# The area under the ROC curve as a variable selection criterion for multiclass classification problems

M. de Figueiredo<sup>1,2</sup> C. B. Y. Cordella<sup>3</sup> D. Jouan-Rimbaud Bouveresse<sup>1,3</sup> X. Archer<sup>2</sup> J.-M. Bégué<sup>2</sup> D. N. Rutledge<sup>1,\*</sup>

<sup>1</sup> UMR Ingénierie Procédés Aliments, AgroParisTech, Inra, Université Paris-Saclay, 91300 Massy, France

<sup>2</sup> Laboratoire Central de la Préfecture de Police, 39bis rue de Dantzig, 75015 Paris, France

<sup>3</sup> UMR Physiologie de la Nutrition et du Comportement Alimentaire, AgroParisTech, Inra, Université Paris-Saclay, 75005 Paris, France

\*Correspondence : rutledge@agroparistech.fr



- ROC analysis is mostly used for two-class problems
  - Receiver Operating Characteristic (ROC) analysis uses the area under the ROC curve (AUC) as criterion to evaluate the discrimination between two classes
  - The AUC value is a direct measure of a classifier performances representing how well two classes are separated
  - Maximizing the AUC value is equivalent to maximize the separation between two classes
  - Classifier performances can be improved by selecting an appropriate subset of variables that combined together provide an optimal AUC
- The objective
  - Nowadays many fields are confronted with multiclass classification problems
  - Demonstrate how two-class ROC analysis and its associated AUC can be used to perform variable selection within a multiclass problem framework

### Lets consider two vectors of continuous responses whose values range from 0 to 1

- Each vector represents a group (positive or negative) that is to be discriminated from the other one
- The better the separation between the two distributions, the higher the AUC value is



**PREFECTURE DE POLICE** 

CPP

- Lets consider two vectors of continuous responses whose values range from 0 to 1
  - Example of decision threshold defining classification rates ٠
  - How is the AUC calculated? Force variation of the threshold between 0 and 1 •
  - At each step, the TPR and FPR are recorded ٠

4

**PREFECTURE DE POLICE** 

CPP

The ROC curve is the representation of the TPR recorded as a function of the FPR •



### **ROC** analysis: basic principles



### Lets consider now that the continuous responses are Euclidean distances

- Each vector element corresponds to the distance between two samples for a given set of variables
- Either pairs of samples belong to the same class (related samples) and should be close to each other or the pairs of samples belong to different classes (unrelated samples) and should be far from each other
- A similarity measurement quantifies how much

### alike two samples are

 Elements in the vector of (un)related samples correspond to the Euclidean distance between pairs of samples (not) belonging to the same class



# The variable selection method

### The goals of the variable selection method based on the AUC criterion

- Select a subset of variables maximizing the AUC value
- Maximizing the AUC value is equivalent to maximize the separation between the distributions of related and unrelated samples
- Maximizing the separation between the two distributions is equivalent to bringing closer together in the multivariate space samples belonging to the same class and, at the same time, separating groups of samples belonging to different classes
- This point is true whether we have a two-class problem at hand or a multiclass one
- As long as the separation between the two distributions is maximized, no matter the number of classes, it is always possible to reduce the problem to a two-class one
- If the AUC represents how well two classes are separated, it can also represent a global value of how well several classes are separated







Chimiométrie XIX, Conservatoire National des Arts et Métiers, Paris, France, 30 et 31 janvier 2018

**DD** THEFECTURE DE POLICE

LCPP

UNIVERSITE

AgroParisTech

### The variable selection method



### AUC estimation

- Everytime an AUC value is estimated according to Hand & Till (2001), vectors of intra- and inter-similarities are calculated beforehand for a given subset of variables
- This estimation is not based on the effective representation of the ROC curve

$$\widehat{AUC} = \frac{S_0 - n_1(n_1 + 1)/2}{n_0 n_1}$$

Hand, D. J., & Till, R. J. (2001). A Simple Generalisation of the Area Under the ROC Curve for Multiple Class Classification Problems. Machine Learning, 45(2), 171-186.

Chimiométrie XIX, Conservatoire National des Arts et Métiers, Paris, France, 30 et 31 janvier 2018

**PREFECTURE DE POLICE** 

LCPP

UNIVERSITE PARIS-SACLAY

AgroParisTech

**M R** 145

# A practical example : the dataset, samples preparation and analysis

#### Dataset

190 gasoline samples representing 4 qualities (SP95, SP98, SP95-E10 and SP-E85) and bought over 4 seasons (fall, winter, spring and summer) in 20 different gas stations

#### QC sample

Mix of 98 gasoline samples prepared and analyzed in parallel with regular samples

#### **Samples preparation**

Passive headspace thermal extraction onto Tenax TA<sup>®</sup> tubes

Samples prepared in triplicates



Data samples acquired

**570** gasoline samples **86** QC samples

Chimiométrie XIX, Conservatoire National des Arts et Métiers, Paris, France, 30 et 31 janvier 2018

## A practical example : data preparation



Fingerprinting. Environmental Forensics, 3(3-4), 263-278.

Chimiométrie XIX, Conservatoire National des Arts et Métiers, Paris, France, 30 et 31 janvier 2018

**PREFECTURE DE POLICE** 

LCPP

AgroParisTech

### Cross-validation PLS-DA

- Using the SAISIR toolbox (Cordella, C. B. Y. et Bertrand, D., 2014)
- 100 holdout/testset cross-validation PLS-DA using random subsets of 2/3 of the samples for the calibration and the remaining 1/3 as a validation set
- Cross-validation PLS-DA was performed on the gasoline dataset with and without variable selection to evaluate the predictability of the models
- The accuracy (percentage of correct classification) was used to evaluate the predictability of the models

Cordella, C.B.Y. & Bertrand, D. SAISIR: A new general chemometric toolbox. TrAC Trends in Analytical Chemistry 54, 75-82, 2014.

Note that columns marked by (\*) concern results with a prior variable selection process

Property	Increment	Ratios selected	AUC	LVs*	Accuracy*	LVs	Accuracy
Quality	0.001	4	0.9877	2	79.45±3.56 %	3	77.20±3.82 %



Note that columns marked by (\*) concern results with a prior variable selection process

Property	Increment	Ratios selected	AUC	LVs*	Accuracy*	LVs	Accuracy
Quality	0.001	4	0.9877	2	79.45±3.56 %	3	77.20±3.82 %



**PREHECTURE DE POLICE** 

Note that columns marked by (\*) concern results with a prior variable selection process

Property	Increment	Ratios selected	AUC	LVs*	Accuracy*	LVs	Accuracy
Quality	0.001	4	0.9877	2	79.45±3.56 %	3	77.20±3.82 %
Quality (no SP-E85)	0.001	9	0.9862	2	96.87±1.14 %	5	79.15±2.67 %



**PREHECTURE DE POLICE** 

ENIA

Note that columns marked by (\*) concern results with a prior variable selection process

Property	Increment	Ratios selected	AUC	LVs*	Accuracy*	LVs	Accuracy
Quality	0.001	4	0.9877	2	79.45±3.56 %	3	77.20±3.82 %
Quality (no SP-E85)	0.001	9	0.9862	2	96.87±1.14 %	5	79.15±2.67 %
Season	0.001	21	0.9789	4	98.88±0.68 %	9	97.67±1.14 %



E

- An adequate variable selection procedure should in general :
  - Substantially reduce the dimensionality of the data
  - Improve the predictability of the models built
  - Facilitate the interpretation of the chemical systems

### • The variable selection presented here :

- Performs data dimensionality reduction by selecting variables in their orginal space and not by building linear combinations of the original variables (reduction close to 100% in the given example)
- Proposes a way to use the traditional two-class ROC analysis within a multiclass context
- Can implement different distance measures other than Euclidean distance
- Can be sped up by avoiding the calculation of similarities between all possible pairs of samples
- Because of the evermore growing amounts of data produced by modern analytical procedures, being able to reduce the dimensionality of the data while maximizing the predictive power of the models and their interpretability if of utmost importance



# **Additional information – Preparation of samples**



U

ENIAL

### **Additional information – Analysis of samples**



Chimiométrie XIX, Conservatoire National des Arts et Métiers, Paris, France, 30 et 31 janvier 2018