

Novelty Detection in Data Streams

Rita P. Ribeiro João Gama

LIAAD-INESC TEC, University of Porto, Portugal



Outline

Introduction

- Novelty

- Applications

One-Class Classification

- Problem Definition

- Common Approaches

- Case Study: Predict Train Door Failures

Novelty Detection

- Problem Definition

- Key Aspects

- MINAS algorithm

Evaluation Issues

Challenges and Future Work

Outline

Introduction

- Novelty
- Applications

One-Class Classification

- Problem Definition
- Common Approaches
- Case Study: Predict Train Door Failures

Novelty Detection

- Problem Definition
- Key Aspects
- MINAS algorithm

Evaluation Issues

Challenges and Future Work

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

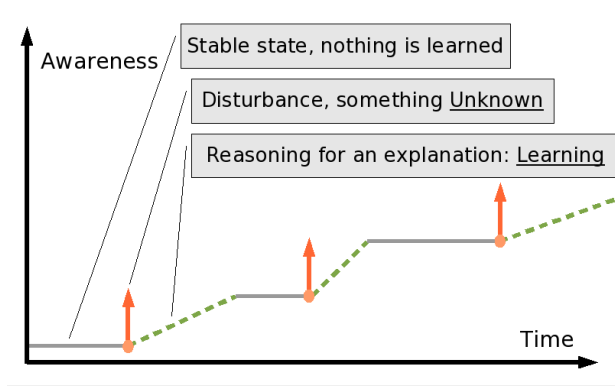
Key Aspects

MINAS algorithm

Evaluation Issues

Challenges and Future Work

How do we learn?



Novelty

- ▶ **Novelty** is a relative concept defined in the context of a representation of our current knowledge
- ▶ **Novelty Detection** refers to the automatic identification of unforeseen phenomena embed in a large amount of normal data
- ▶ Specially useful when novel concepts represent abnormal or unexpected conditions
 - ▶ expensive to obtain abnormal examples
 - ▶ probably impossible to simulate all possible abnormal conditions

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

MINAS algorithm

Evaluation Issues

Challenges and Future Work

Applications

- ▶ Intrusion Detection
- ▶ Fault Detection
- ▶ Fraud Detection
- ▶ Medical Diagnosis
- ▶ Spam Filter
- ▶ Text Classification

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

MINAS algorithm

Evaluation Issues

Challenges and Future Work

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

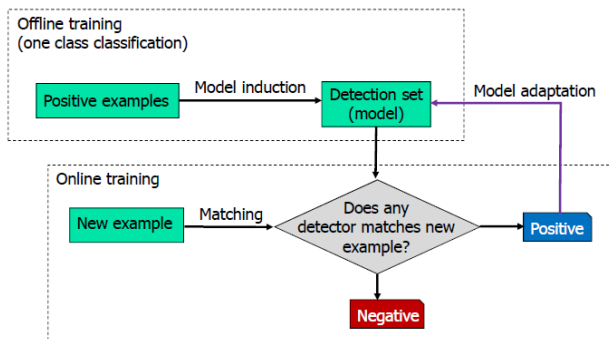
MINAS algorithm

Evaluation Issues

Challenges and Future Work

One-Class Classification: Problem Definition

- ▶ Offline Phase
 - ▶ Normal concept is composed by one class.
 - ▶ All training examples belong to the normal class.
- ▶ Online Phase
 - ▶ Examples not explained by the normal concept are labeled as abnormal.



Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

MINAS algorithm

Evaluation Issues

Challenges and Future Work

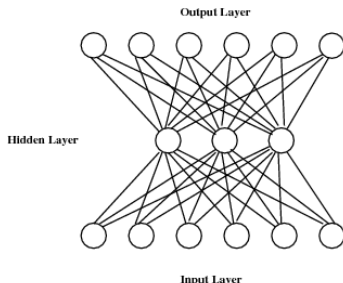
One-Class Classification: Common Approaches

- ▶ Some techniques use:
 - ▶ Artificial Neural Networks
 - ▶ Support Vector Machines
 - ▶ kNN based approaches
 - ▶ Kernel based approaches
 - ▶ Parzen windows

One-Class Classification: Common Approaches - II

Autoencoders [Japkowicz, 1999]

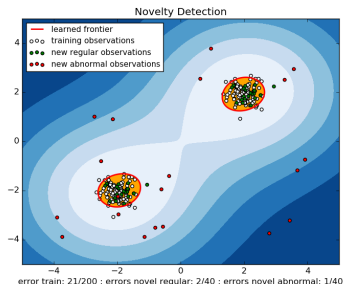
- ▶ three layer network
- ▶ nr. of neurons in the output layer is equal to the input layer
- ▶ the network is trained with backpropagation to reproduce the input at the output layer
- ▶ difference between the input example and the output:
 - ▶ $< \textit{threshold}$: example is from normal class
 - ▶ otherwise: is a counter-example of normal class



One-Class Classification: Common Approaches - III

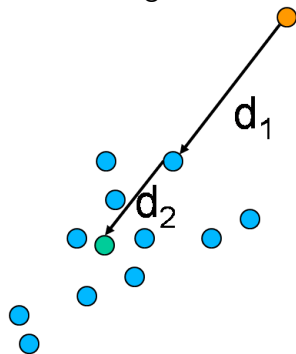
Support Vector Data Description [Tax and Duin, 2004]

- ▶ obtains a spherical boundary, in the feature space, around the data
- ▶ the volume of this hypersphere is minimized, to reduce the effect of incorporating outliers in the solution
- ▶ examples lying outside the hypersphere are considered abnormal



One-Class Classification: Nearest neighbor

Nearest neighbor for novelty detection (Tax, 2001)



If $d_1 / d_2 > 1 \rightarrow$ reject

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

MINAS algorithm

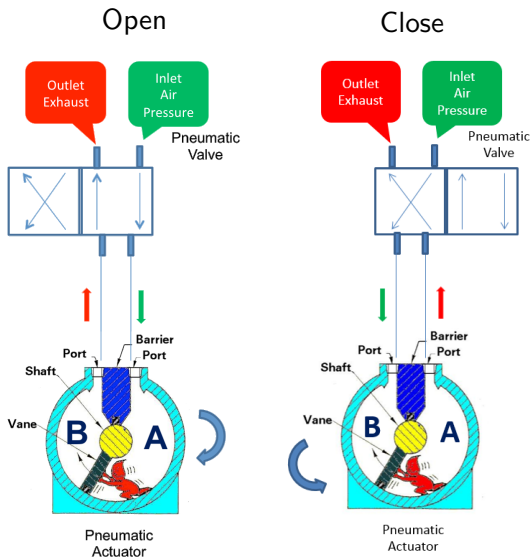
Evaluation Issues

Challenges and Future Work

One-Class Classification: Case Study: Predict Train Door Failures - II

Original Data

- ▶ Pressure readers at both chambers.
- ▶ Interval: 1/10 sec.
- ▶ Episodes time-series data set.
- ▶ ≈ 60 observations / episode.
- ▶ Sept. to Dec. 2012.
- ▶ ≈ 500.000 observations.
- ▶ Movement duration not steady.
- ▶ Unlabelled data set.



One-Class Classification: Case Study: Predict Train Door Failures -

III

- ▶ Problem Setting
 - ▶ Data: sequence of cycles (Open or Close).
 - ▶ Goal: predict the structural state of the door: Normal or Failure.
- ▶ Two-Step Approach:
 1. **Abnormal Cycle Detection**
 - ▶ classify each cycle as Normal (1) or Abnormal (0)
 2. **Failure Sequence Detection**
 - ▶ classify sequence of cycles as Normal or Failure.

One-Class Classification: Case Study: Predict Train Door Failures - IV

- ▶ Abnormal cycle detection assumes independence between observations.
- ▶ Structural failure detection must take into account sequence of observations.
- ▶ For that, we use a **Low-Pass Filter** to post-process cycle classification output.

$$y_i = \begin{cases} 1 & \text{if } i = 0 \\ y_{i-1} + \alpha * (x_i - y_{i-1}) & \text{if } i > 0 \end{cases}$$

where, for instant i , y_i is filter output and x_i is original signal.

- ▶ The α parameter smoothes abrupt changes in the original signal.
- ▶ Lower values of α cause more inertia.
- ▶ Failure Threshold: $y_i < 0.5$

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

MINAS algorithm

Evaluation Issues

Challenges and Future Work

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

MINAS algorithm

Evaluation Issues

Challenges and Future Work

Novelty Detection: Problem Definition

- ▶ Training set (Offline Phase)
 - ▶ $D_{tr} = (X_1, y_1), (X_2, y_2), \dots, (X_m, y_m)$
 - ▶ X_i : vector of input attributes for the i th example
 y_i : target attribute
 - ▶ $y_i \in Y_{tr}$ where $Y_{tr} = c_1, c_2, \dots, c_L$

- ▶ When new data arrive (Online Phase)
 - ▶ Given a sequence of unlabelled examples X_{new}
 - ▶ Goal: Classify X_{new} in Y_{all} where $Y_{all} = c_1, c_2, \dots, c_L, \dots, c_K$ and $K > L$

Open-set Recognition

Anuran species recognition using a hierarchical classification approach

Juan G. Colonna^{1,2}, João Gama², and Eduardo F. Nakamura¹

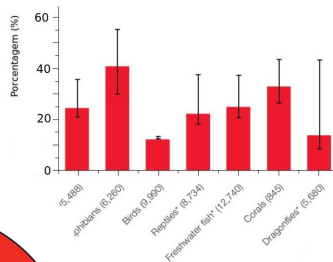
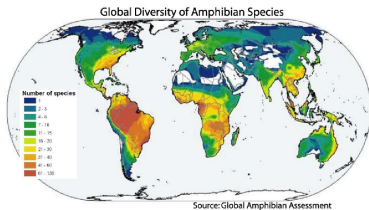
¹Federal University of Amazonas (UFAM), Institute of Computing (Icomp)
²Laboratory of Artificial Intelligence and Decision Support (LIAAD), INESC Tec
{juancolonna, nakamura}@icomp.ufam.edu.br
jgama@fep.up.pt



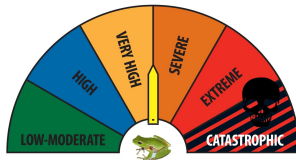
Getting more from family, genus and species of frogs

Introduction - Why frogs?

- **Anura** is the name of an order of animals in the **Amphibian** class which lack a tail, this includes **frogs** and **toads**.



- **Frogs are very sensitive to environmental changes**



Why monitor populations of frogs?

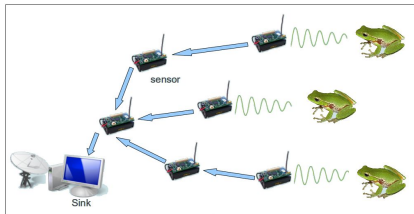
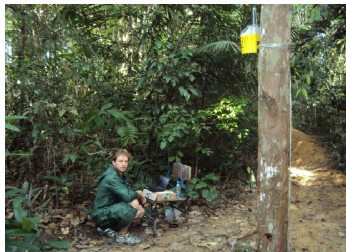
Hypothesis: Tracking the changes in the anuran populations can help us to determine ecological problems in early stages.



It involves several manual tasks!

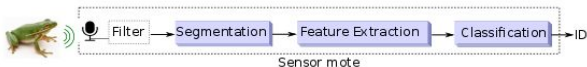
Proposal

Signal processing (SP) + Wireless Sensor Networks (WSN) + Machine Learning (ML)



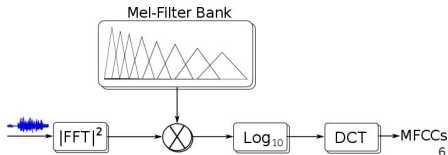
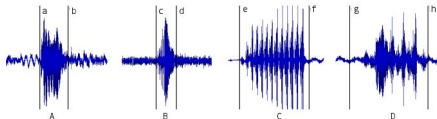
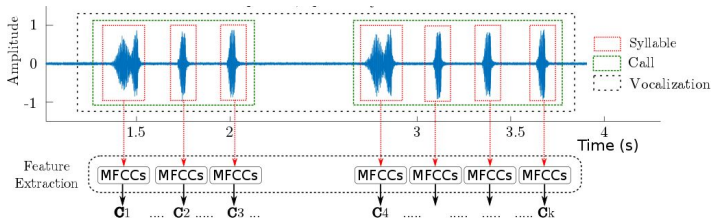
Advantages: It is Automatic, less intrusive and allows long term monitoring.

How to do that?



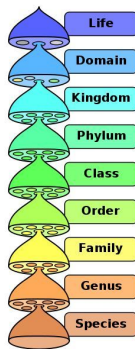
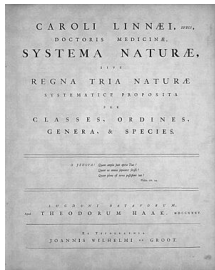
- 1) Pre-processing:
 - a) Filter: band-pass filter, wavelet decomposition, etc.
 - b) Segmentation: **syllable-based approach** (x_k)
- 2) Feature Extraction: that maps $x_k \rightarrow c_k$
 - a) Mel-frequency cepstral coefficients (MFCCs)
 - b) Spectral centroid, Spectral bandwidth, Pitch, etc.
- 3) Recognition: ML technique to classify $c_k \rightarrow \text{ID}$ (species ID)
 - a) Support Vector Machine, kNN, Tree, etc.

Segmentation and feature extraction



Knowledge organization

Carl Linnaeus has defined a particular form of biological organization called *taxonomy* in his work *Systema Naturae* (1735).



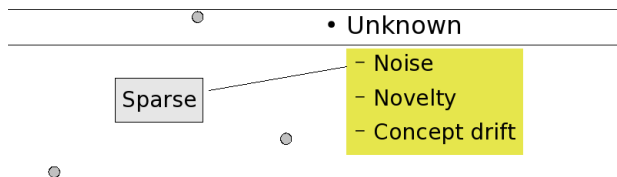
Novelty Detection: Problem Definition

- ▶ Offline Phase
 - ▶ All training examples belong to the *known* classes.
- ▶ Online Phase
 - ▶ Examples not explained by the current model are labeled as *unknown*.
 - ▶ Cohesive group of *unknown* examples are used to detect *novel* classes or extensions to the *known* classes.

Novelty Detection: Problem Definition - II

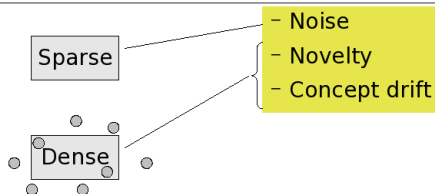
- ▶ In data streams, concepts are hardly ever constant.
- ▶ It is important **to distinguish novelty from:**
 - ▶ noise and outliers
 - ▶ concept drift
 - ▶ concept evolution
 - ▶ recurring concepts

Novelty Detection: Problem Definition - III

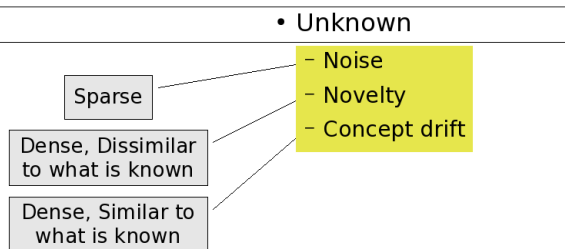


Novelty Detection: Problem Definition - III

- Unknown



Novelty Detection: Problem Definition - III



Novelty Detection: Problem Definition - IV

- ▶ In data streams scenarios:
 - ▶ new concepts may appear
 - ▶ known concepts may evolve, disappear or reappear
- ▶ By monitoring the data stream, emerging concepts may be discovered
- ▶ Emerging concepts may represent
 - ▶ an extension to a known concept (extension)
 - ▶ a novel concept (novelty)

Novelty Detection: Problem Definition - VI

- ▶ Novelty Detection Systems
 - ▶ **OLINDDA**: OnLine Novelty and Drift Detection Algorithm [Spinosa et al., 2007]
 - ▶ **ECSMiner**: Enhanced Classifier for data Streams with novel class Miner [Masud et al., 2011]
 - ▶ **MINAS**: Multi-class learNing Algorithm for data Streams [de Faria et al., 2016]

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

MINAS algorithm

Evaluation Issues

Challenges and Future Work

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

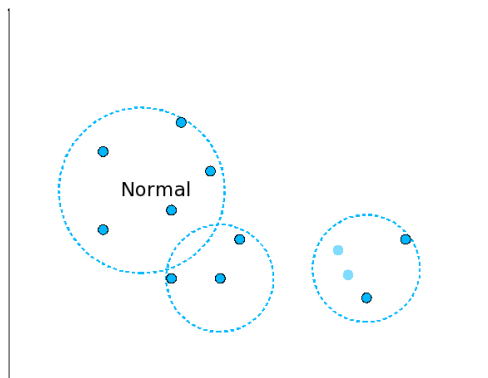
MINAS algorithm

Evaluation Issues

Challenges and Future Work

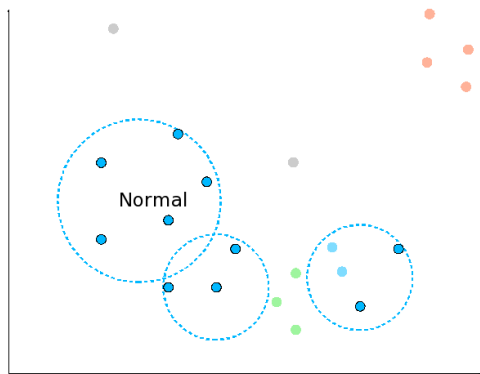
Novelty Detection: MINAS algorithm

MINAS: Multi-class learNing Algorithm for data Streams
[de Faria et al., 2016]



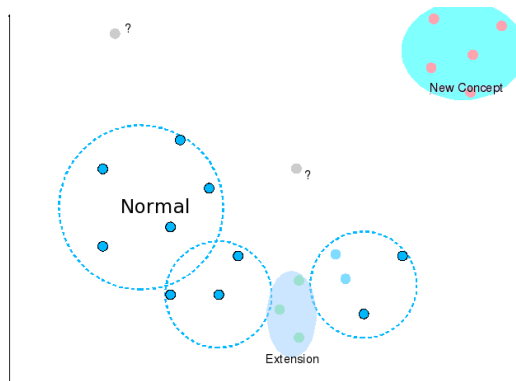
Novelty Detection: MINAS algorithm

MINAS: Multi-class learniNg Algorithm for data Streams
[de Faria et al., 2016]



Novelty Detection: MINAS algorithm

MINAS: Multi-class learnNing Algorithm for data Streams
[de Faria et al., 2016]



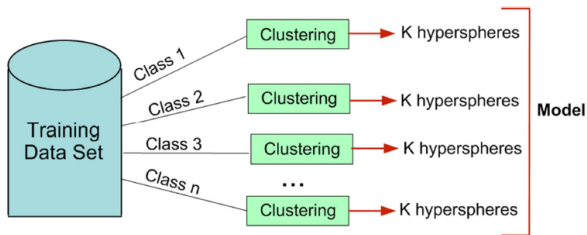
Novelty Detection: MINAS algorithm - II

- ▶ unsupervised algorithm for novelty detection in data streams multi-class problems
 - ▶ training examples are composed by many classes
 - ▶ there may be also several novel classes
- ▶ use of offline (training) and online phases
 - ▶ in each phase learns one or more classes
- ▶ cohesive set of examples is necessary to learn new concepts or extensions
 - ▶ isolated examples are not considered as novelty

Novelty Detection: MINAS algorithm- III

Offline Phase

- ▶ learns a decision model based on the known concepts about the problem (k-means or Clustream)
- ▶ runs only once
- ▶ each class is represented by a set of clusters (hyperspheres)



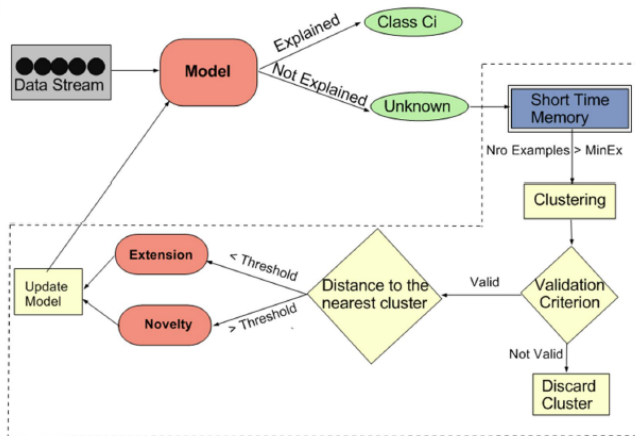
Novelty Detection: MINAS algorithm- IV

Online Phase

- ▶ receives new unlabelled examples from the stream
- ▶ classifies each new example as one of known classes or as unknown
- ▶ unknown examples are stored in the Short Term Memory
- ▶ from time to time
 - ▶ finds clusters in the examples stored in the Short Term Memory
 - ▶ clusters far away from existing ones: novel concept.
 - ▶ clusters close to existing ones: extend known concepts.

Novelty Detection: MINAS algorithm - V

Online Phase



Novelty Detection: MINAS algorithm - VI

Treatment of Outliers

- ▶ clustering is applied to the unknown examples
- ▶ each cluster is validated by the evaluation of its representativeness and cohesiveness
- ▶ clusters with low value are considered invalid and removed
- ▶ however, their examples stay in a temporary memory
- ▶ if there is no space available, the oldest example is removed
- ▶ there is a high chance that the removed examples are noise or outliers

Outline

Introduction

- Novelty

- Applications

One-Class Classification

- Problem Definition

- Common Approaches

- Case Study: Predict Train Door Failures

Novelty Detection

- Problem Definition

- Key Aspects

- MINAS algorithm

Evaluation Issues

Challenges and Future Work

Evaluation Issues: Adaptation of binary classification metrics

- ▶ *Precision* and *Recall* [Albertini and de Mello, 2007]

$$Precision = \frac{\# \text{ true detected novelties}}{\# \text{ detected novelties}} \quad Recall = \frac{\# \text{ true detected novelties}}{\# \text{ novelties}}$$

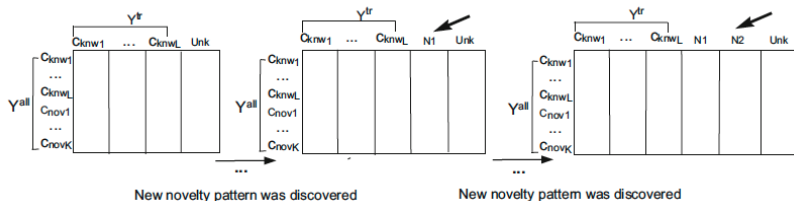
- ▶ *Mnew* and *Fnew* [Spinosa et al., 2007]

$$M_{new} = \frac{100 * \# \text{ false detected novelties}}{\# \text{ novelties}} \quad F_{new} = \frac{100 * \# \text{ false detected normalities}}{\# \text{ total} - \# \text{ novelties}}$$

- ▶ How to reflect the *unknown* label in the evaluation?
- ▶ How to extend the evaluation to a multi-class scenario?

Evaluation Issues: A Rectangular Confusion Matrix [de Faria et al., 2016]

- ▶ At the beginning of online phase, the model is composed by the classes learned offline.
- ▶ When a new concept is discovered, the model is updated and a new column is added to the confusion matrix.



- ▶ rows: true classes (known + novelty)
- ▶ columns: predicted classes (known + novelty patterns + unknown)

Outline

Introduction

Novelty

Applications

One-Class Classification

Problem Definition

Common Approaches

Case Study: Predict Train Door Failures

Novelty Detection

Problem Definition

Key Aspects

MINAS algorithm

Evaluation Issues

Challenges and Future Work

Challenges and Future Work

- ▶ Most of the techniques use one-class classification
 - ▶ many real-world applications are, in fact, multi-class scenarios
- ▶ Assumption that true labels will be available
 - ▶ time consuming
 - ▶ support from domain expert
- ▶ Outliers, noise, changing environments
 - ▶ depends on the data set
 - ▶ concepts may evolve gradually or abruptly
 - ▶ distinguish noise, outliers from a novel concept

Challenges and Future Work - II

- ▶ Recurring contexts
 - ▶ an important phenomenon observed in many real-world applications (e.g. climate change, electricity demand)
 - ▶ systems typically use a forgetting mechanism of old concepts;
 - ▶ a recurring class may be confused with the emergence of a new class → it leads to high false positive rates
 - ▶ relearn an old concept is a waste of effort
 - ▶ ideally, they should be saved and reexamined at some later time
 - ▶ identify when a concept is reappearing

Challenges and Future Work - III

- ▶ When to apply novelty detection in data streams
 - ▶ whenever new example arrives is time consuming
 - ▶ define the time interval
- ▶ Algorithms to induce the decision model
 - ▶ supervised algorithms need labeled examples
 - ▶ unsupervised algorithms (e.g. kmeans) assume that classes constitute hyperspheres, need nr. clusters as input., handle only numerical attributes
- ▶ Evaluation issues and experimental methodology
 - ▶ lack of standards

Some References I



Abdallah, Z. S., Gaber, M. M., Srinivasan, B., and Krishnaswamy, S. (2016).
Anynovel: detection of novel concepts in evolving data streams.
Evolving Systems, 7(2):73–93.



Albertini, M. K. and de Mello, R. F. (2007).
A self-organizing neural network for detecting novelties.
In *Proceedings of the 2007 ACM Symposium on Applied Computing, SAC '07*, pages 462–466, New York, NY, USA. ACM.



Costa, B. S. J., Angelov, P. P., and Guedes, L. A. (2014).
Real-time fault detection using recursive density estimation.
Journal of Control, Automation and Electrical Systems, 25(4):428–437.



de Faria, E. e., Ponce de Leon Ferreira Carvalho, A., and Gama, J. a. (2016).
MINAS: multiclass learning algorithm for novelty detection in data streams.
Data Mining and Knowledge Discovery, 30(3):640–680.



Faria, E. R., Gonçalves, I. J. C. R., de Carvalho, A. C. P. L. F., and Gama, J. (2016).
Novelty detection in data streams.
Artificial Intelligence Review, 45(2):235–269.



Japkowicz, N. (1999).
Concept-learning in the absence of counter-examples: An autoassociation-based approach to classification.



Krawczyk, B. and Woźniak, M. (2015).
One-class classifiers with incremental learning and forgetting for data streams with concept drift.
Soft Computing, 19(12):3387–3400.



Masud, M. M., Gao, J., Khan, L., Han, J., and Thuraisingham, B. M. (2011).
Classification and novel class detection in concept-drifting data streams under time constraints.
IEEE Trans. Knowl. Data Eng., 23(6):859–874.

Some References II



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

Failure prediction - an application in the railway industry.

In *Discovery Science - 17th International Conference, DS 2014*, pages 264–275.



Pimentel, M. A. F., Clifton, D. A., Clifton, L., and Tarassenko, L. (2014).

A review of novelty detection.

Signal Process., 99:215–249.



Ribeiro, R. P., Pereira, P., and Gama, J. (2016).

Sequential anomalies: a study in the railway industry.

Machine Learning, pages 1–27.



Spinosa, E. J., de Leon Ferreira de Carvalho, A. C. P., and Gama, J. (2007).

Olindda: a cluster-based approach for detecting novelty and concept drift in data streams.

In Cho, Y., Wainwright, R. L., Haddad, H., Shin, S. Y., and Koo, Y. W., editors, *SAC*, pages 448–452.

ACM.



Tax, D. M. J. and Duin, R. P. W. (2004).

Support vector data description.

Machine Learning, 54(1):45–66.