



# Adaptive Mixture Discriminant Analysis for Supervised Learning with Unobserved Classes

Charles Bouveyron

## ► To cite this version:

Charles Bouveyron. Adaptive Mixture Discriminant Analysis for Supervised Learning with Unobserved Classes. *Journal of Classification*, 2014, 31 (1), pp.49-84. 10.1007/s00357-014-9147-x. hal-00392297v4

**HAL Id: hal-00392297**

**<https://hal.science/hal-00392297v4>**

Submitted on 23 Jun 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Adaptive Mixture Discriminant Analysis for Supervised Learning with Unobserved Classes

Charles BOUVEYRON

`charles.bouveyron@univ-paris1.fr`

Laboratoire SAMM, EA 4543

Université Paris 1 Panthéon-Sorbonne

90 rue de Tolbiac, 75013 Paris, France

## Abstract

In supervised learning, an important issue usually not taken into account by classical methods is the possibility of having in the test set individuals belonging to a class which has not been observed during the learning phase. Classical supervised algorithms will automatically label such observations as belonging to one of the known classes in the training set and will not be able to detect new classes. This work introduces a model-based discriminant analysis method, called adaptive mixture discriminant analysis (AMDA), which is able to detect several unobserved groups of points and to adapt the learned classifier to the new situation. Two EM-based procedures are proposed for parameter estimation and model selection criteria are used for selecting the actual number of classes. Experiments on artificial and real data demonstrate the ability of the proposed method to deal with complex and real word problems. The proposed approach is also applied to the detection of unobserved communities in social network analysis.

**Key-words:** supervised classification, unobserved classes, adaptive learning, multi-class novelty detection, model-based classification, social network analysis.

## 1 Introduction

The usual framework of supervised classification assumes that all existing classes in the data have been observed during the learning phase and does not take into account the possibility of having in the test set individuals belonging to a class which has not been observed. In particular, such a situation could occur in the case of rare classes or in the case of an evolving population. For instance, an important problem in Biology is the detection of novel species which could appear at any time resulting from structural or physiological modifications. In the same manner, the detection of unobserved communities is a major issue in social network analysis for security or commercial reasons. Unfortunately, classical supervised algorithms,

like support vector machines or linear discriminant analysis, will automatically label observations from a novel class as belonging to one of the known classes in the training set and will not be able to detect new classes. It is therefore important to find a way to allow the supervised classification methods to detect unobserved situations and to adapt themselves to the new configurations.

In statistical learning, the problem of classification with unobserved classes is a problem which has received very few attention. Indeed, both supervised and unsupervised classification contexts have been widely studied but intermediate situations have received less attention. We would like however to mention two related topics in statistical learning called semi-supervised classification and novelty detection. Semi-supervised classification focuses on learning with partially labeled data whereas novelty detection tries to detect new or unknown data points in the test set. Unfortunately, both approaches are unable to detect several unobserved groups of points in the test set and to adapt the classifier to the new situation.

To overcome this problem, this work introduces an approach based on the mixture model which combines unsupervised and supervised learning for detecting unobserved groups of observations in the test set and for adapting the supervised classifier to the new situation. The adapted classifier could be then used to correctly classify new observations in the future. Two EM-based approaches are proposed for parameter estimation: an inductive approach, which is made of a learning and a discovering phase, and a transductive approach which considers all available observations for learning in a unique step. The detection of the number of unobserved classes is done using model selection criteria. Finally, once the classifier adapted, the classification of new observations can be then done through the classical *maximum a posteriori* rule.

The paper is organized as follows. A brief review on generative supervised classification is given in Section 2 as well as a presentation of related works in supervised learning with unobserved classes. Section 3 introduces an adaptive discriminant analysis method based on the mixture model which is able to detect unobserved classes and considers parameter estimation as well. Experimental results highlighting the main features of the proposed method on simulated and real datasets are presented in Section 4 as well as comparison with novelty detection methods. Section 5 presents an application of the proposed method to the detection of unobserved communities in social networks. Finally, Section 6 proposes some concluding remarks and discusses further works.

## 2 Related works

This section briefly reviews the supervised classification problem and solutions based on the mixture model before to present related works on supervised learning with unobserved classes.

## 2.1 Generative supervised classification

Supervised classification, also known as discriminant analysis in the literature, aims to build a supervised classifier from a complete set of learning observations  $\{(x_1, z_1), \dots, (x_n, z_n)\}$  where  $x_i$  is an observation described by  $p$  variables and  $z_i \in \{1, \dots, K\}$  indicates the class of  $x_i$ . The learned classifier is then used for assigning a new observation  $x^*$  to one of the  $K$  known classes. Among all existing approaches, generative discriminant analysis is very popular because of its probabilistic background and its efficiency. Generative supervised classification assumes that the observations  $\{x_1, \dots, x_n\}$  and their labels  $\{z_1, \dots, z_n\}$  are respectively independent realizations of random vectors  $X \in \mathbb{R}^p$  and  $Z \in \{1, \dots, K\}$  and that the conditional density of  $X$  given that  $Z = k$  is a parametric density  $f_k$  parametrized by  $\theta_k$ . Consequently, the marginal density of  $X$  is given by:

$$f(x) = \sum_{k=1}^K \pi_k f_k(x; \theta_k),$$

where  $\pi_k$  is the prior probability of the  $k$ th class. The classification of a new observation  $x^*$  is done afterward using the *maximum a posteriori* (MAP) rule which assigns  $x^*$  to the class with the highest posterior probability  $P(Z = k|X = x)$ :

$$P(Z = k|X = x) = \frac{\pi_k f_k(x; \theta_k)}{\sum_{k=1}^K \pi_k f_k(x; \theta_k)}.$$

We refer the reader to [30] for more details on generative discriminant analysis. The following paragraphs review the most used parametric densities in generative supervised classification.

**Mixture of Gaussians** Among all parametric densities, the Gaussian model is probably the most used in classification. The Gaussian mixture model has been extensively studied in the last decades and used in many situations (see [3] and [31] for a review). Therefore, if the Gaussian model is chosen,  $f_k(x; \theta_k)$  will denote the density of a multivariate Gaussian density parametrized by  $\theta_k = \{\mu_k, \Sigma_k\}$  where  $\mu_k$  and  $\Sigma_k$  are respectively the mean and covariance matrix of  $k$ th component of the mixture.

**Mixture of parsimonious Gaussians** In some situations, modeling the data with a full covariance matrix can be too expensive in terms of number of parameters to estimate. In such a case, it is possible to make additional assumptions on the structure of the covariance matrix. For example, in the well-known Linear Discriminant Analysis (LDA) method, the covariance matrices of the different components are supposed to be equal to a unique covariance matrix. It is also possible to assume that the covariance matrix of each mixture component is diagonal or proportional to the identity matrix. These models are known as parsimonious Gaussian models in the literature since they require to estimate less parameters than the classical Gaussian

model. Celeux and Govaert proposed in [12] a family of parsimonious Gaussian models based on an eigenvalue decomposition of the covariance matrix including the previous models. These parsimonious Gaussian models were then applied in [5] to supervised classification.

**Mixture of HD Gaussians** Nowadays, many scientific domains produce high-dimensional data like medical research, image analysis or Biology. Classifying such data is a challenging problem since the performance of classifiers suffers from the *curse of dimensionality* [4]. Classification methods based on Gaussian mixture models are directly penalized by the fact that the number of parameters to estimate grows up with the square of the dimension. Unfortunately, parsimonious models are usually too constrained to correctly fit the data in a high-dimensional space. To overcome this problem, Bouveyron *et al.* proposed recently in [11] a family of Gaussian models adapted to high-dimensional data. This approach, based on the idea that high-dimensional data live in low-dimensional spaces, assumes that the covariance matrix of each mixture component has only  $d_k + 1$  different eigenvalues where  $d_k$  is the dimension of the subspace of the  $k$ th mixture component. These Gaussian models were then used in [10] for high-dimensional data clustering.

**Mixture with a noise component** Banfield and Raftery have introduced in [3] a mixture model with a noise component in order to improve the robustness of the cluster analysis on noisy datasets. The original work proposed to add to the mixture model a uniform distribution over the convex hull of the data as an additional component. Good results of the use of this specific mixture model were observed in different situations. Hennig and Coretto [22] proposed recently to use an improper uniform distribution that does not depend on the data for improving the robustness and provide a better approximation of the likelihood than the one proposed in the original work. An application of noise detection is proposed in Section 4.

## 2.2 Semi-supervised classification

The first related topic to supervised classification with unobserved classes is semi-supervised classification. Semi-supervised classification is a topic which has been well studied for several years and which focuses on supervised classification with partially labeled data. Usually, unlabeled data are added to the learning data in order to improve the efficiency of the final classifier. Such an approach is particularly useful when only few labeled observations are available for learning (applications with a high supervision cost like image recognition or cancer detection). A good review on semi-supervised classification can be found in [38] and [25]. Earlier approaches [29, 33] used the EM algorithm to assign unlabeled observations to known classes. Most recent approaches include co-training algorithms [8] and graph-based techniques [25] which used prior informations on unlabeled observations. However, all semi-supervised classification methods are not able to detect unobserved groups of points. More

importantly, they will use those new points to re-estimate the model parameters of known classes and the estimates of known classes parameters will be therefore deteriorated.

### 2.3 Novelty detection

The second and most related topic to supervised classification with unobserved classes is novelty detection. Novelty detection focuses on the identification of new or unknown data for which the learned classifier was not aware during the learning phase. This approach has become very popular in several application fields such as fault detection [14], medical imaging (mass detection in mammograms) [40] or e-commerce [26]. In the last years, many methods have been proposed to deal with this problem. An excellent review on novelty detection methods can be found in [27] and [28] which splits novelty detection methods into two main categories: statistical and neural network based approaches. Approaches based on statistical assumptions usually model the data on their statistical properties and use this information for deciding whether test data comes from the same distribution or not. Among parametric techniques, Chow [13] was the first to propose a threshold for outlier rejection which has been improved in [19] by introducing the classification confidence in the rejection. Gaussian densities were also used in [35] for modeling the learning data and detect outliers using a measure based on the Mahalanobis distance. Extreme value theory was also used in [34] for novelty detection by searching for low or high values in the tails of data distributions. Non-parametric approaches include  $k$ -NN based techniques [21, 32] or Parzen windows [43] for estimating the distribution of the data. Neural networks and kernel methods have been also widely used for novelty detection. Bishop [7] used parametric statistics by post-processing neural networks for detecting new data distribution whereas a probability density estimation of neural network outputs is used in [16] as a measure of novelty. Another approach based on neural networks was proposed in [36] which used a thresholding on the neural network output for detection new samples. Kohonen proposed also in [24] two types of novelty detectors based on self-organizing maps (SOM). More recently, Tax and Duin [41] as well as Schölkopf [39] used support vector machines (SVM) for distinguishing know and unknown objects. However, even though all these methods are able to detect new or unobserved data points, no one of them is able to recognize unobserved homogeneous groups of points and to adapt the classifier to the new situation for classifying future observations.

## 3 Adaptive mixture discriminant analysis

We introduce in this section an adaptive model-based classifier able to detect novel classes which have not been observed during the learning phase. Parameter estimation, model selection and classification of future observations will be discussed as well.

### 3.1 The mixture model

Let us consider a classical parametric mixture model of  $K$  components: the observations  $\mathcal{X} = \{x_1, \dots, x_n\} \in \mathbb{R}^p$  are assumed to be independent realizations of a random vector  $X \in \mathbb{R}^p$  with density:

$$f(x; \Theta) = \sum_{k=1}^K \pi_k f_k(x; \theta_k), \quad (1)$$

where  $\pi_k \geq 0$  for  $k = 1, \dots, K$  are the mixing proportions (with the constraint  $\sum_{k=1}^K \pi_k = 1$ ),  $f_k(x; \theta_k)$  is the density of the  $k$ th component of the mixture parametrized by  $\theta_k$  and finally  $\Theta = \{\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K\}$ . We refer to the previous section regarding the choice of the mixture densities. For the mixture model (1), the complete log-likelihood has the following form:

$$\ell(x_1, \dots, x_n; \Theta) = \sum_{i=1}^n \sum_{k=1}^K s_{ik} \log(\pi_k f_k(x_i; \theta_k)),$$

where  $s_{ik} = 1$  if  $x_i$  belongs to the  $k$ th class and  $s_{ik} = 0$  otherwise. However, this work considers a specific learning situation in which only  $C$  classes are represented in the learning dataset  $\mathcal{X} = \{x_1, \dots, x_n\}$  with  $1 \leq C \leq K$ , *i.e.* one or several classes could be not represented in  $\mathcal{X}$ . Therefore, the mixture parameter estimation can not be done using the classical way and two alternative estimation procedures are proposed below.

### 3.2 Parameter estimation: transductive approach

The most intuitive way to identify unobserved classes in the test set is certainly the transductive approach which works on the gathering of learning and test sets. Indeed, since the learning sample  $\mathcal{X} = \{x_1, \dots, x_n\}$  and the test sample  $\mathcal{X}^* = \{x_1^*, \dots, x_{n^*}^*\}$  are assumed to come from the same population, both samples can be used to estimate model parameters. This would be the general framework of semi-supervised classification if  $C = K$  but semi-supervised classification methods can not be used in our context since  $K$  can be strictly larger than  $C$ . We therefore propose to adapt the classical EM algorithm [15] used in semi-supervised classification to the detection of unobserved classes. In the transductive learning case, the log-likelihood of model (1) has the following form:

$$\ell(\mathcal{X}, \mathcal{X}^*; \Theta) = \sum_{i=1}^n \sum_{k=1}^C s_{ik} \log(\pi_k f_k(x_i; \theta_k)) + \sum_{i=1}^{n^*} \sum_{k=1}^K s_{ik}^* \log(\pi_k f_k(x_i^*; \theta_k)).$$

A constrained version of the EM algorithm is presented below to jointly estimate model parameters while searching for new classes. The joint estimation procedure alternates between the following E and M steps at each iteration  $q$ :

- **E step:** on the one hand, the conditional probabilities  $P(Z = k | X = x_i)$  remain fixed for the learning observations  $\{x_1, \dots, x_n\}$  and are equal to  $s_{ik}$ , for  $i = 1, \dots, n$  and

$k = 1, \dots, K$ , where  $s_{ik} = 1$  if  $x_i$  belongs to the  $k$ th class and  $s_{ik} = 0$  otherwise. On the other hand, the conditional probabilities  $t_{ik}^{*(q)} = P(Z = k | X = x_i^*)$  are updated for the test sample  $\{x_1^*, \dots, x_{n^*}^*\}$ , *i.e.* for  $i = 1, \dots, n^*$  and  $k = 1, \dots, K$ , according to the mixture parameters as follows:

$$t_{ik}^{*(q)} = \frac{\hat{\pi}_k^{(q-1)} f_k(x_i^*; \hat{\theta}_k^{(q-1)})}{f(x; \hat{\Theta}^{(q-1)})},$$

where  $\hat{\pi}_k^{(q-1)}$  and  $\hat{\theta}_k^{(q-1)}$  are the mixture parameters estimated in the M step at the step  $(q - 1)$ .

- **M step:** the parameters of the  $C$  observed classes and of the  $K - C$  unobserved classes are estimated by maximizing the conditional expectation of the completed likelihood. Therefore, this step updates now the estimates of parameters  $\pi_k$  and  $\theta_k$  for  $k = 1, \dots, K$ . In the case of the Gaussian mixture, the update formulas for the parameter estimates are, for  $k = 1, \dots, K$ :

$$\begin{aligned} \hat{\pi}_k^{(q)} &= \frac{n_k^{(q)} + n_k^{*(q)}}{n + n^*}, & \hat{\mu}_k^{(q)} &= \frac{1}{n_k^{(q)} + n_k^{*(q)}} \left( \sum_{i=1}^n s_{ik} x_i + \sum_{i=1}^{n^*} t_{ik}^{*(q)} x_i^* \right), \\ \hat{\Sigma}_k^{(q)} &= \frac{1}{n_k^{(q)} + n_k^{*(q)}} \left( S_k^{(q)} + S_k^{*(q)} \right). \end{aligned}$$

where  $S_k^{(q)} = \sum_{i=1}^n s_{ik} (x_i - \hat{\mu}_k^{(q)})^t (x_i - \hat{\mu}_k^{(q)})$ ,  $S_k^{*(q)} = \sum_{i=1}^{n^*} t_{ik}^{*(q)} (x_i^* - \hat{\mu}_k^{(q)})^t (x_i^* - \hat{\mu}_k^{(q)})$ ,  $n_k^{(q)} = \sum_{i=1}^n s_{ik}$  and  $n_k^{*(q)} = \sum_{i=1}^{n^*} t_{ik}^{*(q)}$ .

Proofs of these results are given in Appendix A.1.

### 3.3 Parameter estimation: inductive approach

The inductive learning context is, conversely to the previous situation, a more classical situation in supervised classification because it is more convenient to keep only model parameters to classify new observations than keeping all learning observations. In particular, the inductive approach is the only tenable approach for large dataset classification and real-time dynamic classification. However, the inductive approach poses a more complex problem since the mixture parameter estimation can not be done using the classical way. We therefore propose hereafter an inductive approach made of a learning phase and a discovery phase.

**The learning phase** In this first phase, only learning observations are considered and, since the data of the learning set are complete, *i.e.* a label  $z_i \in \{1, \dots, C\}$  is associated to each observation  $x_i$  of the learning set ( $i = 1, \dots, n$ ), we fall into the classical estimation framework of model-based discriminant analysis. In such a case, the maximization of the likelihood reduces to separately estimate the parameters of each class density by maximizing the associated conditional log-likelihood  $\mathcal{L}_k(\mathcal{X}; \Theta) = \sum_{i=1}^n s_{ik} \log(\pi_k f_k(x_i; \theta_k))$ , for  $k = 1, \dots, C$ ,

and this conduces to an estimation of  $\pi_k$  by  $\hat{\pi}_k = \frac{n_k}{n}$  where  $n_k = \sum_{i=1}^n s_{ik}$  is the number of observations of the  $k$ th class and to an estimation of  $\theta_k$  by  $\hat{\theta}_k$  which depends on the chosen component density. For instance, in the case of a Gaussian density, the maximization of  $\mathcal{L}_k(\mathcal{X}; \Theta)$  conduces to an estimation of  $\theta_k = \{\mu_k, \Sigma_k\}$  by  $\hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^n s_{ik} x_k$  and  $\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^n s_{ik} (x_i - \hat{\mu}_k)^t (x_i - \hat{\mu}_k)$ , for  $k = 1, \dots, C$ . We refer respectively to [12] and [11] for parameter estimation in the case of parsimonious and HD Gaussian models, and to [3] in the case of a mixture with a noise component.

**The discovery phase** Usually, in discriminant analysis, the classification phase consists only in assigning new unlabeled observations to one of known classes. However, in this work, it is assumed that some classes could not be observed during the learning phase. It is therefore necessary to search for new classes before to classify the new observations for avoiding the misclassification of observations from an unobserved class (by assigning them to one of the observed classes). Using the model and the notations introduced above, it remains to find  $K - C$  new classes in the set of  $n^*$  new unlabeled observations  $\mathcal{X}^* = \{x_1^*, \dots, x_{n^*}^*\}$ . Since these new observations are unlabeled, we have to fit the mixture model in a partially unsupervised way. In this case, the log-likelihood has the following form:

$$\ell(\mathcal{X}^*; \Theta) = \sum_{i=1}^{n^*} \left( \sum_{k=1}^C s_{ik}^* \log(\pi_k f_k(x_i^*; \theta_k)) + \sum_{k=C+1}^K s_{ik}^* \log(\pi_k f_k(x_i^*; \theta_k)) \right),$$

where the parameters  $\theta_k$  for  $k = 1, \dots, C$  have been estimated in the learning phase and the parameters  $\theta_k$  for  $k = C + 1, \dots, K$  remain to estimate. Due to the constraint  $\sum_{k=1}^K \pi_k = 1$  on the parameters  $\pi_k$ , the mixture proportions of the  $C$  known classes have to be re-normalized according to the proportions of the  $K - C$  new classes which will be estimated on the new sample  $\{x_1^*, \dots, x_{n^*}^*\}$ . However, the test set  $\{x_1^*, \dots, x_{n^*}^*\}$  is an incomplete dataset since the labels  $z_i^*$  are missing and the  $s_{ik}^*$  are consequently unknown for all observations of this set. In such a situation, the direct maximization of the likelihood is an intractable problem and the EM algorithm can be used to estimate the mixture parameters by iteratively maximizing the likelihood. We propose below a constrained EM algorithm for estimating the parameters of the  $K - C$  unobserved classes which alternates between the following E and M steps at each iteration  $q$ :

- **E step:** the conditional probabilities  $t_{ik}^{*(q)} = P(Z = k | X = x_i^*)$ , for  $i = 1, \dots, n^*$  and  $k = 1, \dots, K$ , are updated according to the mixture parameters as follows:

$$t_{ik}^{*(q)} = \frac{\hat{\pi}_k^{(q-1)} f_k(x_i^*; \hat{\theta}_k^{(q-1)})}{f(x_i^*; \hat{\Theta}^{(q-1)})},$$

where  $\hat{\pi}_k^{(q-1)}$  and  $\hat{\theta}_k^{(q-1)}$  are the mixture parameters estimated in the M step at step  $(q - 1)$ .

- **M step:** the parameters of the  $K - C$  unobserved classes are estimated by maximizing the conditional expectation of the completed likelihood whereas the estimated parameters of the observed classes remain fixed to the values obtained in the learning phase except for the proportions which are re-estimated. Therefore, this step only updates the estimates of parameters  $\pi_k$  for  $k = 1, \dots, K$  and  $\theta_k$  for  $k = C + 1, \dots, K$ . In the case of the Gaussian mixture, the update formulas for the parameter estimates are:

$$\begin{cases} \text{for } k = 1, \dots, C & \hat{\pi}_k^{(q)} = \left(1 - \sum_{\ell=C+1}^K \frac{n_\ell^{*(q)}}{n^*}\right) \frac{n_k}{n}, \\ \text{for } k = C + 1, \dots, K & \hat{\pi}_k^{(q)} = \frac{n_k^{*(q)}}{n^*} \end{cases}$$

where  $n_k^{*(q)} = \sum_{i=1}^{n^*} t_{ik}^{*(q)}$  and for  $k = C + 1, \dots, K$ :

$$\hat{\mu}_k^{(q)} = \frac{1}{n_k^{*(q)}} \sum_{i=1}^{n^*} t_{ik}^{*(q)} x_i^*, \quad \hat{\Sigma}_k^{(q)} = \frac{1}{n_k^{*(q)}} \sum_{i=1}^{n^*} t_{ik}^{*(q)} (x_i^* - \hat{\mu}_k^{(q)})(x_i^* - \hat{\mu}_k^{(q)})^t.$$

Proofs of these results are given in Appendix A.2.

### 3.4 Model selection: determining the number of components

Conversely to usual supervised classification, the total number  $K$  of classes is assumed to be unknown and has to be chosen. Therefore, this step is naturally a critical step in the search for unobserved classes. Classical tools for model selection in the mixture model framework are penalized likelihood criteria and include the AIC [1], BIC [37] and ICL [6] criteria. The Bayesian Information Criterion (BIC) is certainly the most popular and consists in selecting the model which maximizes the quantity:

$$BIC(\mathcal{M}) = \ell(x_1, \dots, x_n; \Theta) - \frac{\nu(\mathcal{M})}{2} \log(n),$$

where  $\nu(\mathcal{M})$  is the the number of parameters in model  $\mathcal{M}$  and  $n$  is the number of observations. The AIC criterion penalizes the log-likelihood by  $\nu(\mathcal{M})$  and the ICL criterion add the penalty  $\sum_{i=1}^n \sum_{k=1}^K t_{ik} \log(t_{ik})$  to the one of the BIC criterion in order to favour well separated models. The values of  $\gamma(\mathcal{M})$  and  $\nu$  are off course specific to the model proposed in this paper and depend on the chosen estimation procedure. For instance, if the classical Gaussian model is used within the transductive approach,  $\gamma(\mathcal{M})$  is equal to  $(K - 1) + Kp + Kp(p + 1)/2$  whereas it is equal to  $(K - 1) + (K - C)p + (K - C)p(p + 1)/2$  with the inductive approach. An evaluation of both criteria in the context of unobserved class detection is presented in the next section.

### 3.5 Classification with the adapted classifier

The previous paragraphs introduced a model-based discriminant analysis method which adapts its mixture model to a new situation including unobserved classes. Therefore, the adapted model can be used to classify new observations in the future. In the classical discriminant analysis framework, new observations are usually assigned to a class using the *maximum a posteriori* (MAP) rule. The MAP rule assigns a new observation  $x \in \mathbb{R}^p$  to the class for which  $x$  has the highest posterior probability. Therefore, the classification step mainly consists in calculating the posterior probability  $P(Z = k|X = x)$  for each class  $k = 1, \dots, K$ . In the case of the model described in this section, this posterior probability can be expressed classically using the Bayes' rule as follows:

$$P(Z = k|X = x) = \frac{\pi_k f_k(x; \theta_k)}{f(x; \Theta)},$$

where  $f(x; \Theta) = \sum_{k=1}^K \pi_k f_k(x; \theta_k)$ . Therefore, the posterior probabilities of the new observations depend on both the classes observed in the learning phase and the classes discovered in the test set.

## 4 Experimental results

This section presents experiments on toy and simulated datasets in order to highlight the main features of the method introduced in the previous section.

### 4.1 An introductory example: the iris dataset

The dataset considered in this first experiment is a classical one: the iris dataset made famous by its use by Fisher in [17] as an example for discriminant analysis. This dataset, in fact collected by Edgar Anderson [2] in the Gaspé Peninsula (Canada), is made of three classes corresponding to different species of iris (*setosa*, *versicolor* and *virginica*) among which the classes *versicolor* and *virginica* are difficult to discriminate (they are at least not linearly separable). The dataset consists of 50 samples from each of three species and four features were measured from each sample. The four measurements are the length and the width of sepal and petal. This dataset is used here as a toy dataset because of its popularity and its biological nature.

**Detection of one unobserved class** First, let suppose that botanists are studying iris species and have only observed the two species *setosa* and *versicolor*. For this experiment, the dataset has been randomly split into a learning dataset without *virginica* examples and a test dataset with several *virginica* examples. The top-left panel of Figure 1 shows what the botanists are supposed to have observed in the past. The top-center panel of the same

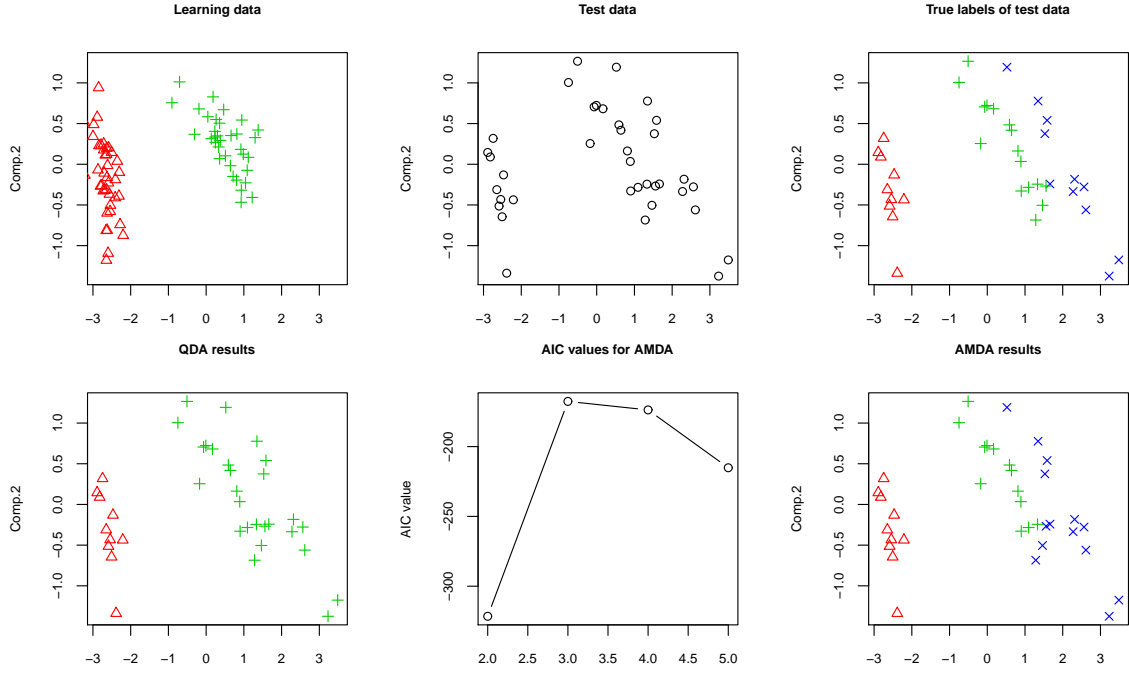


Figure 1: Detection of 1 unobserved class with AMDA on the Iris dataset: the classes “setosa” (red triangles) and “versicolor” (green plus-es) have been observed during the learning phase whereas the class “virginica” (blue crosses) has not.

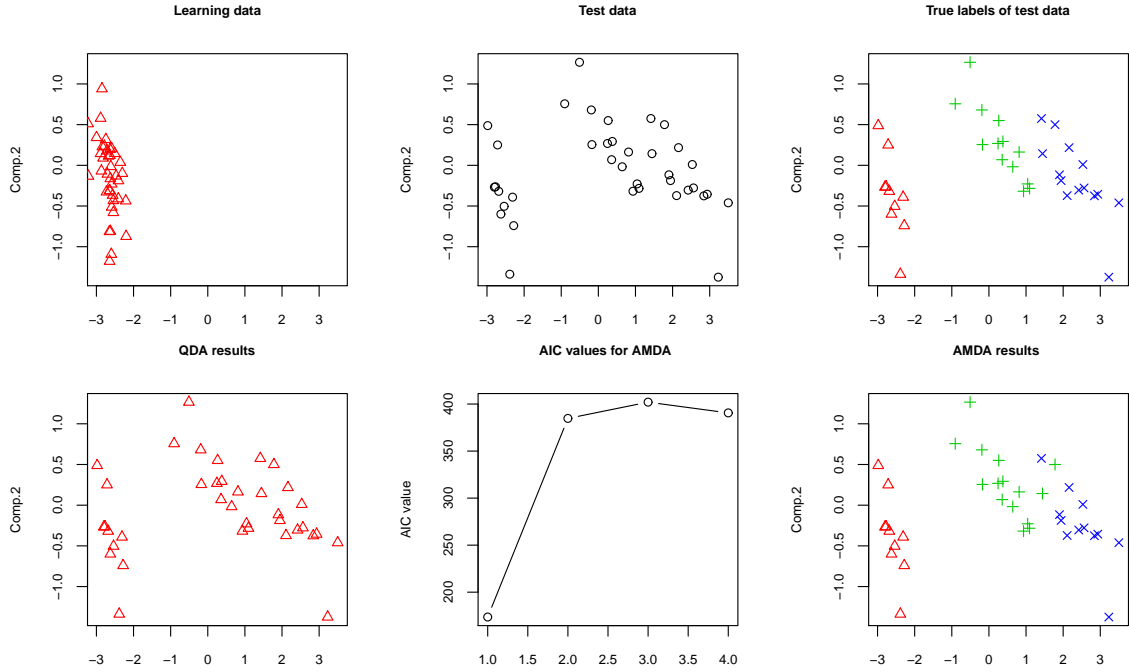


Figure 2: Detection of 2 unobserved classes with AMDA on the Iris dataset: the class “setosa” (red triangles) has been observed during the learning phase whereas the classes “versicolor” (green plus-es) and “virginica” (blue crosses) have not.

figure presents a sample of new observations of iris for which the botanists are asked to classify. However, as the top-right panel indicates, this new sample contains individuals from a class which has not been observed by the botanists in the past and the iris experts will very likely classify all these new observations as belonging to either the class *setosa* or the class *versicolor*. The bottom-left panel of Figure 1 shows the result of such a *scenario*, using Quadratic Discriminant Analysis (QDA) in place of the iris experts, which yields to the classification of all *virginica* observations in the class *versicolor*. Remark that, even though this result is disappointing from our point of view, it is understandable both for an human expert and a classification method since the classes *versicolor* and *virginica* are indeed very difficult to discriminate. The strategy proposed in the previous section, hereafter referred to by Adaptive Model-based Discriminant Analysis (AMDA), was applied to this dataset. The bottom-center and right panels of Figure 1 presents the results provided by AMDA (with the inductive approach). On the one hand, it turns out that model selection criteria (AIC here) succeed in identifying a new group of points in the test set. On the other hand, once the number  $K$  of mixture components chosen, AMDA classifies almost perfectly (only 2 errors on this example) the observations of the unobserved class *virginica*.

**Detection of two unobserved classes** Here, the toy example turns to be a serious problem because the botanists are now assumed to have only observed one species, the species *setosa*, and will have therefore to discover two unobserved classes, the species *versicolor* and *virginica*. For this second experiment, the dataset has been randomly split into a learning dataset without *versicolor* and *virginica* examples and a test dataset with several *versicolor* and *virginica* examples. The top-left panel of Figure 2 shows what the botanists are supposed to have observed in the past whereas the center panel shows the new and unlabeled observations. As one can observe, the new observations are clearly different from the data observed in the past but it is actually not obvious to detect that these new observations come from two different iris species (*cf.* top-right panel of Figure 2). If a supervised classifier like QDA is used, the classifier will assign all the new observations to the only one known class, the class *setosa*, and will make an important error (*cf.* bottom-left panel). In such a situation, there is no doubt that novelty detection methods presented in Section 2 are able to detect that the new observations do not belong to the species *setosa*. However, these techniques are not able to detect that the unlabeled observations are made of two homogeneous groups corresponding to two new iris species. The bottom-center and right panels of Figure 2 demonstrate that AMDA is actually able to detect the two unobserved iris species and can take this information into account to adapt the classifier for classifying future observations.

## 4.2 Detection of an unobserved noise class

This second experiment aims to evaluate the ability of AMDA to detect an unobserved non Gaussian class of noise. For this, data were simulated in a 2-dimensional space according

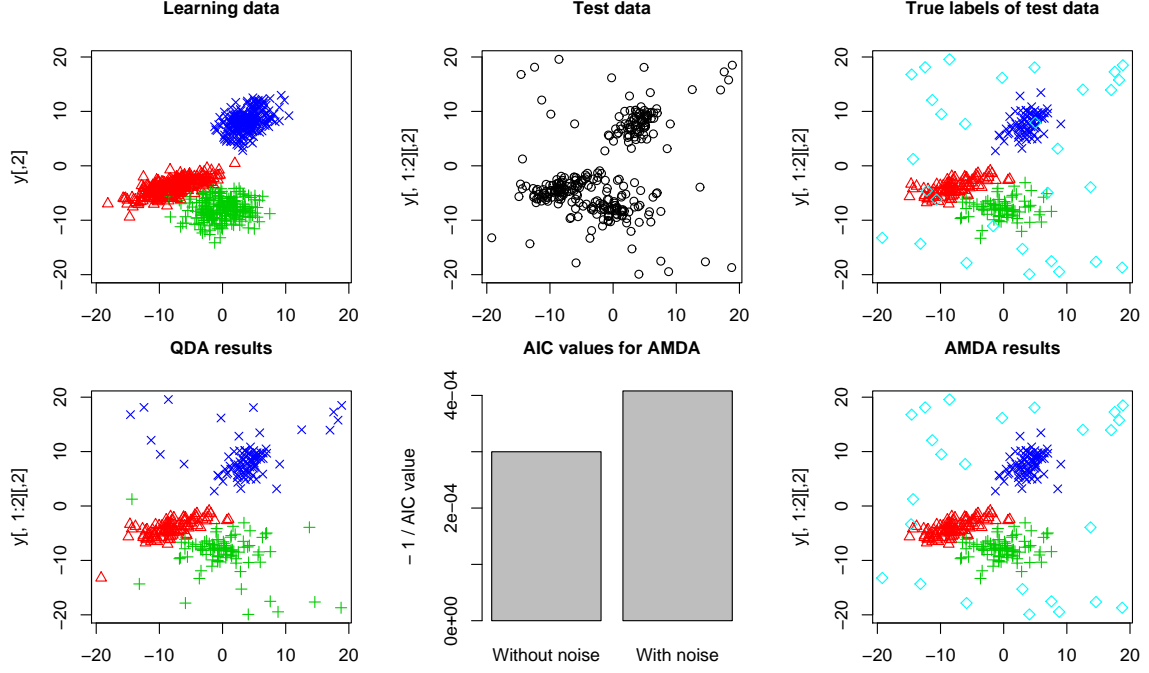


Figure 3: Detection of 1 unobserved noise class with AMDA on 2-dimensional simulated data: 3 observed classes and 1 unobserved noise class (light blue diamonds).

QDA				
Classif.	Truth			
	1	2	3	4
1	75	2		
2		78		10
3			65	21
4				
Correct classif. rate = 0.78				

AMDA				
Classif.	Truth			
	1	2	3	4
1	73	2		
2		77		
3			64	
4		1	1	31
Correct classif. rate = 0.9				

Table 1: Confusion tables for QDA and AMDA on the test dataset for the simulated data with one unobserved noise class (class #4).

a mixture model made of 4 components: 3 Gaussian components and one uniform noise component. Means and covariance matrices of Gaussians were chosen in order to obtain separated enough groups. The learning set was made of 750 observations from the three Gaussian classes. The top-left panel of Figure 3 shows the observations of the learning set. The test set was made of 250 observations from the three Gaussian classes (observed during the learning) and 31 observations from the unobserved uniform noise class. The top-center panel of this figure presents the unlabeled test observations and the top-left panel indicates the true labels of these observations. The bottom-left panel of Figure 3 shows the classification of the test observations with the supervised classifier Quadratic Discriminant Analysis (QDA). Unsurprisingly, QDA classifies all the observations from the noise class to one of the three known Gaussian classes. Table 1 presents confusion tables for QDA and AMDA (inductive

QDA-ND				
<i>Classif.</i>	<i>Truth</i>			
	1	2	3	4
1	66			
2		51		
3			54	
4	18	25	10	26
<i>Correct classif. rate = 0.78</i>				

SVM-ND				
<i>Classif.</i>	<i>Truth</i>			
	1	2	3	4
1	73			
2		68		
3			58	
4	11	8	6	26
<i>Correct classif. rate = 0.9</i>				

AMDA				
<i>Classif.</i>	<i>Truth</i>			
	1	2	3	4
1	83	5		
2	1	71		
3			64	
4				26
<i>Correct classif. rate = 0.98</i>				

Table 2: Confusion tables for QDA-ND, SVM-ND and AMDA on the test dataset for the simulated data with one unobserved class (class #4).

approach) on the test dataset and shows that all noise observations were classified into the classes #2 and #3. The bottom-center and right panels of Figure 3 show respectively the AIC values for AMDA with and without a noise component and the classification provided by AMDA with a detected noise component (as indicated by the largest AIC value). We can observe on this quite complex example that AMDA succeeds in both detecting the unobserved noise class and modeling it through a uniform component. Table 1 confirms that AMDA recognizes all noise observations as belonging to one unobserved class in the past and makes only 2 false noise detections which is very satisfying. Naturally, it could be also possible to detect both unobserved classes and a noise component by comparing AIC curves with and without a noise component for different numbers of Gaussian components.

### 4.3 Comparison with novelty detection methods

This experiment aims to compare, on simulated data, AMDA with two of the most efficient novelty detection methods. The two considered novelty detection methods are those proposed respectively by [41], referred to by QDA-ND in the sequel, and by [39], referred to by SVM-ND.

**Detection of one unobserved class** The first situation considered here is the detection of one unobserved class. For this first simulation, data were simulated according a Gaussian mixture model made of 4 groups (3 observed groups and one unobserved group) in a 2-dimensional space. Means and covariance matrices were chosen in order to obtain separated enough groups. Figure 4 presents the simulated learning dataset (top left panel), the test dataset (top center panel) and the true labels of the test observations (top right panel). The unobserved class is represented by light blue diamonds. The figure shows as well the classification results provided for the validation set by the three studied methods. Table 2 provides the confusion tables for QDA-ND, SVM-ND and AMDA as well as the associated correct classification rates. First, it appears that the three methods have successfully identified the new class since they both classified the 26 observations of the 4th class as novelties. On the one hand, QDA-ND and SVM-ND turn out to be too sensitive since they classify several observations belonging to the 3 observed classes as novelties. Their false positive rates are respectively 0.23 and 0.11. The high sensitivity of QDA-ND and SVM-ND partially explain

QDA-ND						SVM-ND						AMDA					
Classif.	Truth					Classif.	Truth					Classif.	Truth				
	1	2	3	4	5		1	2	3	4	5		1	2	3	4	5
1	48					1	55					1	61	1			
2		56				2		64				2		71			
3			56			3			65			3			69		2
4	13	16	15	24	22	4	6	8	6	24	22	4			1	24	
5						5						5			1		20
Correct classif. rate = 0.74						Correct classif. rate = 0.83						Correct classif. rate = 0.98					

Table 3: Confusion tables for QDA-ND, SVM-ND and AMDA on the test dataset for the simulated data with two unobserved classes (classes #4 and #5).

the very good performances of these two methods in terms of novelty detection. On the other hand, AMDA makes no false novelty detection. This satisfying result is mainly due to the fact that AMDA models all classes, including the unobserved one, before classifying new observations.

**Detection of two unobserved classes** This second situation is certainly the most realistic since there is no reason to limit the number of novelty classes to one. Off course, this situation will not be favorable to novelty detection methods which do not model unobserved classes. For this second simulation, data were simulated according a Gaussian mixture model made of 5 groups (3 observed groups and two unobserved groups) in a 2-dimensional space. Means and covariance matrices were again chosen in order to obtain separated enough groups. Figure 4 presents as before the simulated learning dataset, the test dataset and the true labels of the test observations. The two unobserved classes are respectively represented by light blue diamonds and purple triangles. The figure shows as well the classification results provided for the validation set by the three studied methods. Table 2 provides the confusion tables and the associated correct classification rates for QDA-ND, SVM-ND and AMDA. Unsurprisingly, QDA-ND and SVM-ND recognize the observations from the two unobserved classes as novelties but are unable to separate them into two different classes. In addition, the two novelty detection appear as before to be too sensitive since they make respectively 44 and 20 false novelty detections. Conversely, AMDA succeeds in identifying the two unobserved groups and makes only 1 false detection. This experiment has therefore highlighted the limitations of novelty detection methods and shown that AMDA can be considered, from this point of view, as their extension for the detection of multi-class novelties.

#### 4.4 Monte Carlo simulations

This paragraph presents Monte Carlo experiments on simulated data in order to both compare inductive and transductive approaches, evaluate model selection criteria in the context of unobserved class detection and determinate the breakdown group size for the detection of new classes. For the three following experiments, data were simulated according a Gaussian

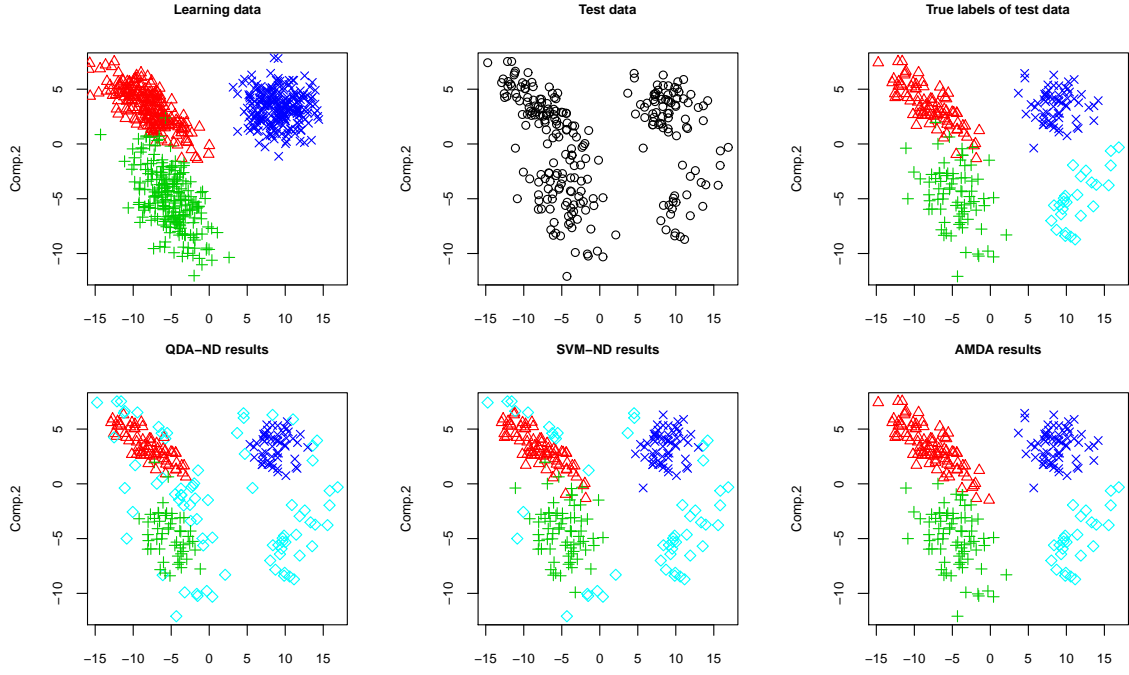


Figure 4: Detection of 1 unobserved class on 2-dimensional simulated data: 3 observed classes and 1 unobserved class (light blue diamonds).

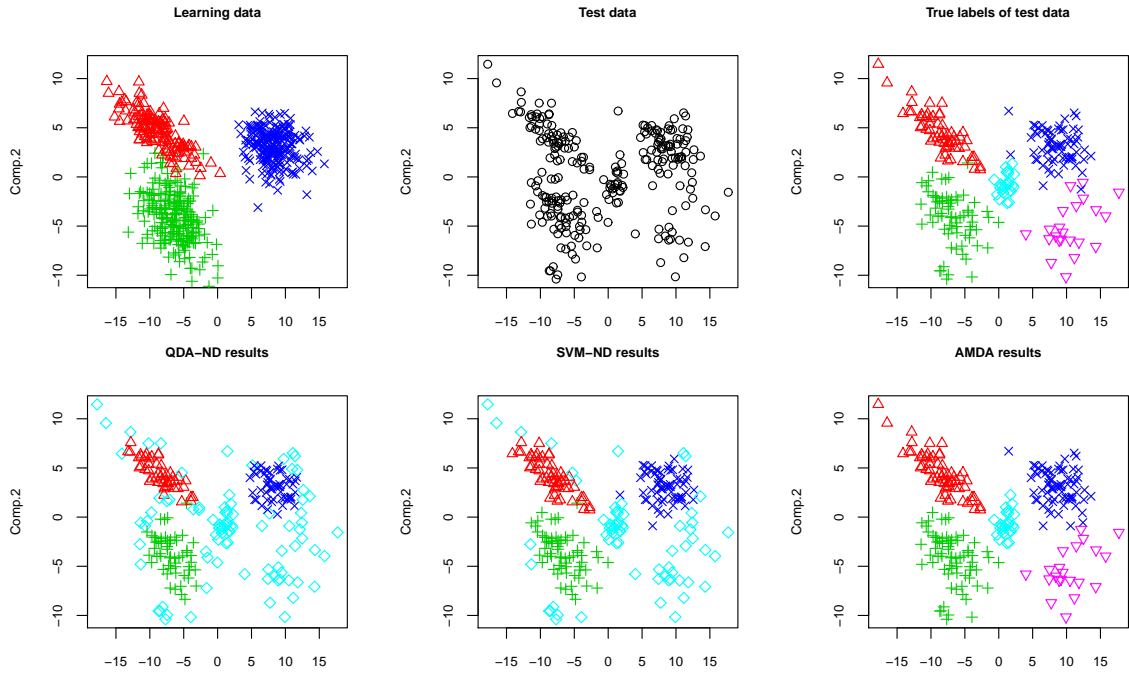


Figure 5: Detection of 2 unobserved classes on 2-dimensional simulated data: 3 observed classes and 2 unobserved classes (light blue diamonds and purple triangles).

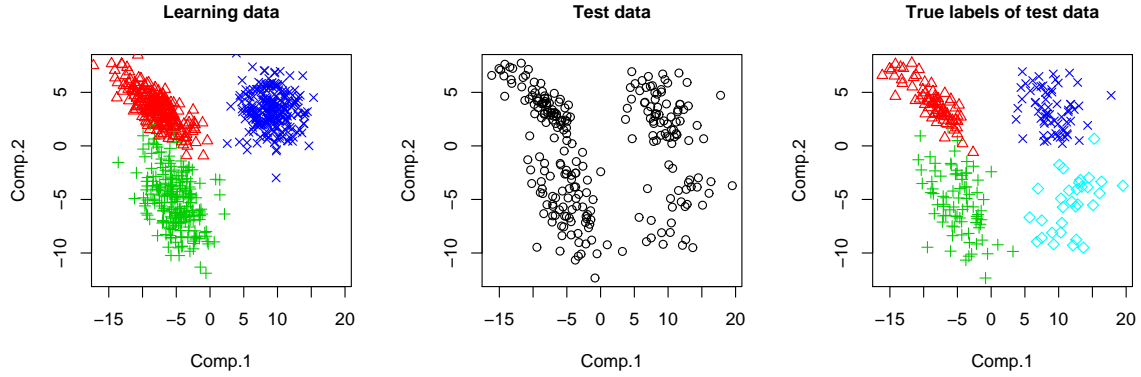


Figure 6: Classification with AMDA of simulated data: 3 observed classes and 1 unobserved class (light blue diamonds) in  $\mathbb{R}^2$ .

mixture model made of 4 groups (3 observed groups and one unobserved group) in a 2-dimensional space. Means and covariance matrices were chosen in order to obtain separated enough groups. Figure 6 presents the simulated learning dataset (left panel), the test dataset (center panel) and the true labels of the test observations (right panel). For each of the 50 replications of the Monte Carlo studies, 750 observations were simulated according to a Gaussian mixture model (250 obs. for each of the observed classes) in order to form the learning set and  $250 + \eta$  observations were simulated for the test set where  $\eta$  is the number of observations of the unobserved class. For each replication, the number  $\eta$  varied from 2 to 50.

**Inductive vs. transductive approaches** This first Monte-Carlo simulation aims to compare the inductive and transductive approaches proposed in Section 3. We choose to compare both approaches on modeling and classification criteria since supervised classification has two main objectives: efficiently classify new observations and correctly model the data for facilitating the understanding of classification results. For this simulation, the actual number of components was provided to the algorithms in order to focus on modeling and classification abilities of both approaches. The left panel of Figure 7 shows the log-likelihood value of the whole dataset (training + test sets) divided by the number of observations for the inductive and transductive approaches according to the size of the unobserved class. In this figure, the information in which we are interested in is the relative behaviour of the inductive approach compared to the transductive one and not the the global behaviour of both curves. Indeed, for each situation, the approach with the highest log-likelihood value per point is the one which provides the best modeling of the data. On the one hand, it appears that the log-likelihood curve of the transductive approach is above the one of inductive approach for sizes of the unobserved class larger than 10. This indicates that, for large unobserved groups of points, the use of the all available observations allows to better model the data than using only the test observations. On the other hand, Figure 7 indicates as well that for small unobserved groups

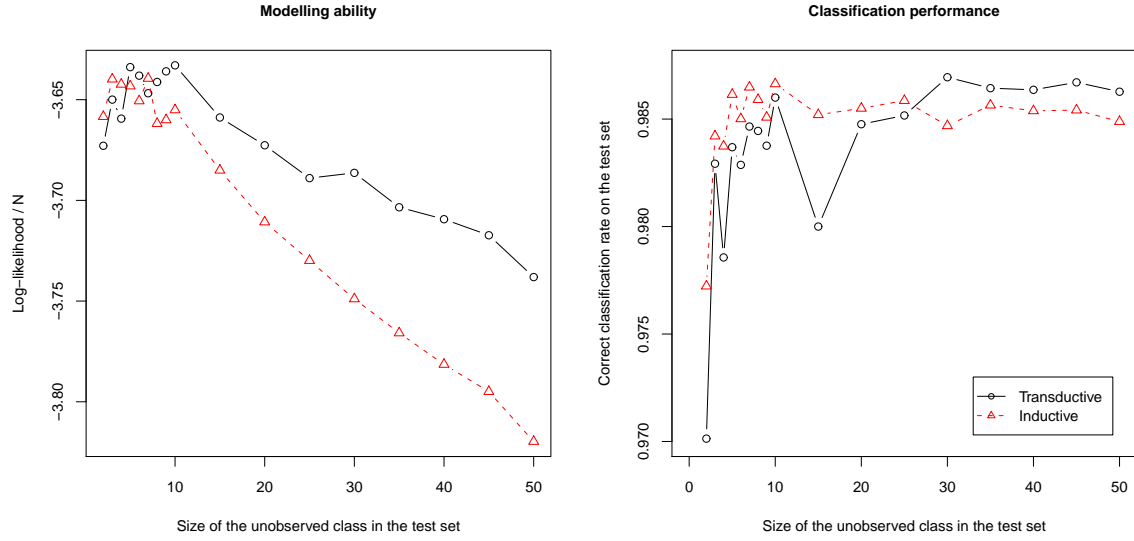


Figure 7: Modeling ability and classification performance of the inductive and transductive versions of AMDA according to the size of the unobserved class on simulated data: 3 observed classes and 1 unobserved class in  $\mathbb{R}^2$ .

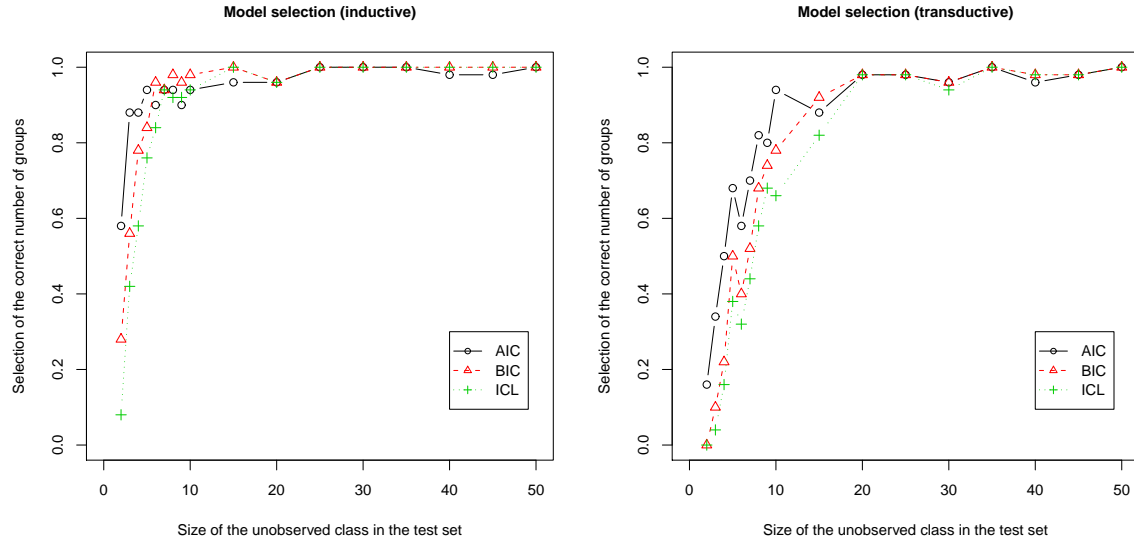


Figure 8: Successful selection of the actual number of groups using AIC, BIC and ICL with the inductive (left) and transductive (right) AMDA according to the size of the unobserved class on simulated data: 3 observed classes and 1 unobserved class in  $\mathbb{R}^2$ .

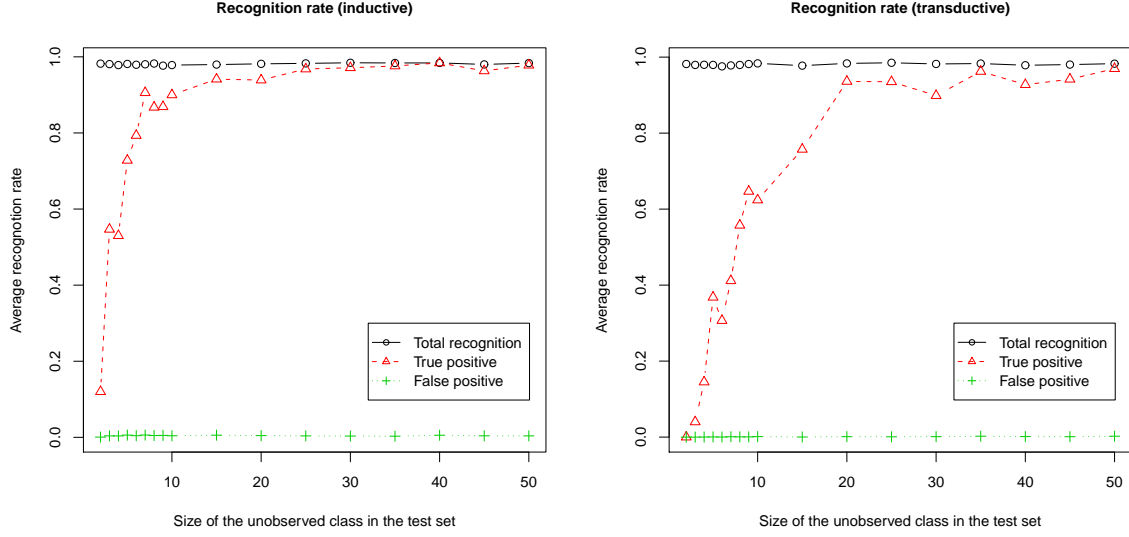


Figure 9: Recognition rates with the inductive (left) and transductive (right) AMDA according to the size of the unobserved class on simulated data: 3 observed classes and 1 unobserved class in  $\mathbb{R}^2$ .

(smaller than 10) the inductive approach seems to better model the data than the transductive version of AMDA. This can be explained by the will of the transductive approach to consider small unobserved groups of points as extreme values of the observed classes. The right panel of Figure 7 shows the correct classification on a second test data set (different from the test set used for detecting new classes) for the two studied approaches according to the size of the unobserved class. A test set different from the test set used for detecting new classes is used here in order to evaluate the ability of both approaches to classify future unlabeled data with the adapted classifier including the discovered classes. One can observe that both classification rates are very good (between 0.97 and 0.99) and that the inductive version of AMDA appears to be slightly more efficient and stable than the transductive one to classify new data with the adapted classifier. In view of this results, we can recommend to use the transductive version for modeling purpose on large datasets and to use the inductive approach for classification purpose or modeling of small datasets.

**Evaluation of model selection criteria** This second Monte Carlo study aims to evaluate model selection criteria in the context of unobserved class detection with AMDA. Figure 9 presents the rate of successful selection of the actual number of groups by the three selection model criteria AIC, BIC and ICL for both the inductive (left panel) and transductive (right panel) versions of AMDA. It appears that the three studied selection model criteria select always the correct number of groups when the unobserved group size is large (larger than 10 for the inductive approach and larger than 20 for the transductive one). For smaller sizes of the unobserved group, AIC turns out to be the more stable criterion since it selects the correct

number of groups more frequently than the two other criteria. We therefore recommend the use of AIC as model selection criterion in the context of unobserved class detection.

**Determination of the breakdown group size** The two panels of Figure 9 shows three recognition rates averaged on the Monte Carlo replications for both the inductive (left panel) and transductive (right panel) versions of AMDA: total recognition rate, true positive rate and false positive rate. The total recognition rate measures the overall correct classification rate for the four classes (the three observed classes and the unobserved one). The true positive rate measures the correct classification rate for observations of the unobserved class (class #4). Conversely, the false positive rate evaluates how many observations of the three observed classes are classified as belonging to the new class. In a satisfying situation, the total recognition rate and the true positive rate should be close to 1 whereas the false positive rate should be close to 0. Both recognition rates were computed on a test dataset. Figure 9 shows that the three recognition rates are very good for sizes of the unobserved class larger than 10 in the case of inductive AMDA (left panel) and larger than 20 in the case of transductive AMDA (right panel). We observe as well that for sizes of the unobserved class smaller than 5–10 the true positive rate is very unstable and this means that the unobserved class is not well modeled. This confirms the observation made in the previous paragraph and the inductive approach seems more robust than transductive AMDA in the case of unobserved classes of small sizes.

To summarize, these Monte Carlo experiments have first demonstrated that the transductive version of AMDA should be used for modeling purpose on large dataset whereas inductive AMDA should be used for detecting and modeling small unobserved classes. They have also shown that AIC is the most efficient criterion for detecting unobserved classes and that the inductive version of AMDA is able to detect and model unobserved classes in the test set for unobserved classes as small as 5–10 observations.

## 5 Application to social network analysis

Graph-structured networks are nowadays widely used to represent relationships between persons in organizations or communities. Recently, the need of classifying and visualizing such data has suddenly grown due to the emergence of internet and of a large number of social network websites. Indeed, increasingly, it is becoming possible to observe “network informations” in a variety of contexts, such as email transactions, connectivity of web pages, protein-protein interactions and social networking. A number of scientific goals can apply to such networks, ranging from unsupervised problems such as describing network structure, to supervised problems such as predicting node labels with information on their relationships. One of the main features of network data is that they are not “frozen” and can be observed over the time. Therefore, in the supervised context, the detection of unobserved communities in networks is

a major issue, for instance in commercial or security applications. In this section, AMDA is applied to the detection of unobserved communities in a real-world network.

### 5.1 The latent space model and its supervised version

Among existing probabilistic social network models, we choose to consider the latent space model proposed by [23] and extended in [18] to the unsupervised classification problem. This model provides probabilistic inference for the visualization and analysis of a social network. A social network is usually represented by a  $n \times n$  socio-matrix  $Y$  where its elements  $Y_{ij}$  indicates an existing relation ( $Y_{ij} = 1$ ) or not ( $Y_{ij} = 0$ ) between the nodes  $i$  and  $j$ , for  $i, j = 1, \dots, n$ . For example, later in the section, we consider data in which  $Y_{ij} = 1$  indicates friendship between individuals  $i$  and  $j$ . The latent space model assumes the the probability of a tie between two nodes mainly depends on the distance between the nodes in a latent space. The model has the following form:

$$\text{logit}(P(Y_{ij} = 1|\theta)) = \alpha - \|Z_i - Z_j\|,$$

where  $\text{logit}(P) = \log(P/(1 - P))$ ,  $\theta = \{\alpha, Z\}$  are the parameters of the model,  $\alpha$  determines the prior probability of an existing link between two nodes and  $Z_i$  is the position of the  $i$ th node in the  $p$ -dimensional latent space. Thus, using this model, nodes  $i$  and  $j$  have a high probability to be connected if  $\alpha$  is large or if they are close in the latent space, *i.e.*  $\|Z_i - Z_j\|$  is close to 0. To learn the latent space model, parameters  $\alpha$  and  $Z_1, \dots, Z_n$  have to be estimated for a fixed value of the latent space dimension  $p$  which can be chosen by cross-validation or using a criterion such as BIC. Parameter estimation can be done by iterative likelihood maximization or MCMC techniques (see [23] for details). Recently, an extension of this model for the supervised classification problem has been proposed in [9]. The main idea of the approach is to introduce the supervised information within the latent space model through a covariate term  $\beta X_{ij}$ . The supervised latent model has therefore the following form:

$$\text{logit}(P(Y_{ij} = 1|\theta)) = \alpha - \beta X_{ij} - \|Z_i - Z_j\|,$$

where  $X_{ij}$  is equal to 1 if the nodes  $i$  and  $j$  are in the same class and  $-1$  if they are not. The parameter  $\beta$  is an hyper-parameter which tunes the importance given to the supervision in the model. Inclusion of  $\beta X_{ij}$  forces the model to provide latent positions which respect the class memberships. Once the latent model parameters are estimated, it is possible to learn any supervised classifier in the latent space. Afterward, new nodes can be projected into the learned latent space using their observed links with learning nodes and classified according to their latent position. We refer to [9] for technical details on the projection and the classification of new nodes in the supervised latent space.

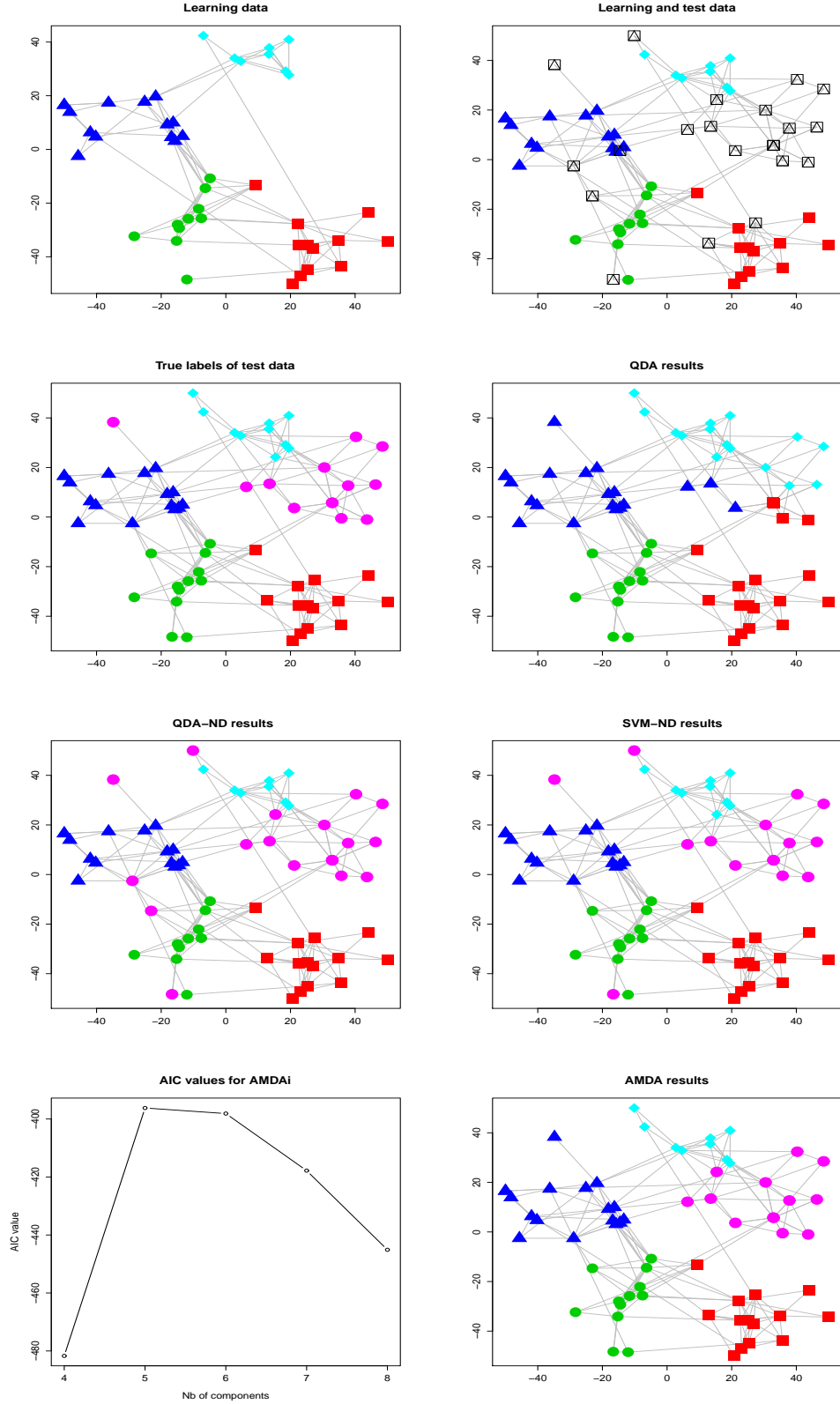


Figure 10: Detection of 1 unobserved community on the Add-Health network: 4 observed classes and 1 unobserved class (purple circles in the top-right corner). See text for details.

QDA					
<i>Classif.</i>	<i>Truth</i>				
	1	2	3	4	5
1	2				3
2		2			
3			2		4
4				2	5
5					
<i>Correct classif. rate = 0.40</i>					
QDA-ND					
<i>Classif.</i>	<i>Truth</i>				
	1	2	3	4	5
1	2				
2					
3			1		
4					
5		2	2	2	12
<i>Correct classif. rate = 0.60</i>					
SVM-ND					
<i>Classif.</i>	<i>Truth</i>				
	1	2	3	4	5
1	2				
2		1			
3			2		
4				1	
5		1		1	12
<i>Correct classif. rate = 0.90</i>					
AMDA					
<i>Classif.</i>	<i>Truth</i>				
	1	2	3	4	5
1	2				
2		2			
3			2		
4				1	1
5				1	11
<i>Correct classif. rate = 0.90</i>					

Table 4: Confusion tables for QDA, QDA-ND, SVM-ND and AMDA on the test dataset for the Add-Health dataset with one unobserved class (class #5).

## 5.2 The Add-Health dataset

The social network studied here is from the National Longitudinal Study of Adolescent Health and it is a part of a big dataset, usually called the “Add-Health” dataset. The data were collected in 1994-95 within 80 high-schools and 52 middle schools in the USA. The whole study is detailed in [20] and [42]. In addition to personal and social information, each student was asked to nominate his best friends. We consider here the social network based on the answers of 67 students from a single school, treating the grade of each student as the class variable. Two adolescents who nominated nobody were removed from the network. We therefore consider a whole dataset made of 65 students distributed into 5 classes: grade 7 to grade 11.

## 5.3 Detection of one unobserved community

Among the 65 nodes of the network, we randomly selected a learning dataset of 55 nodes and a test dataset of 20 nodes among which all the 13 students in grade 11 (5th class). For facilitating the visualization of the results, the latent space dimension  $p$  has been fixed to 2 for this experiment and the following one. Once the supervised latent space learned and the test nodes projected into the latent space, QDA, QDA-ND, SVM-ND and AMDA have been applied within the latent space for classifying the test nodes and trying to detect unobserved classes. Figure 10 presents the classification results obtained in the latent space with the four considered methods. The panels of the first row respectively present the learning nodes organized into 4 (observed) groups and the 20 test nodes for which the class membership is assumed to be unknown. The left panel of the second row provides the actual labels of

QDA					
<i>Classif.</i>	<i>Truth</i>				
	1	2	3	4	5
1	2			13	4
2		3		1	4
3			2		5
4					
5					
<i>Correct classif. rate = 0.21</i>					

QDA-ND					
<i>Classif.</i>	<i>Truth</i>				
	1	2	3	4	5
1	2				
2		1			
3			2		
4		2		14	13
5					
<i>Correct classif. rate = 0.56</i>					

SVM-ND					
<i>Classif.</i>	<i>Truth</i>				
	1	2	3	4	5
1					
2					
3					
4	2	3	2	14	13
5					
<i>Correct classif. rate = 0.41</i>					

AMDA					
<i>Classif.</i>	<i>Truth</i>				
	1	2	3	4	5
1	2				
2		3			
3			2		
4				14	1
5					12
<i>Correct classif. rate = 0.97</i>					

Table 5: Confusion tables for QDA, QDA-ND, SVM-ND and AMDA on the test dataset for the Add-Health dataset with two unobserved classes (classes #1 and #5).

all the nodes. The three following panels show the classification results provided by QDA, QDA-ND and SVM-ND respectively. Finally, the last row presents the results of AMDA (inductive version). Table 4 provides the confusion tables for the four classification methods computed on the test dataset. First, it appears that the different classes of the network are quite homogeneous and that the network has a natural structure. Indeed, when considering the full network (left panel of 2nd row), only two nodes have latent positions which do not agree with the class membership of the nodes. Regarding the classification results, QDA classifies unsurprisingly all the test nodes between the known classes and does not detect the unobserved class (light blue diamonds). QDA-ND and SVM-ND both detect the presence of novelties but their recognition rates are very unsatisfying. Indeed, the true positive and false positive rates are respectively equal to 1 and 0.85 for QDA-ND, and to 0.07 and 0 for SVM-ND. We recall that the optimal result is simultaneously 1 for the true positive rate and 0 for the false positive rate. Therefore, both QDA-ND and SVM-ND are clearly not optimal in this situation. Conversely, the true positive and false positive rates are respectively equal to 0.92 and 0 for AMDA. This means that AMDA has missed only one nodes from the new class without false positive detections. Furthermore, the missed node could be considered has an outlier for the new class when considering both its latent position or its friendship links.

#### 5.4 Detection of two unobserved communities

This second study focuses on the detection of two unobserved communities in the Add-Health network. The experimental setup is the same as before except that the test set contains 34 nodes among which all the nodes from the two unobserved classes. Figure 11 presents the classification results obtained in the latent space with the four considered methods whereas

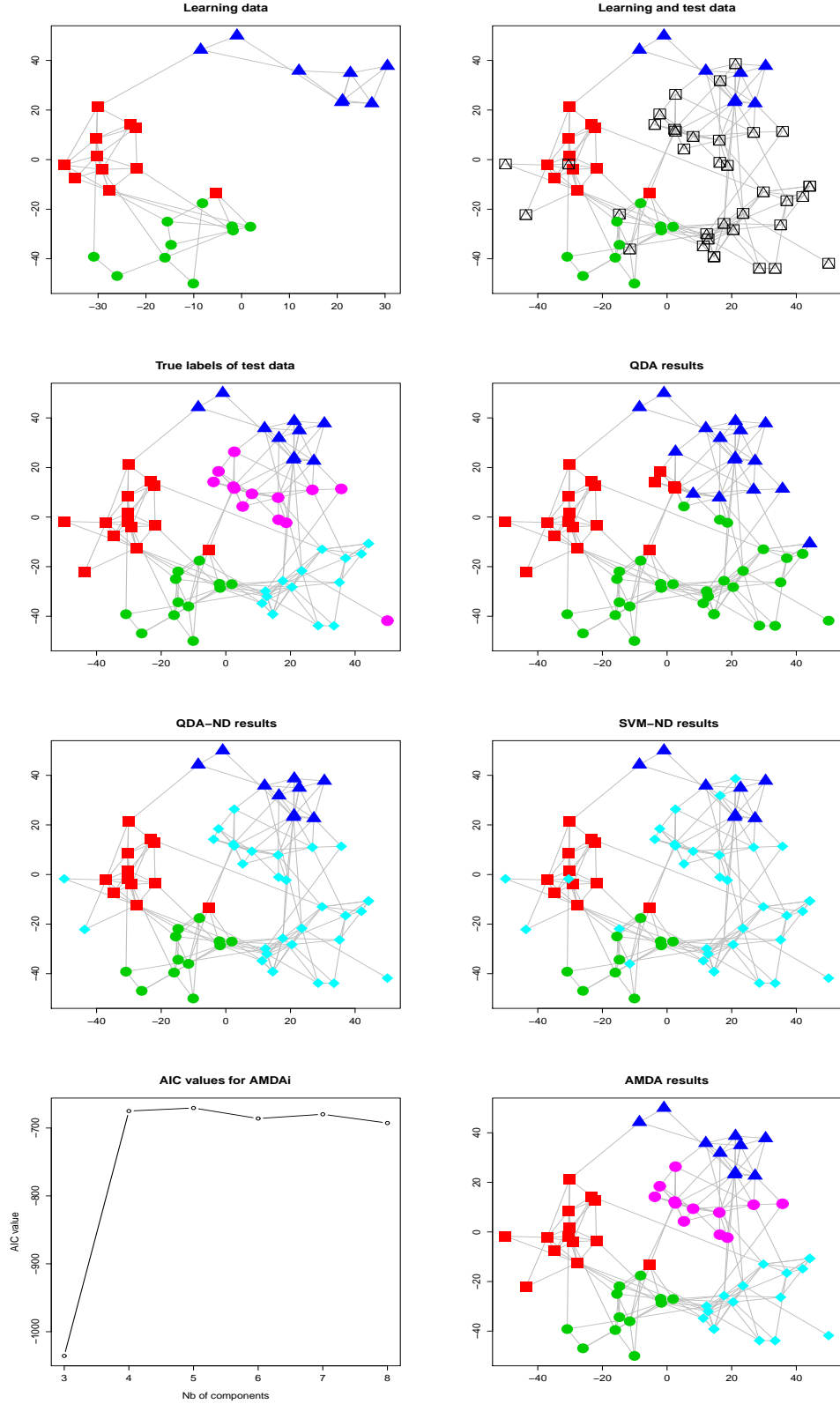


Figure 11: Detection of 2 unobserved communities on the Add-Health network: 3 observed classes and 2 unobserved classes (purple circles and light blue diamonds). See text for details.

Table 5 provides their confusion tables computed on the test dataset. As one can expect, QDA does not detect the unobserved classes (light blue diamonds and purple points) and assigns all nodes from these two classes to known classes. QDA-ND and SVM-ND both detect the presence of novelties but are unable to see the two new communities. QDA-ND and SVM-ND appear again too sensitive since they make several false novelty detection. Finally, AMDA both succeeds in detecting novelties and identifies two unobserved communities. As previously, AMDA fails in assigning the bottom-right node to its actual class but this is again understandable. This second experiment is particularly important since it demonstrates that AMDA is able to detect and model several unobserved classes in a real-world problem.

## 6 Conclusion and further works

This work has focused on the problem of learning a supervised classifier with unobserved classes. An adaptive model-based discriminant analysis method has been presented in this paper which is able to both detect unobserved groups of points in a new set of observations and to adapt the supervised classifier to the new situation. Two EM-based approaches have been proposed for parameter estimation: a transductive approach which considers all available observations for learning in a unique step and an inductive approach, which is made of a learning and a discovering phase. The detection of the number of unobserved classes is done using model selection criteria. Experiments on simulated and real datasets have shown that the proposed method is able to detect different kinds of unobserved classes (Gaussian, uniform noise, ...) and overcomes the drawbacks of novelty detection methods which are unable to detect several unobserved classes. AMDA has also been applied with success to the detection of unobserved communities in social network analysis.

It remains however to deal in the future with the problem of label switching when  $C - K > 1$ . A way to solve this problem could be to ask domain experts to classify some observations of the new detected groups in order to associate a class name with the detected groups. Parsimonious Gaussian models could be used as well for modeling small groups in order to detect unobserved groups of points smaller than 5-10 observations. In the same way, a mixture of Gaussians could be used to model each class in order to model more precisely the data. Finally, it could be very interesting to study the evolution of the proposed strategy in the context of dynamic classification.

## A Appendix

This ultimate section presents the proofs of parameter estimators given in Section 3 for both transductive and inductive approaches.

## A.1 Transductive approach

At the iteration  $q$  of the M step, the expectation of the complete log-likelihood  $Q(\mathcal{X}, \mathcal{X}^*; \Theta)$  conditionally to the posterior probabilities  $t_{ik}^*$  has the following form:

$$Q(\mathcal{X}, \mathcal{X}^*; \Theta) = \sum_{i=1}^n \sum_{k=1}^K \tilde{s}_{ik} \log(\pi_k f_k(x_i; \theta_k)) + \sum_{i=1}^{n^*} \sum_{k=1}^K t_{ik}^* \log(\pi_k f_k(x_i^*; \theta_k)),$$

where  $\log(\pi_k f_k(x_i; \theta_k))$  is given above. We recall that  $\tilde{s}_{ik} = s_{ik}$  if  $k = 1, \dots, C$  and  $\tilde{s}_{ik} = 0$  for  $k = C + 1, \dots, K$  and that for  $i = 1, \dots, n$ .

**ML estimator for parameter  $\pi_k$**  The maximization of  $Q(\mathcal{X}, \mathcal{X}^*; \Theta)$  according to the mixture proportion  $\pi_k$  under the constraint  $\sum_{k=1}^K \pi_k = 1$  is equivalent to find a saddle point of the Lagrangian  $\mathcal{L}(\Theta, \omega)$ :

$$\mathcal{L}(\Theta, \omega) = Q(\Theta) - \omega \left( \sum_{k=1}^K \pi_k - 1 \right),$$

where  $\omega$  is the Lagrangian coefficient. The partial derivative of the Lagrangian  $\mathcal{L}(\Theta, \omega)$  according to  $\pi_k$  is:

$$\frac{\partial}{\partial \pi_k} \mathcal{L}(\Theta, \omega) = \frac{(n_k + n_k^*)}{\pi_k} + \omega,$$

where  $n_k = \sum_{i=1}^n \tilde{s}_{ik}$  and  $n_k^* = \sum_{i=1}^{n^*} t_{ik}^*$ . The relation  $\frac{\partial}{\partial \pi_k} \mathcal{L}(\Theta, \omega) = 0$  implies that, for all  $k = 1, \dots, K$ :

$$(n_k + n_k^*) + \omega \pi_k = 0, \tag{2}$$

and summing up this quantity over  $k$  provides the value of the Lagrangian coefficient  $\omega$ :

$$\omega = -(n + n^*),$$

where  $n = \sum_{k=1}^K n_k$  and  $n^* = \sum_{k=1}^K n_k^*$ . Finally, replacing  $\omega$  by its value in (2) allows to find the ML estimate of  $\pi_k$ :

$$\hat{\pi}_k = \frac{(n_k + n_k^*)}{(n + n^*)}.$$

**ML estimator for parameter  $\mu_k$**  The partial derivative of  $Q(\mathcal{X}, \mathcal{X}^*; \Theta)$  according to  $\mu_k$  has the following form:

$$\frac{\partial}{\partial \mu_k} Q(\mathcal{X}, \mathcal{X}^*; \Theta) = -\Sigma_k^{-1} \left( \sum_{i=1}^n \tilde{s}_{ik} (x_i - \mu_k) + \sum_{i=1}^{n^*} t_{ik}^* (x_i^* - \mu_k) \right).$$

The relation  $\frac{\partial}{\partial \mu_k} Q(\mathcal{X}, \mathcal{X}^*; \Theta) = 0$  implies that:

$$\sum_{i=1}^n \tilde{s}_{ik}(x_i - \mu_k) + \sum_{i=1}^{n^*} t_{ik}^*(x_i^* - \mu_k) = 0,$$

which is equivalent to:

$$\mu_k \left( \sum_{i=1}^n \tilde{s}_{ik} + \sum_{i=1}^{n^*} t_{ik}^* \right) = \sum_{i=1}^n \tilde{s}_{ik} x_i + \sum_{i=1}^{n^*} t_{ik}^* x_i^*,$$

and this finally yields to the ML estimate of  $\mu_k$ :

$$\hat{\mu}_k = \frac{1}{n_k + n_k^*} \left( \sum_{i=1}^n \tilde{s}_{ik} x_i + \sum_{i=1}^{n^*} t_{ik}^* x_i^* \right),$$

where  $n_k = \sum_{i=1}^n \tilde{s}_{ik}$  and  $n_k^* = \sum_{i=1}^{n^*} t_{ik}^*$ .

**ML estimator for parameter  $\Sigma_k$**  At the optimum for parameter  $\mu_k$ , the partial derivative of  $Q(\mathcal{X}, \mathcal{X}^*; \Theta)$  according to  $\Sigma_k$  has the following form:

$$\begin{aligned} \frac{\partial}{\partial \Sigma_k} Q(\mathcal{X}, \mathcal{X}^*; \Theta) = & -\frac{1}{2} \frac{\partial}{\partial \Sigma_k} \left[ \sum_{i=1}^n \tilde{s}_{ik} (\log(|\Sigma_k|) + (x_i - \hat{\mu}_k)^t \Sigma_k^{-1} (x_i - \hat{\mu}_k)) \right. \\ & \left. + \sum_{i=1}^{n^*} t_{ik}^* (\log(|\Sigma_k|) + (x_i^* - \hat{\mu}_k)^t \Sigma_k^{-1} (x_i^* - \hat{\mu}_k)) \right]. \end{aligned}$$

Using the classical trick of the trace of the  $1 \times 1$  matrix, we can write that  $(x_i - \hat{\mu}_k)^t \Sigma_k^{-1} (x_i - \hat{\mu}_k) = \text{tr}((x_i - \hat{\mu}_k)^t \Sigma_k^{-1} (x_i - \hat{\mu}_k))$  and, using the identity  $\text{tr}(AB) = \text{tr}(BA)$ , we get:

$$\frac{\partial}{\partial \Sigma_k} Q(\mathcal{X}, \mathcal{X}^*; \Theta) = -\frac{1}{2} \frac{\partial}{\partial \Sigma_k} [(n_k + n_k^*) \log(|\Sigma_k|) + \text{tr}(\Sigma_k^{-1} S_k) + \text{tr}(\Sigma_k^{-1} S_k^*)],$$

where  $S_k = \sum_{i=1}^n \tilde{s}_{ik}(x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^t$  and  $S_k^* = \sum_{i=1}^{n^*} t_{ik}^*(x_i^* - \hat{\mu}_k)(x_i^* - \hat{\mu}_k)^t$ . Using the additivity property of the trace of square matrices, we end up with:

$$\frac{\partial}{\partial \Sigma_k} Q(\mathcal{X}, \mathcal{X}^*; \Theta) = -\frac{1}{2} \frac{\partial}{\partial \Sigma_k} [(n_k + n_k^*) \log(|\Sigma_k|) + \text{tr}(\Sigma_k^{-1} (S_k + S_k^*))].$$

Finally, using the matrix derivative formula of the logarithm of a determinant,  $\frac{\partial}{\partial A} \log(|A|) = (A^{-1})^t$ , and of the trace of a product,  $\frac{\partial}{\partial A} \text{tr}(A^{-1}B) = -(A^{-1}BA^{-1})^t$ , the equality of  $\frac{\partial}{\partial \Sigma_k} Q(\mathcal{X}, \mathcal{X}^*; \Theta)$  to the  $p \times p$  zero matrix yields to the relation:

$$(n_k + n_k^*) \Sigma_k^{-1} = \Sigma_k^{-1} (S_k + S_k^*) \Sigma_k^{-1},$$

and, by multiplying on the left and on the right by  $\Sigma_k$ , we find out the ML estimate of  $\Sigma_k$ :

$$\hat{\Sigma}_k = \frac{1}{(n_k + n_k^*)} (S_k + S_k^*).$$

## A.2 Inductive approach

At the iteration  $q$  of the M step, the expectation of the completed log-likelihood  $Q(\mathcal{X}^*; \Theta)$  conditionally to the posterior probabilities  $t_{ik}^*$  has the following form:

$$Q(\mathcal{X}^*; \Theta) = \sum_{i=1}^{n^*} \left( \sum_{k=1}^C t_{ik}^* \log(\pi_k f_k(x_i^*; \theta_k)) + \sum_{k=C+1}^K t_{ik}^* \log(\pi_k f_k(x_i^*; \theta_k)) \right),$$

where  $\log(\pi_k f_k(x_i; \theta_k))$  can be written as follows in the case of the multivariate Gaussian density:

$$\log(\pi_k f_k(x_i; \theta_k)) = -\frac{1}{2} (\log(\pi_k) + \log(|\Sigma_k|) + (x_i - \mu_k)^t \Sigma_k^{-1} (x_i - \mu_k)) + C,$$

where  $C = -p \log(2\pi)/2$  is a constant which does not depend on mixture parameters. In the case of the discovery phase of the inductive approach, the maximization of  $Q(\mathcal{X}^*; \Theta)$  according to the parameters  $\mu_k$  and  $\Sigma_k$  can be done classically except that parameters  $\mu_k$  and  $\Sigma_k$  have only to be estimated for  $k = C + 1, \dots, K$ . We therefore refer to [31] for ML inference for  $\mu_k$  and  $\Sigma_k$  in finite mixture models.

The estimation of the mixture proportions  $\pi_k$  can unfortunately not be done classically and must be done sequentially. On the one hand, for  $k = C + 1, \dots, K$ , the maximization of  $Q(\mathcal{X}, \mathcal{X}^*; \Theta)$  according to the mixture proportion  $\pi_k$  under the constraint  $\sum_{k=1}^K \pi_k = 1$  allows to find the ML estimate for  $\pi_k$ :

$$\forall k = C + 1, \dots, K, \quad \hat{\pi}_k = \frac{n_k^*}{n^*},$$

where  $n_k^* = \sum_{i=1}^{n^*} t_{ik}^*$ . On the other hand, ML estimate of  $\pi_k$  for  $k = 1, \dots, C$  must be updated such that the constraint  $\sum_{k=1}^K \pi_k = 1$  still holds. For this, estimators  $\hat{\pi}_1, \dots, \hat{\pi}_C$  can be updated as follows:

$$\forall k = 1, \dots, C, \quad \hat{\pi}_k = \left( 1 - \sum_{\ell=C+1}^K \frac{n_\ell^*}{n^*} \right) \frac{n_k}{n},$$

where  $n_k = \sum_{i=1}^n s_{ik}$ .

## References

- [1] H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723, 1974.

- [2] E. Anderson. The irises of the gaspé peninsula. *Bulletin of the American Iris Society*, 59:2–5, 1935.
- [3] J. Banfield and A. Raftery. Model-based Gaussian and non-Gaussian clustering. *Biometrics*, 49:803–821, 1993.
- [4] R. Bellman. *Dynamic Programming*. Princeton University Press, 1957.
- [5] H. Bensmail and G. Celeux. Regularized Gaussian discriminant analysis through eigenvalue decomposition. *Journal of the American Statistical Association*, 91:1743–1748, 1996.
- [6] C. Biernacki, G. Celeux, and G. Govaert. Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(7):719–725, 2000.
- [7] C. Bishop. Novelty detection and neural network validation. In *IEES Conference on Vision and Image Signal Processing*, pages 217–222, 1994.
- [8] A. Blum and T. Mitchell. Combining labeled and unlabeled data with co-training. In *Conference on Computational Learning Theory*, 1998.
- [9] C. Bouveyron, H. Chipman, and E. Côme. Supervised classification and visualization of social networks based on a probabilistic latent space model. In *7th International Workshop on Mining and Learning with Graphs*, Leuven, Belgium, 2009.
- [10] C. Bouveyron, S. Girard, and C. Schmid. High-Dimensional Data Clustering. *Computational Statistics and Data Analysis*, 52(1):502–519, 2007.
- [11] C. Bouveyron, S. Girard, and C. Schmid. High Dimensional Discriminant Analysis. *Communications in Statistics: Theory and Methods*, 36(14):2607–2623, 2007.
- [12] G. Celeux and G. Govaert. Gaussian parsimonious clustering models. *Pattern Recognition*, 28(5):781–793, 1995.
- [13] C. Chow. On optimum recognition error and reject tradeoff. In *IEEE Transactions on Information Theory*, pages 41–46, 1970.
- [14] D. Dasgupta and F. Nino. A comparison of negative and positive selection algorithms in novel pattern detection. In *IEEE International Conference on Systems and Cybernetics*, pages 125–130, 2000.
- [15] A. Dempster, N. Laird, and D. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39(1):1–38, 1977.

- [16] M. Desforges, P. Jacob, and J. Cooper. Applications of probability density estimation to the detection of abnormal conditions in engineering. In *Proc. Institute of Mechanical Engineers*, pages 687–703.
- [17] R. Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7:179–188, 1936.
- [18] M. Handcock, A. Raftery, and J. Tantrum. Model-based clustering for social networks. *Journal of the Royal Statistical Society, Series A*, 170(2):1–22, 2007.
- [19] L. Hansen, C. Liisberg, and P. Salamon. The error-reject tradeoff. *Open Systems and Information Dynamics*, 4:159–184, 1997.
- [20] K. Harris, F. Florey, J. Tabor, P. Bearman, J. Jones, and R. Udry. The national longitudinal of adolescent health: Research design. Technical report, Carolina Population Center, University of North Carolina, 2003.
- [21] M. Hellman. The nearest neighbour classification with a reject option. In *IEEE Transactions on Systems Science and Cybernetics*,, pages 179–185, 1970.
- [22] C. Hennig and P. Coretto. *The Noise Component in Model-based Cluster Analysis*, pages 127–138. Data Analysis, Machine Learning and Applications. Springer, 2008.
- [23] P. Hoff, A. Raftery, and M. Handcock. Latent spaces approaches to social network analysis. *Journal of the American Statistical Association*, 97(460):1090–1098, 2002.
- [24] T. Kohonen. *Self-organisation and associative memory*. Springer-verlag, berlin edition, 1988.
- [25] B. Krishnapuram, D. Williams, Y. Xue, A. Hartemink, L. Carin, and M. Figueiredo. On semi-supervised classification. In *NIPS*, 2004.
- [26] C. Manikopoulos and S. Papavassiliou. Network intrusion and fault detection: a statistical anomaly approach. *rk intrusion and fault detection: a IEEE Communications Magazine*, 40(10):76–82, 2002.
- [27] M. Markou and S. Singh. Novelty detection: A review - part 1: Statistical approaches. *Signal Processing*, 83(12):2481–2497, 2003.
- [28] M. Markou and S. Singh. Novelty detection: A review - part 2: Neural network based approaches. *Signal Processing*, 83(12):2499–2521, 2003.
- [29] G. McLachlan. Iterative reclassification procedure for constructing an asymptotically optimal rule of allocation in discriminant analysis. *Journal of the American Statistical Association*, (70):365–369, 1975.

- [30] G. McLachlan. *Discriminant Analysis and Statistical Pattern Recognition*. Wiley, New York, 1992.
- [31] G. McLachlan and D. Peel. *Finite Mixture Models*. Wiley, New York, 2000.
- [32] T. Odin and D. Addison. Novelty detection using neural network technology. In *Proc. of COMADEN conference*, 2000.
- [33] T. O'Neill. Normal discrimination with unclassified observations. *Journal of the American Statistical Association*, (73):821–826, 1978.
- [34] S. Roberts. Novelty detection using extreme value statistics. In *IEE Proc. on Vision, Image and Signal Processing*, volume 146, pages 124–129, 1999.
- [35] S. Roberts and L. Tarassenko. A probabilistic resource allocating network for novelty detection. *Neural Computation*, 6:270–284, 1994.
- [36] J. Ryan, M. Lin, and R. Mäkeläinen. Intrusion detection with neural networks. In *Advances in Neural Information Processing Systems*, 1998.
- [37] G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6:461–464, 1978.
- [38] M. Seeger. Learning with labeled and unlabeled data. Technical report, Institute for Adaptive and Neural Computation, University of Edinburgh, 2001.
- [39] B. Schölkopf, R. Williamson, A. Smola, J. Taylor, and J. Platt. Support vector method for novelty detection. In *Advances in Neural Information Processing Systems*, pages 582–588, 2000.
- [40] L. Tarassenko. Novelty detection for the identification of masses in mammograms. In *4th IEE International Conference on Artificial Neural Networks*, volume 4, pages 442–447, 1995.
- [41] D. Tax and R. Duin. Outlier detection using classifier instability. In *Advances in Pattern Recognition*, pages 251–256, 1999.
- [42] R. Udry. The national longitudinal of adolescent health: waves 1 and 2 (1994-1996), wave 3 (2001-2002). Technical report, Carolina Population Center, University of North Carolina, 2003.
- [43] D. Yeung and C. Chow. Parzen window network intrusion detectors. In *Proc. of International Conference on Pattern Recognition*, 2002.