparametric procedures are developed and their probabilistic properties are investigated.

References

Stream data mining algorithms (concept drift)

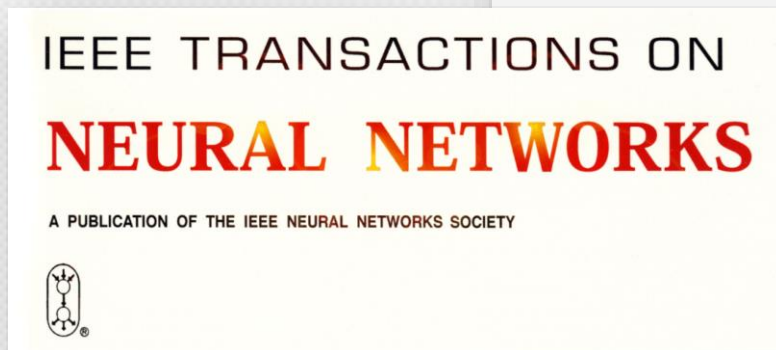
Rutkowski L., Generalized regression neural networks in time-varying environment, IEEE Transactions on Neural Networks, vol. 15, pp. 576-596, 2004.



References

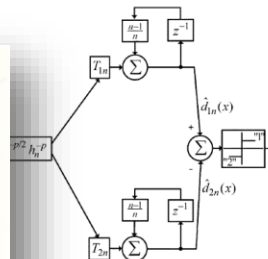
Stream data mining algorithms (concept drift)

Rutkowski L., Adaptive probabilistic neural-networks for pattern classification in time-varying environment, IEEE Transactions on Neural Networks, vol. 15, 2004.



816

IEEE TRANSACTIONS ON NEURAL NETWORKS, VOL. 15, NO. 4, JULY 2004



Theorem 2: If the following conditions are satisfied:

$$\sum_{n=1}^{\infty} a_n^2 \int \text{var} [Y_n K_n(x, X_n)] dx < \infty \quad (37)$$

$$\sum_{n=1}^{\infty} a_n^{-1} \int (r_n(x) - R_n(x))^2 dx < \infty \quad (38)$$

$$\sum_{n=1}^{\infty} a_n^{-1} \int (R_{n+1}(x) - R_n(x))^2 dx < \infty \quad (39)$$

then

$$I_n \xrightarrow{n \rightarrow \infty} 0 \quad \text{with pr. 1.} \quad (40)$$

Theorem 3: If the following conditions are satisfied:

$$\int \text{var} [Y_n K_n(x, X_n)] dx = O(n^{-A}), \quad A > 0 \quad (41)$$

$$\int (R_{n+1}(x) - R_n(x))^2 dx = O(n^{-B}), \quad B > 0 \quad (42)$$

$$\int (R_n(x) - R(x))^2 dx = O(n^{-C}), \quad C > 0 \quad (43)$$

Let us define the following function:

$$R_n(x) \stackrel{\text{df}}{=} f_n(x) E[Y_n | X_n = x] \quad n = 1, 2, \dots \quad (28)$$

References

Stream data mining algorithms (concept drift)

Rutkowski L., Pietruczuk L., Duda P., Jaworski M.,
Decision Trees for Mining Data Streams Based on the
McDiarmid's Bound, IEEE Transactions on Knowledge and
Data Engineering, vol. 25, pp. 1272 – 1279, 2013.

IEEE TRANSACTIONS ON
**KNOWLEDGE AND
DATA ENGINEERING**
A publication of the IEEE Computer Society
Indexed in ISI

4

IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 25, NO. X, XXXXXX 2013

$$= -\frac{1}{N} \log_2 \frac{1}{N} - \frac{N-1}{N} \log_2 \frac{N-1}{N-1} \underbrace{\log_2 \frac{N-1}{N-1}}_0$$

$$\left[\log_2 \frac{N-1}{N} + \frac{1}{N} \log_2 \frac{1}{N} \right] \quad (24)$$

$$\log_2 \frac{N-1}{N} + \log_2 \frac{1}{N}.$$

$$|f(Z) - f(Z')| = \left| \frac{N-1}{N} \log_2 \frac{N}{N-1} + \frac{\log_2 N}{N} \right|. \quad (25)$$

It is easy to prove that

$$\log_2 \left(\frac{n+1}{n} \right) \leq \frac{\log_2 e}{n}. \quad (26)$$

In view of (26), for $N \geq 2$, we can bound (25) as follows:

$$\frac{1}{N} \leq \frac{\log_2 N}{N}$$

Corollary 1. Suppose a and b are attributes for which the values of information gain, calculated from the data sample Z , satisfy $\text{Gain}_a(Z) > \text{Gain}_b(Z)$. For any fixed δ and ϵ given by (29), if $f(Z) > \epsilon$, then with probability $1 - \delta$ attribute a is better to split than attribute b , according to whole data stream. Moreover, if a and b are attributes with the highest and the second highest values of information gain, then with probability $1 - \delta$, a is the best attribute to split according to the whole stream.

Proof. For this particular choice of attributes a and b the assumptions of Theorem 1 are satisfied. Therefore, inequality (30) holds and can be transformed to the form

$$P(E[f(Z)] \geq f(Z) - \epsilon) \geq 1 - \delta. \quad (31)$$

Notice that if $f(Z) > \epsilon$, then $E[f(Z)] > 0$ with probability $1 - \delta$. The inequality $E[f(Z)] > 0$ is equivalent to $E[\text{Gain}_a(Z)] > E[\text{Gain}_b(Z)]$ what means that the attribute a is better than attribute b to split, according to the whole data stream, with probability $1 - \delta$.

Let us consider now the case, where a and b are

References

Stream data mining algorithms (concept drift)

Rutkowski L. Jaworski M., Pietruczuk L., Duda P.,
The CART decision tree for mining data streams,
Information Sciences, vol. 266, pp. 1–15, 2014.



Algorithm 1: The dsCART

Inputs: \mathbf{S} is a sequence of examples,
 \mathfrak{A} is a set of discrete attributes, of the one class
 α is one minus the desired probability of choosing the correct
attribute at any given node,
 θ is the tie breaking parameter.

Output: *dsCART* is a decision tree.

Procedure **dsCART**($\mathbf{S}, \mathfrak{A}, \alpha$)

Let *dsCART* be a tree with a single leaf L_0 (the root).

Let $\mathfrak{A}_0 = \mathfrak{A}$

For each attribute $a^i \in \mathfrak{A}$

For each value a_{λ}^i of attribute a^i

For each class k

$n_{i,\lambda,0}^k = 0$

For each example s in \mathbf{S}

Sort s into a leaf L_q using the current tree.

For each attribute $a^i \in \mathfrak{A}_a$

References

Stream data mining algorithms (concept drift)

Rutkowski L., Jaworski M., Pietruczuk L., Duda P., Decision Trees for Mining Data Streams Based on the Gaussian Approximation, IEEE Transactions on Knowledge and Data Engineering, vol. 26 , no. 1, pp. 108–119, 2014.

IEEE TRANSACTIONS ON
**KNOWLEDGE AND
DATA ENGINEERING**

A publication of the IEEE Computer Society
Indexed in ISI

RUTKOWSKI ET AL.: DECISION TREES FOR MINING DATA STREAMS BASED ON THE GAUSSIAN APPROXIMATION

9

be emphasized that for Gini index (CART decision tree) McDiarmid's method gives much better result (see [29]).

GAUSSIAN DECISION TREE ALGORITHM

1 allows us to propose an algorithm called the *decision tree*. This algorithm is a modification of *ing tree algorithm* proposed in [7]. For clarity of the de (see Table 1) the following notation will be d:

is the information gain computed for the ribute a^i in the leaf L_q .

n_q^k is the number of elements from the k th class in the leaf L_q , for which the value of attribute a^i is equal to a_λ^i .

- n_q^k is the number of elements from the k th class in the leaf L_q .

At the beginning, the algorithm starts with a single leaf (the root) and initializes the input parameters (in particular the parameter *last*, denoting the index of recently created

TABLE 1
The GDT Algorithm

Inputs: S is a sequence of examples,
 \mathcal{A} is a set of discrete attributes,
 th is a minimal fraction of elements belonging to the one class,
 δ is one minus the desired probability of choosing the correct attribute at any given node.

Output: GDT is a decision tree.

Procedure $GDT(S, \mathcal{A}, th, \delta)$

01. Let GDT be a tree with a single leaf L_0 (the root).
02. Let $\mathcal{A}_0 = \mathcal{A}$
03. Let $C' = \frac{1-th}{th}$, $last = 0$
04. For each attribute $a^i \in \mathcal{A}$
05. For each value a_λ^i of attribute a^i
06. $n_{i,\lambda,0}^1 = 0$, $n_{i,\lambda,0}^2 = 0$.
07. For each example s in S

References

Stream data mining algorithms (concept drift)

Rutkowski L., Jaworski M., Pietruczuk L., Duda P.,
A new method for data stream mining based on the
misclassification error, IEEE Transaction on Neural
Networks and Learning Systems, vol.26, no. 5,
pp. 1048-1059, 2015

IEEE TRANSACTIONS ON
**NEURAL NETWORKS AND
LEARNING SYSTEMS**

IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS

1

A New Method for Data Stream Mining Based on the Misclassification Error

Leszek Rutkowski, Fellow, IEEE, Maciej Jaworski, Lena Pietruczuk, and Piotr Duda

... method for constructing
proposed. First a new spli-
sification error is derived.
the best attribute computed
available data sample is
ity, as the attribute derived
Next this result is combined
the Gini Index. It is shown
that such combination provides the highest accuracy among all
studied algorithms.

Index Terms—Classification, data stream, decision trees,
impurity measure, splitting criterion.

I. INTRODUCTION

A. Motivation and Results

In recent years, the amount of data that needs to be analyzed
is growing very fast. Potentially unlimited number of data is

One of the most important techniques used in data mining is
the data classification. Let us assume that the data elements are
described by D attributes. Each attribute can be either nominal
or numerical. If the i th attribute is nominal then the number of
possible values is finite and equal to n_i . Moreover, to each data
element, a class is assigned. The number of different classes
is denoted by K . The aim of the classification is to construct
a function called the classifier, based on the training data set
of elements. The classifier maps the set of values of attributes
into the set of classes. It is further used to classify unlabeled
data elements. There exists a wide variety of methods used for
data classification. The most popular are neural networks [25],
 k -nearest neighbors [26] and decision trees [27]–[29]. The
last one is the main subject of this paper. Decision tree
consists of nodes and leaves. Each node is split according to
some attribute into its children (nodes or leaves). Each child
corresponds to one value of the attribute (in case of nonbinary


The 16th International Conference on Artificial Intelligence and Soft Computing ICAISC 2017

ICAISC 2017 About Important dates Invited talks Program Scope Conference Details ▾ Contact


ICAISC Zakopane

International Conference on Artificial Intelligence and Soft Computing

June 11-15, 2017



1:10



About ICAISC 2017

The 16th International Conference on Artificial Intelligence and Soft Computing ICAISC 2017 will be held in Zakopane (situated in the High Tatra mountains), Poland on June 11-15, 2017 in [Mercure Zakopane Kasprowy Hotel](#). The conference will provide an excellent opportunity for scientists and engineers to present and discuss the latest scientific results and methods. The conference will include keynote addresses, contributed papers, and numerous lectures and tutorials on a wide range of topics.

The working language of the conference is English. Only original, unpublished papers in the aforementioned fields are invited. Authors should submit an electronic version of papers by the conference web page. The papers should be organized in accordance with a common scientific structure (abstract, state of the art in the

**Thank you for your
attention**