



ENSTA PARISTECH - UNIVERSIT  PARIS SUD

FINAL YEAR INTERNSHIP

Report : Operational Feature Selection in Gaussian Mixture Models

Adrien Lagrange

supervised by
Dr. Mathieu FAUVEL

August 4, 2016

This document is non-confidential. Thus, it can be broadcast outside in paper or electronic format.

Lorem ipsum...

Keywords: enter, keyword, here.

Contents

1. Introduction	6
2. Organization	6
2.1. Objectives	6
2.2. Planning	6
3. Features selection with Gaussian Mixture Models	7
3.1. Gaussian Mixture Models	7
3.2. Criterion function	9
3.2.1. Measures of good classification	9
3.2.2. Measures of similarity between distributions	10
4. Ridge Regularization	11
4.1. Principe of Ridge Regularization	11
4.2. Implementation	13
5. Features selection	13
5.1. Sequential forward features selection	14
5.2. Sequential floating forward feature selection	16
5.3. Implementation	16
5.3.1. Update for cross validation	17
5.3.2. Criterion function computation	17
5.4. Computational issues	19
6. Experimental results about features selection with GMM	21
6.1. Aisa dataset	21
6.2. Experimental results	21
7. Implementation of the Orfeo Toolbox module	24
7.1. Integration in Orfeo Toolbox	24
7.2. Code structure	24
7.3. Applications	27
8. Experimental results comparison to classical classifiers	27
9. Conclusion	27
A. Update for cross validation (calculation)	28
A.1. Mean update	28
A.2. Covariance matrix update	28

List of Figures

1.	Gantt diagram of the internship.	7
2.	Confusion matrix with TP, TN, FP and FN relatively to c_1	9
3.	Consequence of regularization on eigenvalues.	12
4.	Example datasets: (a) 2 classes and 2 features but only one informative feature, (b) same dataset with a reduced gap between means of informative feature, (c) similar dataset with unbalanced classes (350/50 samples).	15
5.	Aisa dataset: (a) a single band of the image, (b) groundtruth.	21
6.	table	22
7.	Classification results averaged on 20 trials using forward selection and Jeffries-Matusita distance as criterion with 100 samples by class in training set.	22
8.	Classification results averaged on 20 trials using forward selection and Jeffries-Matusita distance as criterion: (a) with 50, 100, 200 samples by class in training set and (b) with 0.5%, 1%, 2.5% of dataset in training set keeping proportion.	23
9.	Classification results with 20 trials : (a) using Jeffries-Matusita distance as criterion and (b) using Cohen's kappa as criterion.	23
10.	Classification results with different criterion function with 20 trials : (a) with a training composed of 100 samples for each class and (b) with a training set keeping class proportions composed of 1% of the dataset (=xxxx samples).	24
11.	(Figures temporaires en attente de resultats) Processing time for selection with 20 trials : (a) using Cohen's kappa, (b) using overall accuracy, (c) using mean of F1-scores, (d) using KL divergence, (e) using JM distance.	25
12.	UML digram of the developped OTB remote module.	26
13.	table	28

1. Introduction

With the increasing number of remote sensing missions, the quantity of data available for a given landscape becomes larger and larger. In particular, the European satellite Sentinel-2 was launched successfully recently and the hyperspectral and hypertemporal data produced by this mission will be fully available at the end of 2016¹. Therefore, data is more and more difficult to process because computational and also statistical limits. Hence, in many remote sensing applications, the extraction of features from a large amount of available data is required.

For instance, in land-cover classification, given a set of spatial, temporal and spectral features, it is possible to extract those which are the most discriminant for the purpose of classification as explained in [1]. In hyperspectral data analysis from the hundreds of available spectral channels, it is possible to reduce the number of channels to make the processing more efficient.

There are two ways to reduce dimension: features extraction and feature selection. Feature extraction means creating new features in combining the existing ones, for example linear combination as in Principal Component Analysis [2]. To the contrary, features selection select a subset of the existing features which has the advantage to be much more interpretable for the end-user. The selected subset of features corresponds to the most important features for the given task.

There is a large diversity of method for feature selection. However, they usually do not scale well with the number of pixels to be processed as shown in [3]. Nevertheless, methods based on Gaussian Mixture Models (GMM) have several interesting properties described in [4] that make them suitable for feature selection in the context of large amount of data. By taking advantage of their intrinsic properties, it is possible to increase the computational efficiency with respect to standard implementation.

The objectives of the internship is to develop and implement a feature selection method based on GMM. The remaining of this mid-term report is organized as follows.

2. Organization

2.1. Objectives

Several outcomes are expected of this internship. First, based on the previous work [3] of M.Fauvel, we want to **reproduce the SFS selection method and develop an upgrade version** corresponding to the floating forward selection described in the previous section. This first part will results in a Python code freely available on Github².

Then, the aim is collaborate with CNES in order to **produce an external module which could be plugged in the open-source library Orfeo Toolbox (OTB) developed by CNES ([5])**. This module implements the same methods as the Python version but has to be written in C++ and be compatible with the OTB. The module obtained is a fork of a template furnished by OTB developers and is also available freely on Github³.

Finally, we want to **test the efficiency of this selection method and compare it to a kernel-based method ([6])**. The testing is conducted first on hyperspectral images where features are directly the spectral band and secondly on remote sensing images where numerous and various features are used (texture features, HoG, morphological profiles, radiometric indexes, ...).

2.2. Planning

As can be expected, the cycle is conducted from high-level to low-level which means that the Python code is first developed and so allow us to set definitively the structure of the algorithm. The low-

¹<https://sentinel2.cnes.fr/en/sentinel-2-0>

²<https://github.com/Laadr/FFFS>

³<https://github.com/Laadr/otbExternalFastFeaturesSelection>

level C++ code is then developed in order to assure good performances and to respect the interface with the OTB.

Each implementation is followed by a testing period during which the code is tested with artificial and easily manipulable data. Moreover, the development of the C++ code is done in interaction with developers of the OTB and, in particular, a first meeting is organized at the early stage of the development to assure that the best choice are made and a second meeting is set during the validation of the C++ code to get feedbacks.

The experimentation are conducted at the end to demonstrate the interest of the project. The Gantt diagram 1 summarized the expected organization of the project.

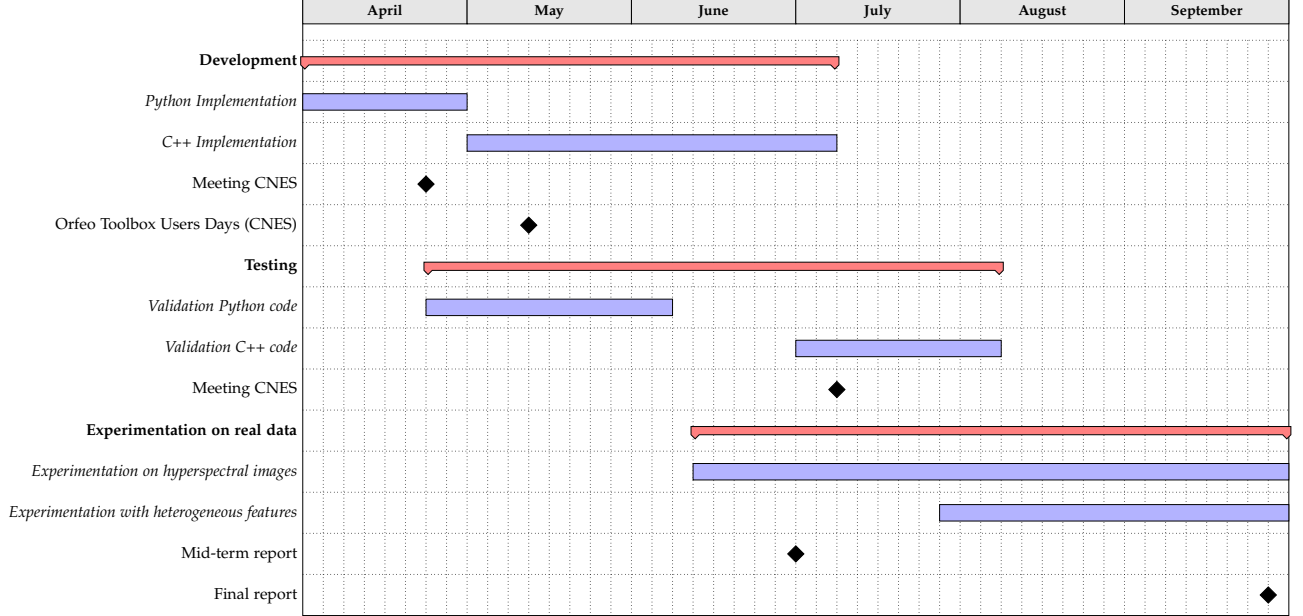


Figure 1: Gantt diagram of the internship.

3. Features selection with Gaussian Mixture Models

In the remaining, the following notations are used. $\mathcal{S} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ is the training set where $\mathbf{x}_i \in \mathbb{R}^d$ is the vector of features of the i^{th} sample, $y_i = 1, \dots, C$ the associated label, C the number of labels, n the number of samples and n_c the number of samples of class c .

3.1. Gaussian Mixture Models

The hypothesis of Gaussian mixture models (GMM) is that a given sample is the realization of a random vector which distribution is a mixture (convex combination) of several class conditioned distribution:

$$p(\mathbf{x}) = \sum_{c=1}^C \pi_c f_c(\mathbf{x}|\theta) \quad (1)$$

where π_c is the prior i.e. the proportion of class c and f_c a probability density function parametrized by θ .

The Gaussian model assumes that each f_c is , conditionally to c , the probability density of a

Gaussian distribution of parameters μ_c and Σ_c and so $f_c(\mathbf{x}|\theta)$ can be written

$$f_c(\mathbf{x}|\theta) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_c|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{x} - \mu_c)^t \Sigma_c^{-1} (\mathbf{x} - \mu_c) \right).$$

Such a model is used in the case of supervised learning and the class parameters μ_c and Σ_c can be estimated using the training samples. In our work, we choose to compute them with the conventional unbiased empirical estimator

$$\hat{\pi}_c = \frac{n_c}{n}, \quad (2)$$

$$\hat{\mu}_c = \frac{1}{n_c} \sum_{\{i|y_i=c\}} \mathbf{x}_i, \quad (3)$$

$$\hat{\Sigma}_c = \frac{1}{(n_c - 1)} \sum_{\{i|y_i=c\}} (\mathbf{x}_i - \mu_c)(\mathbf{x}_i - \mu_c)^t. \quad (4)$$

To determine the class of a samples, the maximum a posteriori (MAP) rule is used and thus, if we simplify with Bayes' law, the decision rule can be written as

$$\begin{aligned} \mathbf{x} \text{ belongs to } c &\Leftrightarrow c = \arg \max_{c \in C} p(c|\mathbf{x}), \\ &\Leftrightarrow c = \arg \max_{c \in C} \frac{p(c)p(\mathbf{x}|c)}{p(\mathbf{x})}, \\ &\Leftrightarrow c = \arg \max_{c \in C} p(c)p(\mathbf{x}|c). \end{aligned}$$

By taking the log, a simplified decision formula is obtained

$$\begin{aligned} Q_c(\mathbf{x}) &= 2 \log (p(c)p(\mathbf{x}|c)) \\ &= 2 \log \left(\pi_c \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_c|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{x} - \mu_c)^t \Sigma_c^{-1} (\mathbf{x} - \mu_c) \right) \right) \\ &= 2 \left(-\frac{1}{2} (\mathbf{x} - \mu_c)^t \Sigma_c^{-1} (\mathbf{x} - \mu_c) \right) + 2 \log \left(\frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_c|^{\frac{1}{2}}} \right) + 2 \log(\pi_c) \\ &= -(\mathbf{x} - \mu_c)^t \Sigma_c^{-1} (\mathbf{x} - \mu_c) - \log(|\Sigma_c|) + 2 \log(\pi_c) - d \log(2\pi). \end{aligned} \quad (5)$$

It is interesting to notice that the covariance, its inverse and its determinant are a key element of the decision function and that the estimation of these elements suffer from the curse of dimensionality. More precisely, the number of parameters to estimate (mean vectors, covariance matrices, proportions) increases quadratically relatively to the number of features and we need at least as many samples as parameters to make an estimation which can be an issue. For example, with hyperspectral data, we face high dimensional samples but very few samples are labeled available because of the difficulty and the cost to collect ground-truth.

A lack of samples induces ill-conditioned covariance matrices and so unstable inversion and computation of the determinant. There are two major solutions to this problem. First option is to use a regularization method to stabilize the inversion of the covariance matrices. Second option is to use a features extraction method in order to reduce the dimensionality of the samples.

In our study based on GMM, a Ridge regularization method described in Section 4 is implemented. Then we focus on a feature selection method named sequential forward features selection presented in Section 5.1 and an improvement of this algorithm, the sequential floating forward features selection introduced in Section 5.2.

In order to evaluate the benefits of such methods for classification, we introduce in the next section several functions called criterion functions.

3.2. Criterion function

Criterion functions aim to evaluate either a rate of good classification or the separability/similarity of class distributions. These functions are used to estimate which sets of variables are the best to represent data, to assure the separability of the classes and to perform classification. The choice of a criterion function is the choice of a way to compare sets of variables.

3.2.1. Measures of good classification

As described in [7] (chapter 4), all these measures of good classification are based on an error matrix M called confusion matrix which is defined so that M_{ij} is the number of samples of class i classified as class j . The confusion matrix allows the computation of several interesting values relatively to each class: the number of True Positive (TP) corresponding to good prediction, True Negative (TN) corresponding to good classification of the other class, False Negative (FN) corresponding to the samples of the class labeled as another class and the False Positive (FP) corresponding to the samples wrongly classified as part of this class. Figure 2 illustrates the definition of these values.

		Prediction outcome		
		c_1	c_2	c_3
Actual value	c'_1	True Positive	False Negative	False Negative
	c'_2	False Positive	True Negative	
	c'_3	False Positive		True Negative

Figure 2: Confusion matrix with TP, TN, FP and FN relatively to c_1

The overall accuracy is simply the rate of the number of samples with the correct predicted label over the number of samples. This metric is easy to interpret but is biased in the case of unbalanced classes.

The Cohen's kappa is a statistic which measures the probabilities of agreement between predictions and ground-truth. This metric tends to give an equal importance to each class.

The mean F1 score is the average of the F1 score for each class and the F1 score is the harmonic mean of the precision (number of True Positive over True Positive plus False Positive) and the recall (number of True Positive over True Positive plus False Negative).

$$\text{Overall Accuracy} = \frac{\text{TP}}{n} \quad (6)$$

$$\text{Kappa} = \frac{pa - pr}{1 - pr} \quad (7)$$

$$\text{F1 Mean} = \frac{2\text{TP}}{2\text{TP} + \text{FN} + \text{FP}} \quad (8)$$

where TP stands for True Positive, FN for False Negative, FP for False Positive, pa is the probability of agreement defined by $pa = \text{Overall Accuracy}$ and pr the probability of random agreement defined by $pr = \sum_{c=1}^C \frac{\text{TP}_c}{\text{FP}_c} \frac{\text{TP}_c}{\text{FN}_c}$.

In order to estimate classification rate, the GMM has to be trained and tested on separate datasets. In our case, we choose to use a cross-validation process over our training set. Thus, the training set is divided in k folds, then, the training is done with $(k - 1)$ folds and the performances are estimated with the remaining fold. For more details about cross-validation, see [8] (Chapter 7.10).

3.2.2. Measures of similarity between distributions

The second type of criterion functions is a measure of distance between two distributions. These measures are called divergence function and are defined so that, if we note S a space of all probability distribution with same support, it verifies

$$\begin{aligned} \forall (p, q) \in S, \text{Div}(p, q) &\geq 0, \\ \text{Div}(p, q) &= 0 \Leftrightarrow p = q. \end{aligned}$$

More specifically, we focus on two particular divergence: the Kullback–Leibler divergence and the Jeffries–Matusita distance. The advantage of these divergence is that they have a simplified expression in the case of Gaussian model. The simplification allows us to get rid of any integration calculus which could be a major problem when dealing with high-dimensional data.

The Kullback–Leibler divergence measures the amount of information lost when the first distribution is approximated by the second one. The formal definition is

$$\text{Div}_{KL}(c_i, c_j) = \int_{\mathbf{x}} p(\mathbf{x}|c_i) \ln\left(\frac{p(\mathbf{x}|c_i)}{p(\mathbf{x}|c_j)}\right) d\mathbf{x}. \quad (9)$$

And in the case of Gaussian model, it can be rewritten as follows

$$\text{Div}_{KL}(c_i, c_j) = \frac{1}{2} \left(\text{Tr}(\Sigma_{c_i}^{-1} \Sigma_{c_j}) + (\mu_{c_i} - \mu_{c_j})^t \Sigma_{c_i}^{-1} (\mu_{c_i} - \mu_{c_j}) - d + \log \left(\frac{|\Sigma_{c_i}|}{|\Sigma_{c_j}|} \right) \right) \quad (10)$$

where Tr is the trace operator and d the dimension of the distribution.

It can be noticed that the KL divergence is not symmetric, i.e. $\text{Div}_{KL}(c_i, c_j) \neq \text{Div}_{KL}(c_j, c_i)$, and so the symmetrized version is used to compute the criterion function. In the case of Gaussian model, the symmetrization induces the following simplification of the formula

$$\begin{aligned} \text{SKL}_{ij} &= \text{Div}_{KL}(c_i, c_j) + \text{Div}_{KL}(c_j, c_i) \\ &= \frac{1}{2} \left(\text{Tr}(\Sigma_{c_i}^{-1} \Sigma_{c_j} + \Sigma_{c_j}^{-1} \Sigma_{c_i}) + (\mu_{c_i} - \mu_{c_j})^t (\Sigma_{c_i}^{-1} + \Sigma_{c_j}^{-1}) (\mu_{c_i} - \mu_{c_j}) - 2d \right) \end{aligned} \quad (11)$$

Moreover, the divergence is computed between two classes and to obtain a unique value the weighted mean of divergence measures is taken

$$C_{SKL} = \sum_{i=1}^C \sum_{j=i+1}^C \pi_{c_i} \pi_{c_j} \text{SKL}_{ij} \quad (12)$$

The Bhattacharyya distance is defined as follows

$$B_{ij} = -\ln \left(\int_{\mathbf{x}} \sqrt{p(\mathbf{x}|c_i)p(\mathbf{x}|c_j)} d\mathbf{x} \right) \quad (13)$$

And in the case of Gaussian model:

$$B_{ij} = \frac{1}{8}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^t \left(\frac{\boldsymbol{\Sigma}_i + \boldsymbol{\Sigma}_j}{2} \right)^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) + \frac{1}{2} \ln \left(\frac{|\boldsymbol{\Sigma}_i + \boldsymbol{\Sigma}_j|}{\sqrt{|\boldsymbol{\Sigma}_i||\boldsymbol{\Sigma}_j|}} \right) \quad (14)$$

The Jeffries–Matusita distance is a measure based on the Bhattacharyya distance but transform in a way that the distance saturate if the separability between the two distribution increases. The JM distance is defined according to

$$JM_{ij} = \sqrt{\int_{\mathbf{x}} \left[\sqrt{p(\mathbf{x}|c_i)} - \sqrt{p(\mathbf{x}|c_j)} \right]^2 d\mathbf{x}} \quad (15)$$

And the Jeffries–Matusita distance can be rewritten according to the Bhattacharyya distance

$$JM_{ij} = \sqrt{2\{1 - \exp[-B_{ij}]\}} \quad (16)$$

As for the KL divergence, a weighted mean of the distance between two classes is computed to aggregate the measures in a single value.

According to [9], it is interesting to notice that the KL divergence increases quadratically with respect to the distance between the mean vectors of the class distributions whereas the measures of good classification we used asymptotically tends to one when distributions are perfectly separable. On the contrary, the JM distance tends to saturate as these measures of good classification.

4. Ridge Regularization

4.1. Principe of Ridge Regularization

As introduced in the previous section, the aim of a regularization method is to stabilize the inversion of the covariance matrix in the context of high dimensionality which often imply a badly-conditioned matrix.

The Ridge regularization is a particular case of the Tikhonov regularization introduced in [10] which aim to penalize the linear problem used to compute the inverse of the covariance matrix by the norm of the inverse. We want to obtain the regularized inverse by solving

$$\hat{\mathbf{A}} = \min_{\mathbf{A}} \|\boldsymbol{\Sigma}\mathbf{A} - \mathbf{I}\|^2 + \|\boldsymbol{\Gamma}\mathbf{A}\|^2. \quad (17)$$

In deriving the expression and looking at the minimum, we find that the explicit solution to this problem is

$$\hat{\mathbf{A}} = (\boldsymbol{\Sigma}^t\boldsymbol{\Sigma} + \boldsymbol{\Gamma}^t\boldsymbol{\Gamma})^{-1}\boldsymbol{\Sigma}. \quad (18)$$

There are then two cases, either we choose $\boldsymbol{\Gamma} = \sqrt{\tau}\mathbf{I}$ which corresponds to the Tikhonov regularization or $\boldsymbol{\Gamma} = \sqrt{\tau}\boldsymbol{\Sigma}^{\frac{1}{2}}$ which corresponds to the Ridge regularization with τ a scalar. Then, we get

$$\text{Tikhonov: } \hat{\mathbf{A}} = (\boldsymbol{\Sigma}^2 + \tau\mathbf{I})^{-1}\boldsymbol{\Sigma}, \quad (19)$$

$$\text{Ridge: } \hat{\mathbf{A}} = (\boldsymbol{\Sigma} + \tau\mathbf{I})^{-1}. \quad (20)$$

It is interesting to see the effect of the regularization method through the eigenvalues of the regularized inverse. Thus, we obtain the following expression

$$\text{Tikhonov: } \hat{\lambda}_i^{-1} = \frac{\lambda_i}{\lambda_i^2 + \tau}, \quad (21)$$

$$\text{Ridge: } \hat{\lambda}_i^{-1} = \frac{1}{\lambda_i + \tau}, \quad (22)$$

where λ_i is the i th eigenvalue of the covariance matrix and $\hat{\lambda}_i^{-1}$ the eigenvalue of the regularized inverse.

With the help of Figure 3, we can see that the difference between these two regularization is the way to treat small eigenvalues. The Ridge regression has two advantages. First, it is monotonous and secondly it does not converge to zero when eigenvalue tends to zero.

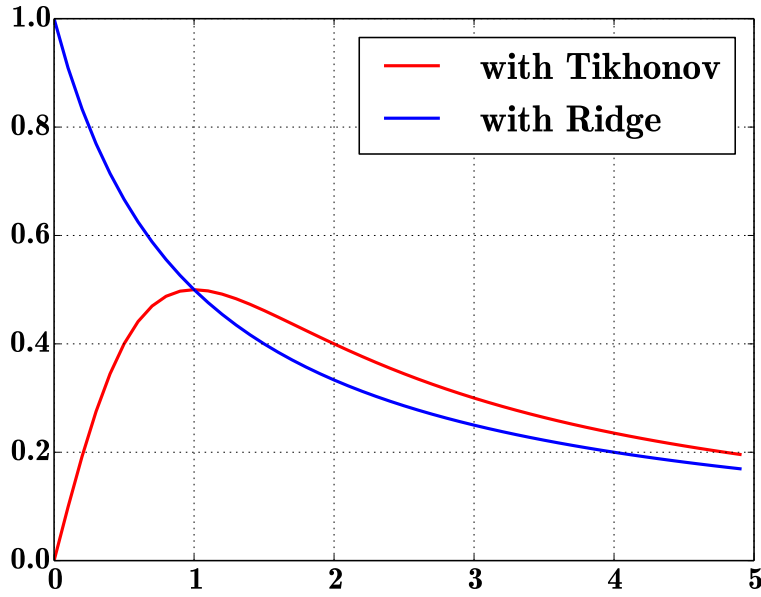


Figure 3: Consequence of regularization on eigenvalues.

In particular, the Ridge regularization proposes to solve the problem in adding a small value τ to each eigenvalues of the covariance matrix. In our case, we choose to add the same constant to each eigenvalues but it is not compulsory in a general case.

Thus, in practice, we choose to use Ridge regularization. To perform the regularization, the covariance matrix is decomposed with a diagonalization algorithm for symmetric matrices s.t. $\Sigma = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^t$ where \mathbf{Q} is the matrix of eigenvectors and $\mathbf{\Lambda}$ is the diagonal matrix of eigenvalues. When the decomposition is available, the regularization almost finish and the decision function is simply rewritten as follow

$$Q_c(\mathbf{x}) = -(\mathbf{x} - \boldsymbol{\mu}_c)^t \mathbf{Q}_c (\boldsymbol{\Lambda}_c + \tau \mathbf{I})^{-1} \mathbf{Q}_c^t (\mathbf{x} - \boldsymbol{\mu}_c) - \log(|\boldsymbol{\Sigma}_c + \tau \mathbf{I}|) + 2 \log(\pi_c) - d \log(2\pi), \quad (23)$$

where \mathbf{I} is the identity matrix. It is to be noticed that $(\boldsymbol{\Lambda} + \tau \mathbf{I})$ is also a diagonal matrix and so the inverse is easily obtained in inverting each diagonal element of the matrix.

4.2. Implementation

The main steps of the training algorithm when using Ridge regularization are summarized in Algorithm 1. Indeed, before training the GMM model, a preliminary step is needed in order to determine a good value for the parameter τ and to do so, a gridsearch is performed with a set of values \mathcal{P} given by the user. To estimate the quality of the classification with a given τ , one of the good classification criterion presented in Section 3.2.1 is used. As explained at the end of the same section, the classification rate is computed with a cross-validation process.

A particularity of this regularization method is used to accelerate the gridsearch. More precisely, the two computationally costly step are the computation of the inverse and the determinant of the regularized covariance matrices and it is important to see that these operations can be performed only once to test all τ . The covariance matrix is diagonalized before selecting τ and the eigenvalues and eigenvectors are stored and the regularized inverse and determinant can be update as follows

$$\begin{aligned}\hat{\Sigma}_c^{-1} &= \mathbf{Q}_c(\Lambda_c + \tau\mathbf{I})^{-1}\mathbf{Q}_c^t, \\ \log(|\hat{\Sigma}_c|) &= \log(|\Sigma_c + \tau\mathbf{I}|) = \sum_{i=1}^d \log(\lambda_i + \tau).\end{aligned}\tag{24}$$

Algorithm 1 Ridge Regularization training steps

Require: $\mathcal{S}, k, \mathcal{P}$

- 1: Learn GMM with \mathcal{S}
 - 2: Randomly cut \mathcal{S} in k subsets such as $\mathcal{S}_1 \cup \dots \cup \mathcal{S}_k = \mathcal{S}$ and $\mathcal{S}_i \cap \mathcal{S}_j = \emptyset$
 - 3: **for all** \mathcal{S}_i **do**
 - 4: Update GMM model for i^{th} fold with Equations 1 and 2
 - 5: Compute decomposition of covariance matrices of each class s.t. $\Sigma_c = \mathbf{Q}_c \Lambda_c \mathbf{Q}_c^t$
 - 6: **for all** $\tau \in \mathcal{P}$ **do**
 - 7: Update inverse and logdet with Equations 24
 - 8: Estimate classification rate on \mathcal{S}_i using Equation 23
 - 9: **end for**
 - 10: **end for**
 - 11: Average the classification rate over the k folds
 - 12: Compute $\arg \max_{\tau}$ of the classification rates to get the best parameter τ^*
 - 13: Parametrize the GMM with τ^*
-

Finally, another trick could be used to optimize the gridsearch. As explained in Section 5.3.1, it is possible to avoid the learning step on the $(k - 1)$ fold and instead learn the model with the whole dataset and deduce the submodel used for cross-validation from the complete model and the mean vector and covariance matrix of the k^{th} fold.

As stated before, we choose to implement a basic Ridge regression method and only the small improvement explained herebefore as been proposed. To the contrary, we choose in our study to focus more on a feature selection method named sequential forward features selection described in following section.

5. Features selection

The objective of feature selection method is to select a subset of variables to describe each sample and so to get rid of high dimensionality and its curse.

Features selection algorithm may be divided into three types. The first type which called filter method is based uniquely on data analysis. Features are ranked according to a statistical analysis of

the data. For example, the Principal Component Analysis (PCA) described in [2] is a typical filter method or else various metrics [11], [12], [13].

The second sort are known as wrapper methods which can be seen as search method to determine the best set of variables for a given learning model. As exhaustive searches are out of question in a practical amount of time, numerous search strategies have been designed some optimal under particular hypothesis [14] and other suboptimal but easier to set up [15], [16]. The advantage of such method compare to filter method is that they are dedicated to a particular model but on the other hand, as they require the training of multiple models to test various set of variables, they tend to be slower.

The third type corresponds to the embedded method which do not separate the features selection process from the learning algorithm and allow interaction between the two processes. The basic example of such method is the decision tree algorithm in which a feature is selected for creation of each node. Embedded methods also exist for other model, e.g. SVM ([17], [18]).

The selection method proposed in this work in a wrapper method associated to GMM models. It is important to underline that this method is a *selection* method and not an *extraction* method. It means the subset of variables obtained is composed of actual variables of the original set and not variables built as combination of several others as it is the case for example with PCA. This choice is made to assure an easier interpretation by the user of the obtained subset of variables.

Two selection algorithms are presented the sequential forward features selection method (Section 5.1) and the sequential floating forward feature selection method (5.2) the second being a more complex variation of the first.

5.1. Sequential forward features selection

The Sequential Forward Selection (SFS) starts with an empty set of selected features. Then it tests at each step for all the remaining features the value of a criterion function J chosen among the ones presented in Section 3.2 when the feature is added to the pool of selected features. The feature that maximizes the criterion function is definitively added to the pool of selected features and it moves to the next iteration. The algorithm stops when a given number of variables $maxVarNb$ has been selected.

To summarize, the features are selected one by one and the set of selected features are never questioned. It results in a reasonable computational time but a suboptimal solution to the selection problem meaning that a better subset of variables could exist.

Algorithm 2 Sequential forward features selection

Require: $S, J, maxVarNb$

```

1:  $S = \emptyset$ 
2:  $F = \{\text{all variables}\}$ 
3: while  $\text{card}(S) \leq maxVarNb$  do
4:   for all  $f_i \in F$  do
5:      $R_i = J(\{S + f_i\})$ 
6:   end for
7:    $j = \arg \max_i R_i$ 
8:    $S = \{S + f_j\}$ 
9:    $F = F \setminus f_j$ 
10: end while
11: return  $S$ 
```

Figure 4 shows a toy dataset with 2 classes and 2 features. The projection of the Gaussian distribution learned for each class can be visualized on the corresponding axis. In this examples, each class has been generated given a Gaussian law but it is not necessary in real case.

In the cases of the example presented in Figure 4, we compute the criterion functions described in 3.2 and summarize it in Table 1. We can see that all the criterion functions are maximum with the

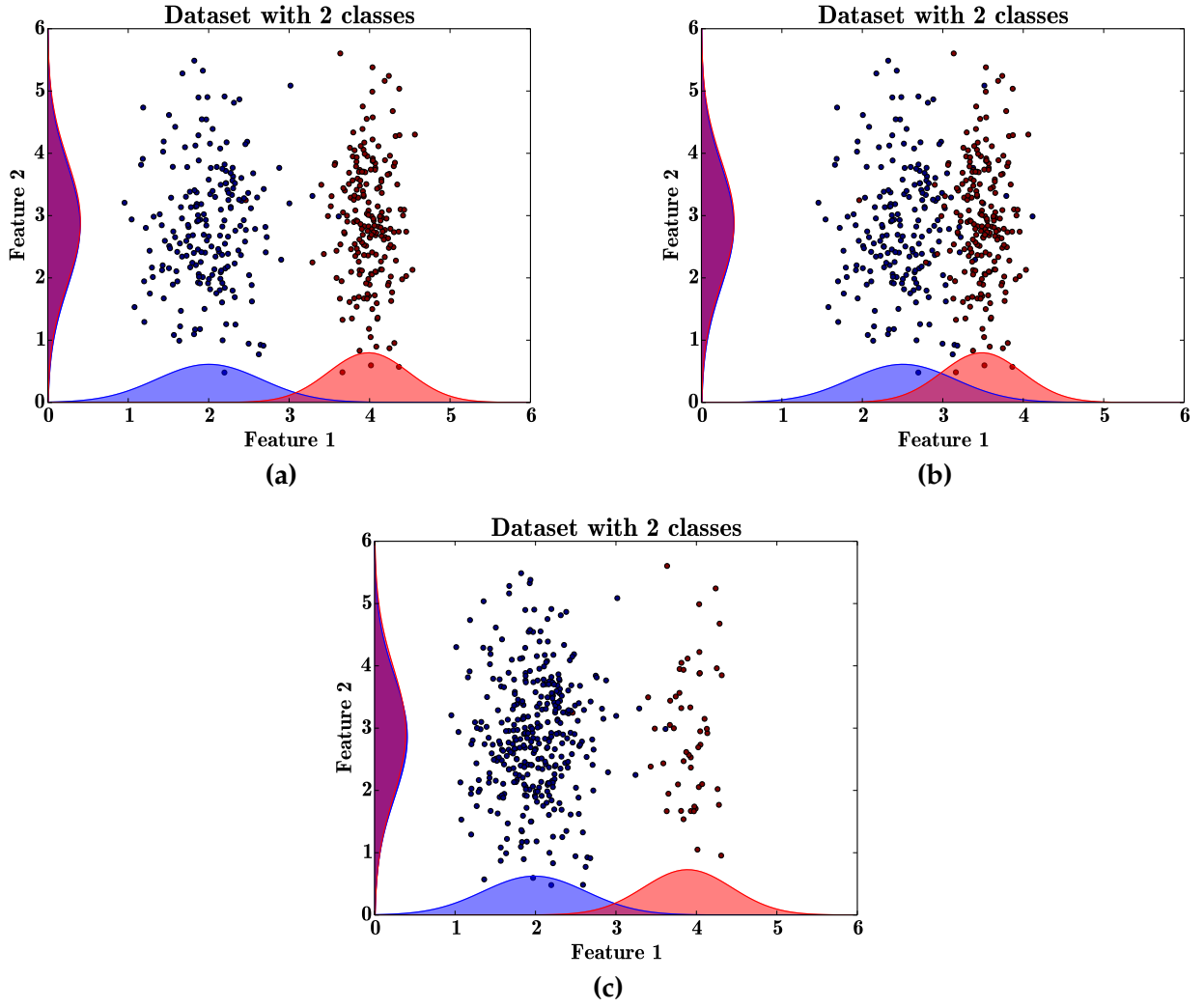


Figure 4: Example datasets: **(a)** 2 classes and 2 features but only one informative feature, **(b)** same dataset with a reduced gap between means of informative feature, **(c)** similar dataset with unbalanced classes (350/50 samples).

feature 1 which makes sense because the projection of the learned Gaussian distribution effectively seems separable only for feature 1. So the algorithm will have expected select the feature 1 as the best features to perform classification. The case (c) also shows that the measures are robust to unbalanced classes except the overall accuracy which still correctly ranks features in this case but makes less difference between features 1 and 2 because of the imbalance.

	Case a			Case b			Case c		
	Feat 1	Feat 2	Feat 1+2	Feat 1	Feat 2	Feat 1+2	Feat 1	Feat 2	Feat 1+2
Overall Accuracy	0.995	0.445	0.995	0.935	0.445	0.925	0.985	0.87	0.985
Cohen's kappa	0.99	-0.11	0.99	0.87	-0.11	0.85	0.95	0	0.95
F1-score mean	0.995	0.43	0.995	0.935	0.43	0.92	0.97	0.47	0.97
KL divergence	13.04	0.25	13.44	3.48	0.25	3.77	2.66	0.12	2.79
Jeffrey-Matusita distance	0.35	0.19	0.35	0.31	0.19	0.32	0.16	0.09	0.16

Table 1: Criterion functions values corresponding to Figure 4 dataset (50% for training/50% for testing).

5.2. Sequential floating forward feature selection

The Sequential Floating Forward Selection (SFFS) is actually based on two algorithms: the SFS described above and the Sequential Backward Selection (SBS). The SBS is the backward equivalent of SFS. The difference is that it starts with every features in the pool of selected features and tries at each step to remove the less significant one in term of the given criterion function.

The SFFS works as the SFS but between each step of the SFS algorithm a backward selection is operated. At the end of the SBS step, the value of the criterion function is compared to the best value ever obtained with a set of features of the same size and if the new value is better the feature puts into question is effectively taken away and the next step is again a SBS but if the new value is not better the SBS step is forgotten and it moves to the next SFS step. The algorithm stops when a given number of features has been selected.

This SFFS algorithm eventually tests more solutions than the SFS algorithm. The results are expected to be better but the trade-off is an increased computational time which is dependent on the complexity of the dataset.

Algorithm 3 Sequential floating forward features selection

Require: S, J, maxVar

```

1:  $S = \emptyset$ 
2:  $F = \{\text{all variables}\}$ 
3:  $k = 0$ 
4: while Repeat till  $S$  contained  $\text{maxVar}$  variables do
5:   for all  $f_i \in F$  do
6:      $R_i = J(\{S(k) + f_i\})$ 
7:   end for
8:    $j = \arg \max_i R_i$ 
9:   if  $R_j \geq J(S^{(k+1)})$  then
10:     $k = k + 1$ 
11:   else
12:     $S^{(k+1)} = \{S(k) + f_i\}$ 
13:     $k = k + 1$ 
14:     $\text{flag} = 1$ 
15:    while  $k > 2$  and  $\text{flag} = 1$  do
16:      for all  $f_i \in S^{(k)}$  do
17:         $R_i = J(\{S(k) \setminus f_i\})$ 
18:      end for
19:       $j = \arg \max_i R_i$ 
20:      if  $R_j < J(S^{(k-1)})$  then
21:         $S^{(k-1)} = \{S(k) \setminus f_i\}$ 
22:         $k = k - 1$ 
23:      else
24:         $\text{flag} = 0$ 
25:      end if
26:    end while
27:   end if
28: end while
29: return  $S$ 

```

5.3. Implementation

A major contribution of our work is the optimization in term of computational efficiency. More precisely, three steps of our method has been upgraded in order to reduce computational time. First of all, the GMM model is learned only once using samples. If a covariance matrix or a mean vector of a reduced set of variables is required, we get it from the global model learned at the beginning.

Secondly, when a cross-validation is performed, we do not learn a submodel, i.e. a GMM model trained with $(n - 1)$ folds, using the $(n - 1)$ folds but, instead, we derive it from the global model

and the covariance matrices and mean vectors of the fold used for validation for this given submodel. The process is described in details in Section 5.3.1.

Finally and most importantly, when variables are tested one by one during a selection step, costly operations are made when computing criterion functions for each variable especially the computation of the inverse of covariance matrices and its determinant. We manage to set up several update rules which allow us to compute this inverse and determinant only once to test all variables. These update rules are presented in Section 5.3.2.

5.3.1. Update for cross validation

Based on [3], a method to accelerate the k-fold cross-validation process in the case of criterion functions based on correct classification measures was implemented. The idea is to estimate the GMM model with the whole training set and then, instead of training a model on $(k - 1)$ folds, the complete model and the mean vector and the covariance matrix of the k^{th} fold is used to derive the corresponding submodel.

The following formulae can be obtained (details of calculation in Appendix A)

Proposition 1. (CV mean update)

$$\boldsymbol{\mu}_c^{n_c - v_c} = \frac{n_c \boldsymbol{\mu}_c^{n_c} - v_c \boldsymbol{\mu}_c^{v_c}}{n_c - v_c}$$

Proposition 2. (CV covariance matrix update)

$$\boldsymbol{\Sigma}_c^{n_c - v_c} = \frac{1}{n_c - v_c - 1} ((n_c - 1) \boldsymbol{\Sigma}_c^{n_c} - (v_c - 1) \boldsymbol{\Sigma}_c^{v_c} - \frac{n_c v_c}{(n_c - v_c)} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t)$$

where v_c is the number of samples of class c removed from the set.

Additionally, the update of the GMM model is the occasion to precompute and store the constant $2 \log(\pi_c)$ which is a term of the decision function constant for a given training set.

5.3.2. Criterion function computation

As explained earlier, at each iteration the SFS and SFFS algorithms compute the value of a criterion function for every possible set of features composed of the selected ones augmented by one of the remaining features. One of the main achievements of this work is to reduce the computational time needed to compute the criterion function of augmented sets.

We note $\boldsymbol{\Sigma}^{(k-1)}$ a covariance matrix of the $k - 1^{th}$ iteration, i.e., a covariance of the selected features and $\boldsymbol{\Sigma}^{(k)}$ the corresponding covariance matrix at the k^{th} iteration, i.e., the covariance of the augmented set. Then, the inverse of the covariance matrix $(\boldsymbol{\Sigma}^{(k)})^{-1}$, the quadratical term $(\mathbf{x}^{(k)})^t (\boldsymbol{\Sigma}^{(k)})^{-1} \mathbf{x}^{(k)}$ and the determinant $\log |\boldsymbol{\Sigma}^{(k)}|$ can be expressed in function of terms of the $(k - 1)^{th}$ iteration. These calculi use the fact that the covariance matrix is a positive definite symmetric matrix to simplify formulae using matrices defined by blocks. These update rules are summed up hereafter and are available in [4] (chapter 9.2).

In order to easily understand which are available and which ones needs to be computed, we choose to underline and set in red the terms already computed in the update rules.

As $\boldsymbol{\Sigma}^{(k)}$ is a positive definite symmetric matrix, we can use the following notation

$$\boldsymbol{\Sigma}^{(k)} = \begin{bmatrix} \boldsymbol{\Sigma}^{(k-1)} & \mathbf{u} \\ \mathbf{u}^t & \sigma_{kk} \end{bmatrix}$$

where \mathbf{u} is the k^{th} column of the matrix without the diagonal element i.e. $\mathbf{u}_i = \Sigma_{i+1,k}^{(k)}$ with $i \in [1, k - 1]$.

Using the formula of the inverse of a block matrix, the following formula expressing $(\boldsymbol{\Sigma}^{(k)})^{-1}$ in function of $(\boldsymbol{\Sigma}^{(k-1)})^{-1}$ is obtained

Proposition 3. (Forward update rule for inverse of covariance matrix)

$$(\Sigma^{(k)})^{-1} = \begin{bmatrix} \frac{(\Sigma^{(k-1)})^{-1}}{\alpha} + \frac{1}{\alpha}(\Sigma^{(k-1)})^{-1}\mathbf{u}\mathbf{u}^t(\Sigma^{(k-1)})^{-1} & -\frac{1}{\alpha}(\Sigma^{(k-1)})^{-1}\mathbf{u} \\ -\frac{1}{\alpha}\mathbf{u}^t(\Sigma^{(k-1)})^{-1} & \frac{1}{\alpha} \end{bmatrix}$$

where $\alpha = \sigma_{kk} - \mathbf{u}^t(\Sigma^{(k-1)})^{-1}\mathbf{u}$.

If we note

$$(\Sigma^{(k)})^{-1} = \begin{bmatrix} \mathbf{A} & \mathbf{v} \\ \mathbf{v}^t & \frac{1}{\alpha} \end{bmatrix}$$

with $\mathbf{A} = (\Sigma^{(k-1)})^{-1} + \frac{1}{\alpha}(\Sigma^{(k-1)})^{-1}\mathbf{u}\mathbf{u}^t(\Sigma^{(k-1)})^{-1}$ and $\mathbf{v} = -\frac{1}{\alpha}(\Sigma^{(k-1)})^{-1}\mathbf{u}$.

In the case of a backward step in SFFS algorithm, we need to invert the formula in order to compute $(\Sigma^{(k-1)})^{-1}$ knowing $(\Sigma^{(k)})^{-1}$. Thus, the update rule becomes

Proposition 4. (Backward update rule for inverse of covariance matrix)

$$(\Sigma^{(k-1)})^{-1} = \underline{\mathbf{A}} - \alpha\mathbf{v}\mathbf{v}^t$$

Proof.

$$\begin{aligned} \mathbf{A} - \alpha\mathbf{v}\mathbf{v}^t &= (\Sigma^{(k-1)})^{-1} + \frac{1}{\alpha}(\Sigma^{(k-1)})^{-1}\mathbf{u}\mathbf{u}^t(\Sigma^{(k-1)})^{-1} - \alpha(-\frac{1}{\alpha}(\Sigma^{(k-1)})^{-1}\mathbf{u})(-\frac{1}{\alpha}\mathbf{u}^t(\Sigma^{(k-1)})^{-1}) \\ &= (\Sigma^{(k-1)})^{-1} \end{aligned}$$

□

A formula can also be deduced for the quadratical term of the decision function. If we note $(\mathbf{x}^{(k)})^t = [(\mathbf{x}^{(k-1)})^t \ x_k]$, we get the following proposition

Proposition 5. (Update rule for quadratical term)

$$(\mathbf{x}^{(k)})^t(\Sigma^{(k)})^{-1}\mathbf{x}^{(k)} = \underline{(\mathbf{x}^{(k-1)})^t(\Sigma^{(k-1)})^{-1}\mathbf{x}^{(k-1)}} + \alpha([\mathbf{v}^t \ \frac{1}{\alpha}] \mathbf{x}^{(k)})^2$$

Proof.

$$\begin{aligned} (\mathbf{x}^{(k)})^t(\Sigma^{(k)})^{-1}\mathbf{x}^{(k)} &= [(\mathbf{x}^{(k-1)})^t \ x_k] \begin{bmatrix} \mathbf{A} & \mathbf{v} \\ \mathbf{v}^t & \frac{1}{\alpha} \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ x_k \end{bmatrix} \\ &= [(\mathbf{x}^{(k-1)})^t \mathbf{A} + x_k\mathbf{v}^t \quad (\mathbf{x}^{(k-1)})^t\mathbf{v} + \frac{x_k}{\alpha}] \begin{bmatrix} \mathbf{v} \\ x_k \end{bmatrix} \\ &= (\mathbf{x}^{(k-1)})^t\mathbf{A}\mathbf{x}^{(k-1)} + x_k\mathbf{v}^t\mathbf{x}^{(k-1)} + (\mathbf{x}^{(k-1)})^t\mathbf{v}x_k + \frac{(x_k)^2}{\alpha} \\ &= (\mathbf{x}^{(k-1)})^t((\Sigma^{(k-1)})^{-1} + \frac{1}{\alpha}(\Sigma^{(k-1)})^{-1}\mathbf{u}\mathbf{u}^t(\Sigma^{(k-1)})^{-1})\mathbf{x}^{(k-1)} + 2x_k\mathbf{v}^t\mathbf{x}^{(k-1)} + \frac{(x_k)^2}{\alpha} \\ &= (\mathbf{x}^{(k-1)})^t((\Sigma^{(k-1)})^{-1} + \alpha\mathbf{v}\mathbf{v}^t)\mathbf{x}^{(k-1)} + 2x_k\mathbf{v}^t\mathbf{x}^{(k-1)} + \frac{(x_k)^2}{\alpha} \\ &= (\mathbf{x}^{(k-1)})^t(\Sigma^{(k-1)})^{-1}\mathbf{x}^{(k-1)} + \alpha((\mathbf{x}^{(k-1)})^t\mathbf{v}\mathbf{v}^t\mathbf{x}^{(k-1)} + 2\frac{x_k}{\alpha}\mathbf{v}^t\mathbf{x}^{(k-1)} + \frac{(x_k)^2}{\alpha^2}) \\ &= (\mathbf{x}^{(k-1)})^t(\Sigma^{(k-1)})^{-1}\mathbf{x}^{(k-1)} + \alpha((\mathbf{x}^{(k-1)})^t\mathbf{v} + \frac{x_k}{\alpha})^2 \\ &= (\mathbf{x}^{(k-1)})^t(\Sigma^{(k-1)})^{-1}\mathbf{x}^{(k-1)} + \alpha([\mathbf{v}^t \ \frac{1}{\alpha}] \mathbf{x}^{(k)})^2 \end{aligned}$$

□

Finally, using the formula of the determinant of a block matrix and after taking the log, we obtain

Proposition 6. (Update rule for logdet)

$$\log(|\Sigma^{(k)}|) = \log(|\Sigma^{(k-1)}|) + \log \alpha$$

Now using all update rules, criterion functions can be rewritten as follows to compute as less times as possible the computationally heavy terms

$$\begin{aligned} Q(\mathbf{x}) = & - \frac{(\mathbf{x}^{(k-1)} - \boldsymbol{\mu}^{(k-1)})^t (\Sigma^{(k-1)})^{-1} (\mathbf{x}^{(k-1)} - \boldsymbol{\mu}^{(k-1)})}{2} \\ & - \alpha \left(\begin{bmatrix} \mathbf{v}^t & \frac{1}{\alpha} \end{bmatrix} (\mathbf{x}^{(k)} - \boldsymbol{\mu}^{(k)}) \right)^2 \\ & - \log(|\Sigma^{(k-1)}|) - \log \alpha + \underline{2 \log(\pi_c)} + \underline{k \log(2\pi)} \end{aligned}$$

$$\begin{aligned} \text{SKL}_{ij} = & \frac{1}{2} \left(\text{Tr}((\Sigma_{c_i}^{(k)})^{-1} \Sigma_{c_j}^{(k)} + (\Sigma_{c_j}^{(k)})^{-1} \Sigma_{c_i}^{(k)}) + \frac{(\boldsymbol{\mu}_{c_i}^{(k-1)} - \boldsymbol{\mu}_{c_j}^{(k-1)})^t ((\Sigma_{c_i}^{(k-1)})^{-1} + (\Sigma_{c_j}^{(k-1)})^{-1}) (\boldsymbol{\mu}_{c_i}^{(k-1)} - \boldsymbol{\mu}_{c_j}^{(k-1)})}{2} \right. \\ & \left. + \alpha \left(\begin{bmatrix} \mathbf{v}_i^t & \frac{1}{\alpha_i} \end{bmatrix} (\boldsymbol{\mu}_{c_i}^{(k)} - \boldsymbol{\mu}_{c_j}^{(k)}) \right)^2 + \alpha \left(\begin{bmatrix} \mathbf{v}_j^t & \frac{1}{\alpha_j} \end{bmatrix} (\boldsymbol{\mu}_{c_i}^{(k)} - \boldsymbol{\mu}_{c_j}^{(k)}) \right)^2 - \underline{2k} \right) \end{aligned}$$

with $(\Sigma_{c_i}^{(k)})^{-1}$ computed with Proposition 3.

$$\begin{aligned} B_{ij} = & \frac{1}{4} \frac{(\boldsymbol{\mu}_i^{(k-1)} - \boldsymbol{\mu}_j^{(k-1)})^t (\tilde{\Sigma}^{(k-1)})^{-1} (\boldsymbol{\mu}_i^{(k-1)} - \boldsymbol{\mu}_j^{(k-1)})}{2} + \frac{1}{4} \tilde{\alpha} \left(\begin{bmatrix} \tilde{\mathbf{v}}^t & \frac{1}{\tilde{\alpha}} \end{bmatrix} (\boldsymbol{\mu}_i^{(k)} - \boldsymbol{\mu}_j^{(k)}) \right)^2 \\ & + \frac{1}{2} \ln \left(\frac{|\tilde{\Sigma}^{(k-1)}|}{\sqrt{|\Sigma_i^{(k-1)}| |\Sigma_j^{(k-1)}|}} \right) + \frac{1}{2} \ln \left(\frac{\tilde{\alpha}}{\sqrt{\alpha_i \alpha_j}} \right) \end{aligned}$$

where $\tilde{\Sigma} = \Sigma_i + \Sigma_j$

It can be noticed that these new formulae allow us to precompute the computationally heavy terms because they are the same for all possible augmented sets at a given iteration of the selection. The Algorithm 4 illustrates the optimization of the Algorithm 3.

5.4. Computational issues

During an iteration of the selection algorithm, it still needs to invert the covariance matrix of the previous iteration. In order to so, we choose to perform a decomposition in eigenvalues and eigenvectors using algorithm specific for symmetric matrices. In other words, we find Λ the diagonal matrix of eigenvalues and \mathbf{Q} the matrix of eigenvectors s.t. $\Sigma = \mathbf{Q}\Lambda\mathbf{Q}^t$.

This decomposition has several advantages. First, it is very simple to compute the determinant once the eigenvalues are available because the determinant is equal to the product of these values. Secondly, it is possible to check the eigenvalues and so to assure a better computational stability. Due to the curse of dimensionality, the covariance matrix is sometimes badly conditioned and the results is that some eigenvalues are really small (below computational precision) and the decomposition helps us to check their actual values and if they are below EPS_FLT which the machine maximum precision for a float, we set these eigenvalues to EPS_FLT. For the same reason, the constant α used in update rules is also thresholded to EPS_FLT. Algorithm 4 details when computational stability is checked.

Algorithm 4 Sequential floating forward features selection with updates

Require: S, J, maxVar

```

1:  $S = \emptyset$ 
2:  $F = \{\text{all variables}\}$ 
3:  $k = 0$ 
4: while Repeat till  $S$  contained maxVar variables do
5:   Diagonalize  $\Sigma^{(k-1)} = \mathbf{Q}\Lambda\mathbf{Q}^t$ 
6:   for all  $\lambda_i$  do  $\lambda_i = \min(\text{EPS\_FLT}, \lambda_i)$ 
7:   Precompute  $(\Sigma^{(k-1)})^{-1}$ ,  $(\mathbf{x}^{(k-1)} - \boldsymbol{\mu}^{(k-1)})^t (\Sigma^{(k-1)})^{-1} (\mathbf{x}^{(k-1)} - \boldsymbol{\mu}^{(k-1)})$  and  $\log(|\Sigma^{(k-1)}|)$ 
8:   for all  $f_i \in F$  do
9:     Compute update constant  $\alpha$ 
10:     $\alpha = \min(\text{EPS\_FLT}, \alpha)$ 
11:     $R_i = J(\{S(k) + f_i\})$  using Properties 3, 5 and 6
12:   end for
13:    $j = \arg \max_i R_i$ 
14:   if  $R_j \geq J(S^{(k+1)})$  then
15:      $k = k + 1$ 
16:   else
17:      $S^{(k+1)} = \{S(k) + f_i\}$ 
18:      $k = k + 1$ 
19:     flag = 1
20:     while  $k > 2$  and flag = 1 do
21:       Diagonalize  $\Sigma^{(k-1)} = \mathbf{Q}\Lambda\mathbf{Q}^t$ 
22:       for all  $\lambda_i$  do  $\lambda_i = \min(\text{EPS\_FLT}, \lambda_i)$ 
23:       Precompute  $(\Sigma^{(k-1)})^{-1}$ ,  $(\mathbf{x}^{(k-1)} - \boldsymbol{\mu}^{(k-1)})^t (\Sigma^{(k-1)})^{-1} (\mathbf{x}^{(k-1)} - \boldsymbol{\mu}^{(k-1)})$  and  $\log(|\Sigma^{(k-1)}|)$ 
24:       for all  $f_i \in S^{(k)}$  do
25:         Compute update constant  $\alpha$ 
26:          $\alpha = \min(\text{EPS\_FLT}, \alpha)$ 
27:          $R_i = J(\{S(k) \setminus f_i\})$  using Properties 3, 5 and 6
28:       end for
29:        $j = \arg \max_i R_i$ 
30:       if  $R_j < J(S^{(k-1)})$  then
31:          $S^{(k-1)} = \{S(k) \setminus f_i\}$ 
32:          $k = k - 1$ 
33:       else
34:         flag = 0
35:       end if
36:     end while
37:   end if
38: end while
39: return  $S$ 

```

6. Experimental results about features selection with GMM

In this section, we discuss the results of the features selection method. The forward selection and the floating forward selection are compared and the results with the different criterion function are also discussed. All the experimentation has been conducted with a dataset presented in the next subsection.

6.1. Aisa dataset

The dataset has been acquired by the AISA Eagle sensor during a flight campaign over Heves, Hungary. It contains 252 bands ranging from 395 to 975 nm. 16 classes have been defined for a total of 361,971 referenced pixels. Figure 5 shows one channel of the multispectral image and the shapefile representing groundtruth and Table 6 presents the repartition of the various classes.

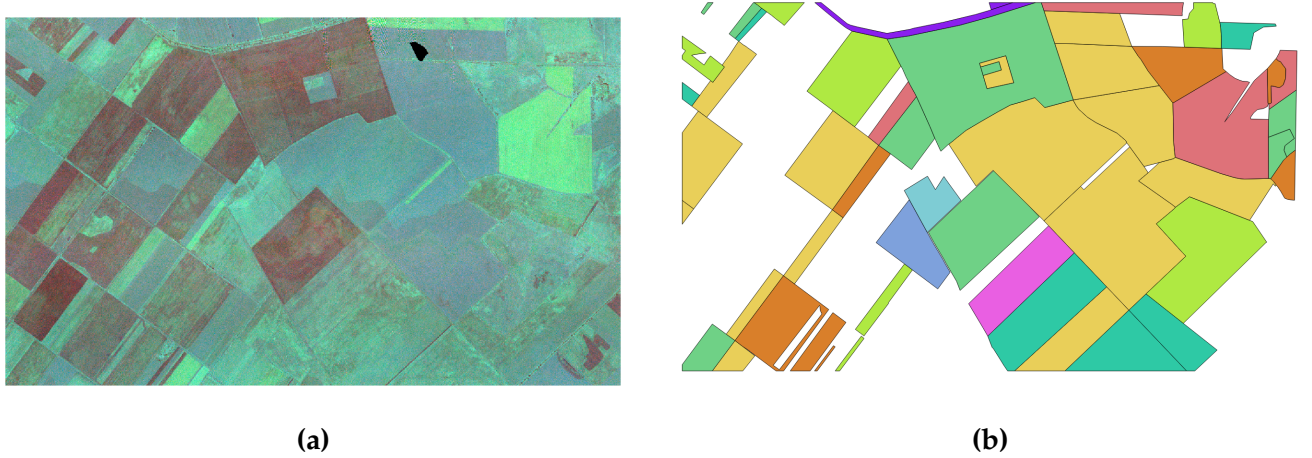


Figure 5: Aisa dataset: (a) a single band of the image, (b) groundtruth.

6.2. Experimental results

For the results discussed in this section, the selection process was repeated 20 times with each time a different training set and the classification rate and time processing presented hereafter are the means over these 20 trials.

Figure 7 corresponds to the evolution of the classification rate expressed with Cohen's kappa in function of the number of selected variables. We observe that at first the classification rate rapidly increase before stabilizing and then slowly decrease. This is the typical behavior expected when using features selection and a GMM classifier and it illustrated the Hughes phenomenon [19] which states that with a given number of samples, prediction accuracy first increases with dimensionality but then decays when the complexity is higher than some optimum value. We want to use the number of variables which maximizes the classification rate even if the maximum is not always easily found for example because of local maxima.

It is also interesting to confirm the influence of the size of the training set. From Figure 8, we see that, either with an equalized training set or with a proportional training set, the classification rate is obviously better with more samples and also that the maximum is reached with more selected variables which that more samples allows to evaluate more precisely the relevant variables.

Using Figure 9, we can compare the results when using forward selection or SFFS. And we see that with this particular dataset, the SFFS does not produce better results than the other method. One should be careful not to draw conclusion too fast. It does not mean that the SFFS method has

Class	Number of samples
Winter wheat	136,524
Sunflower	61,517
Green fallow last year treatment	30,197
Alfalfa	17,626
Maize	18,278
Millet	7,199
Broadleaved forest	10,746
Meadow	23,283
Winter barley	2,799
Reed	4,222
Water course	4,773
Rape	26,566
Green fallow with shrub	9,272
Green fallow last year treated	3,426
Pasture	2,107
Oat	3,436

Figure 6: table
Repartition of classes in Aisa dataset.

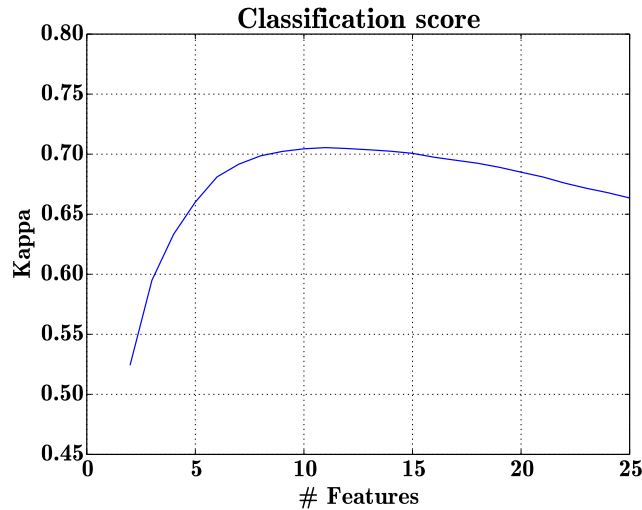


Figure 7: Classification results averaged on 20 trials using forward selection and Jeffries-Matusita distance as criterion with 100 samples by class in training set.

no advantages but rather that the SFFS advantages are highly related to the dataset structure and thus it is always good to check if the basic forward method is sufficient for a given dataset.

Now, if we take a closer look to the influence of criterion function, we can see from Figure 10 the results of classification applied to the validation set for each of criterion presented in Section 3.2. What we expected was that the criterion based on good classification measure would be better because it would optimize directly the classification rate and maybe more flexible regarding to the Gaussian assumption. But it appears that for this given dataset, as good results are obtained with Jeffries-Matusita distance.

Finally, if we compare the various method and criterion function in term of computational time, we first notice that the SFFS method takes more time than the forward method which is actually a fact supported by the theoretical ground. It can also be noticed that the processing time for selection

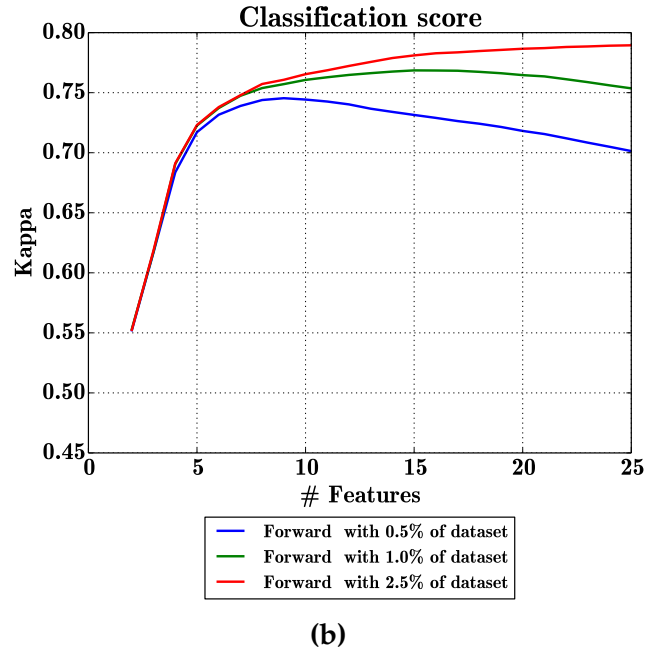
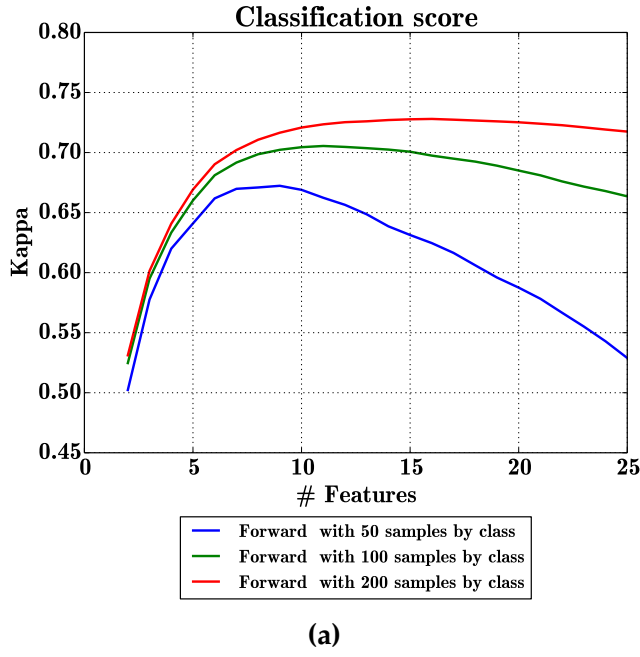


Figure 8: Classification results averaged on 20 trials using forward selection and Jeffries-Matusita distance as criterion: **(a)** with 50, 100, 200 samples by class in training set and **(b)** with 0.5%, 1%, 2.5% of dataset in training set keeping proportion.

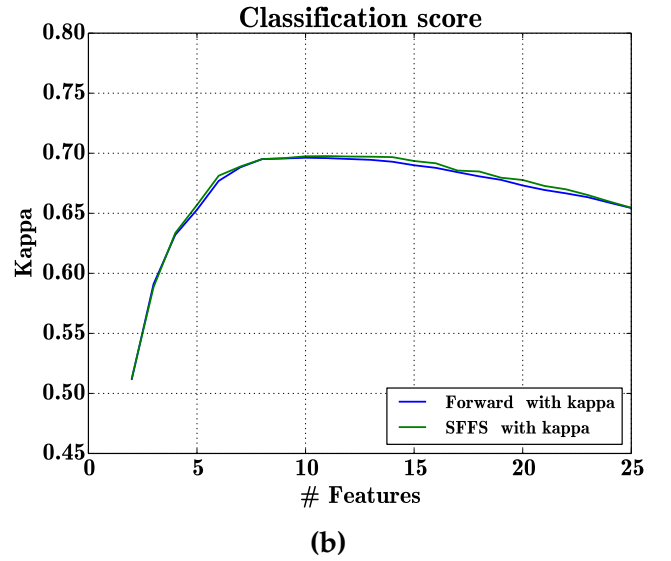
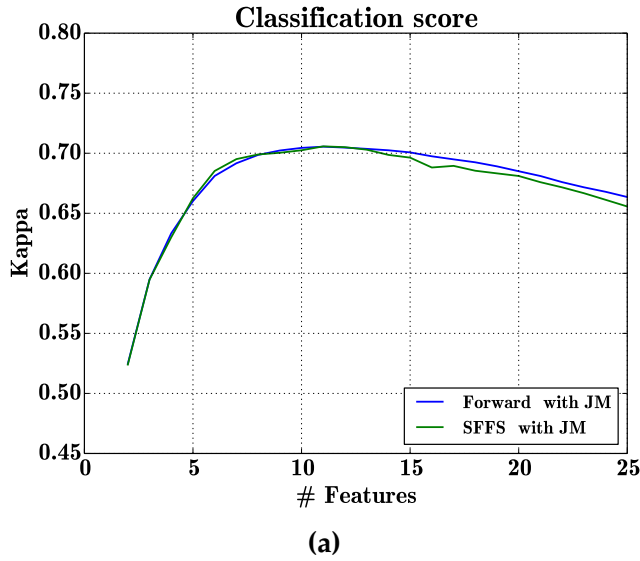


Figure 9: Classification results with 20 trials : **(a)** using Jeffries-Matusita distance as criterion and **(b)** using Cohen's kappa as criterion.

do not increase with the number of samples in training set for metrics based on distance between distribution (Kullback-Leibler divergence, Jeffries-Matusita distance) whereas, in the case of good classification criteria, the processing time rapidly increases. All processing time are presented in Figure 11.

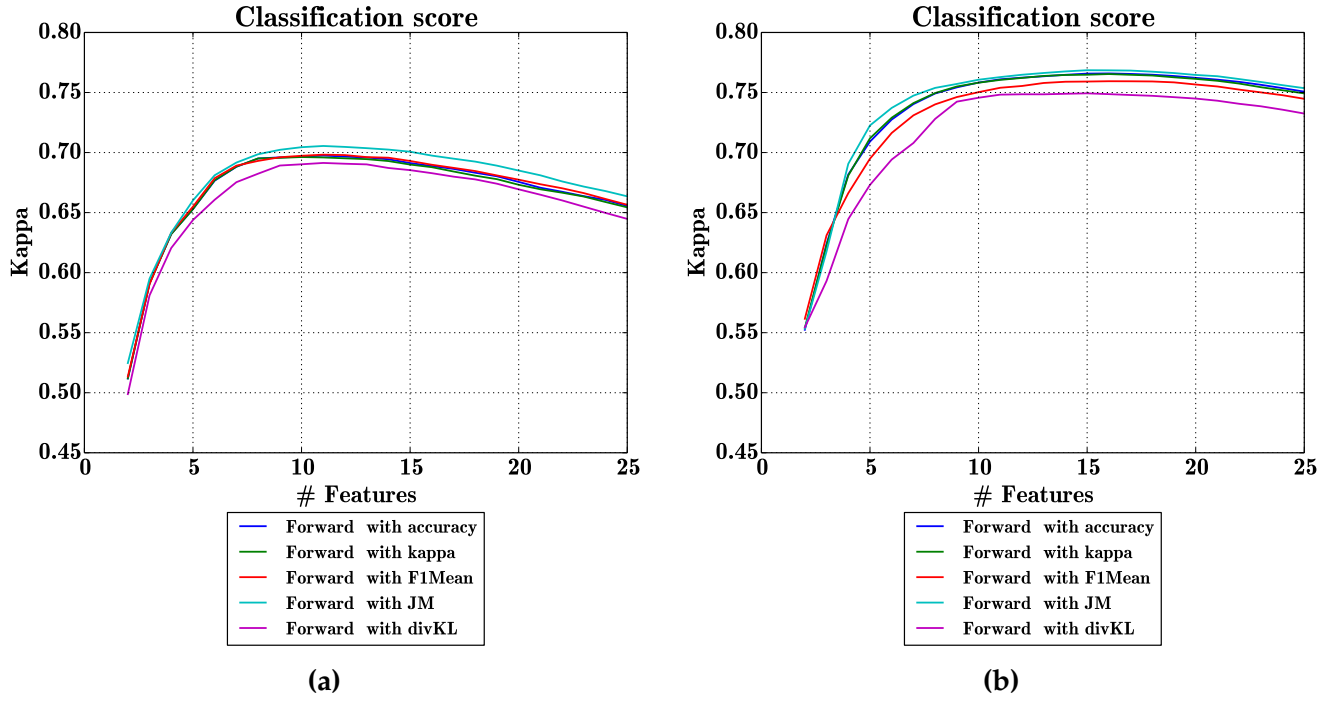


Figure 10: Classification results with different criterion function with 20 trials : **(a)** with a training composed of 100 samples for each class and **(b)** with a training set keeping class proportions composed of 1% of the dataset (=xxxx samples).

7. Implementation of the Orfeo Toolbox module

After the calibration of the algorithm with a Python implementation, we aim to realize a C++ implementation of this selection method. And, to make it easily available, we choose to develop a remote module for the Orfeo Toolbox (OTB) designed by the CNES. The current section describes this C++ implementation.

7.1. Integration in Orfeo Toolbox

The Orfeo Toolbox is an open-source library for remote sensing image processing. This library allows developers to code external modules and make them available to the community. It is to notice that the toolbox already implements several classifiers inherited mostly from OpenCV. Nevertheless, the GMM classifier currently in place is not robust to high dimensionality (see Section 8) and that's why we choose to implement our own model.

The module we developed is available on Github ⁴ and is actually a fork of the template for remote module furnished by OTB developers.

7.2. Code structure

Even if a GMM classifier inherited from Opencv library is already available in the OTB, we choose to implement our own version of a GMM classifier in a class named *GMMMachineLearningModel* in order to assure better performance in term of computational time. The second step is to implement a subclass of the GMM classifier including the features selection method.

The *GMMMachineLearningModel* class inherits from the OTB class *MachineLearningModel* which is a basic class used for all classifier in the OTB. The *MachineLearningModel* class enables the management

⁴<https://github.com/Laadr/otbExternalFastFeaturesSelection>

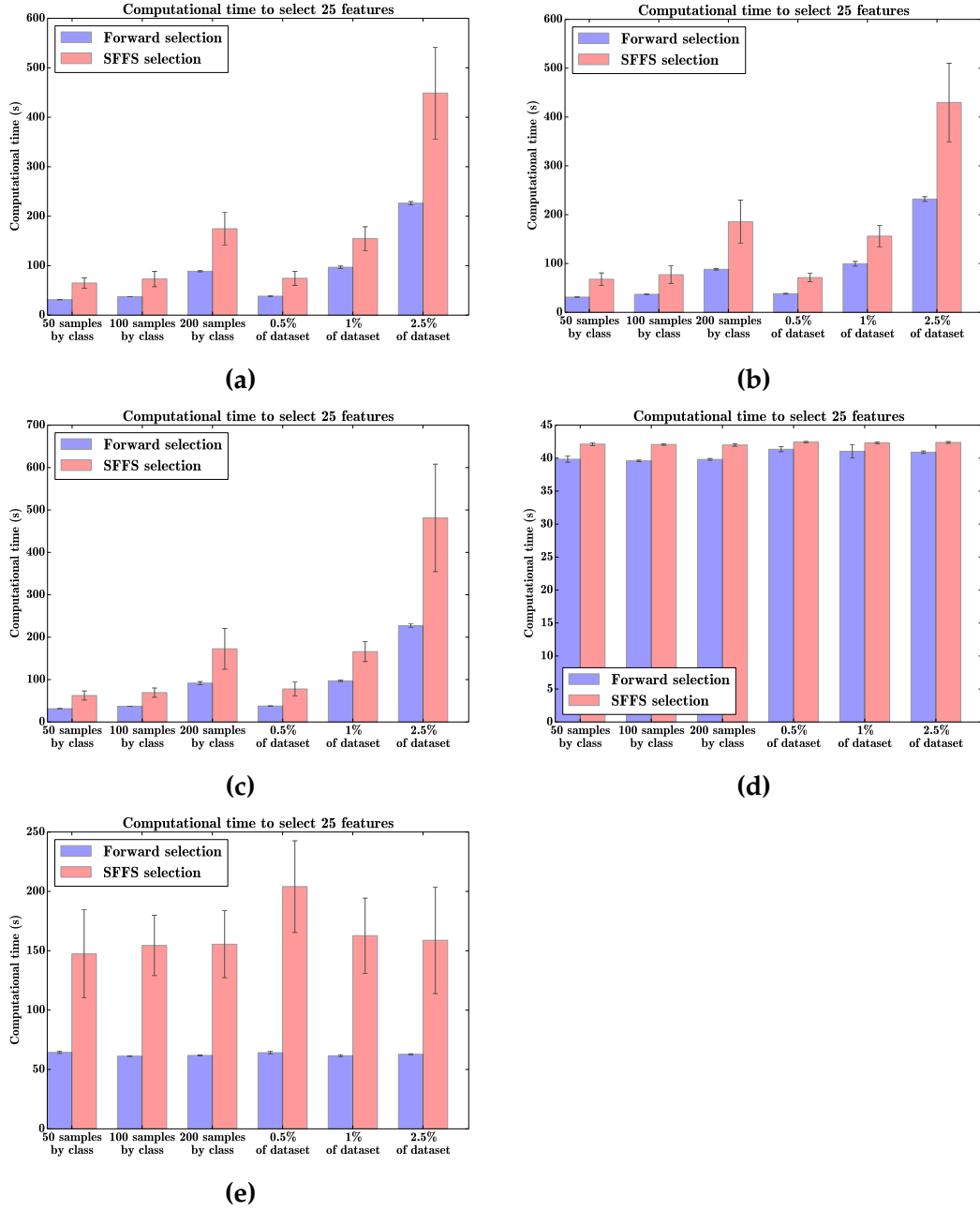


Figure 11: (Figures temporaires en attente de resultats) Processing time for selection with 20 trials : **(a)** using Cohen's kappa, **(b)** using overall accuracy, **(c)** using mean of F1-scores, **(d)** using KL divergence, **(e)** using JM distance.

of a list of samples used for training and also declares the virtual methods inherited: *Train()*, *Predict()*, *Save()*, *Load()*, *CanReadFile()*, *CanWriteFile()*.

The *GMMMachineLearningModel* class implements all these virtual methods. The *Train()* aims to train the classifier using a list of samples set as a member of the class. Then, it is possible to use the method *Predict()* to classify a sample and optional get the confidence in the prediction. The methods *Save()* and *Load()* obviously are used to save and load a trained model. And, finally, *CanReadFile()* and *CanWriteFile()* state the validity of a GMM model file.

In addition to these inherited methods, two noticeable methods are implemented. The *Decomposition()* method perform a decomposition of a symmetric matrix in eigenvalues and eigenvectors

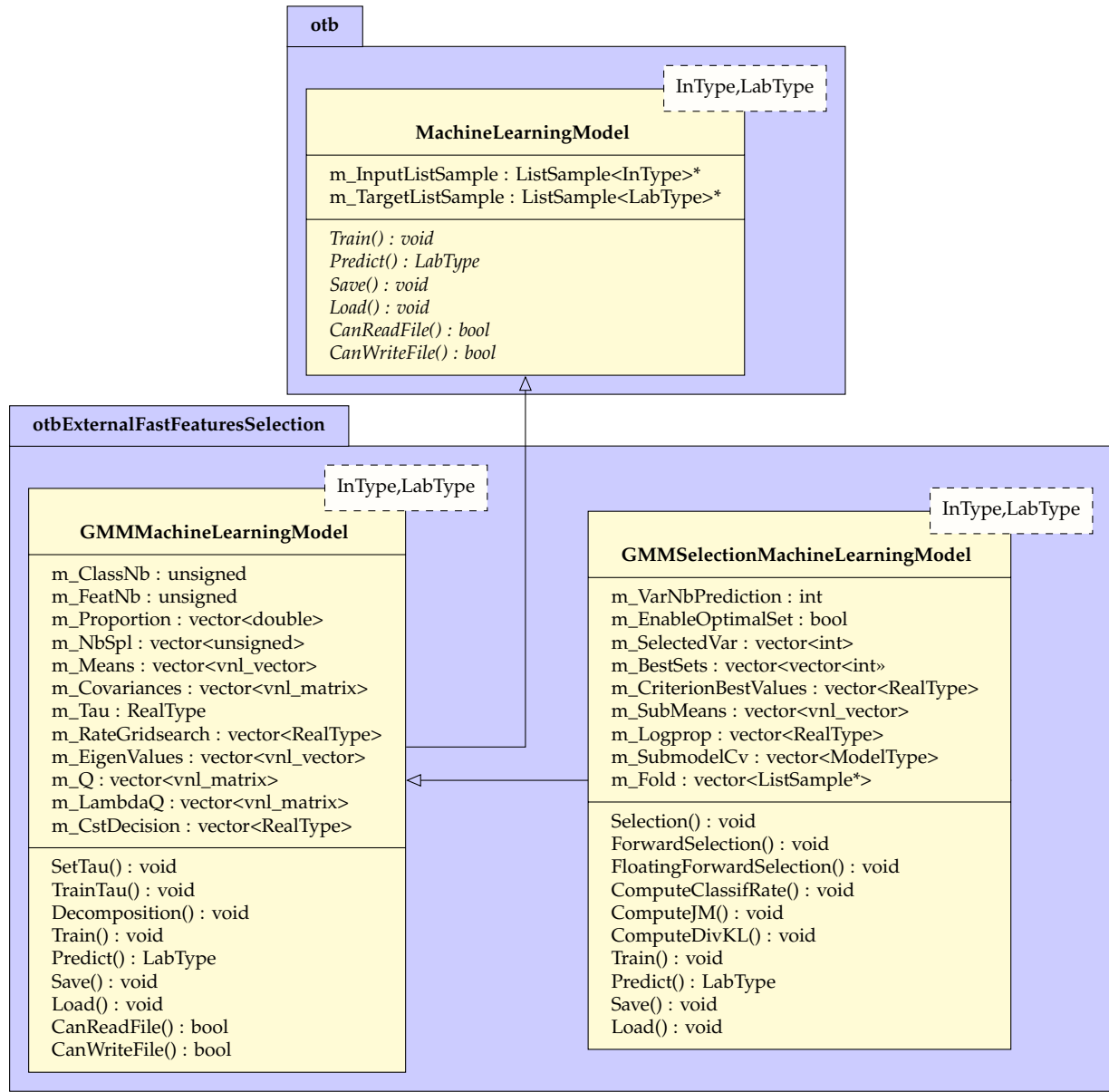


Figure 12: UML diagram of the developed OTB remote module.

using a function of the VNL library. It is to notice that this method check the value of the extracted eigenvalues and minor them to EPSILON_FLT (or EPSILON_DBL) which corresponds to computational precision. The idea is to limit computational errors when parameters are ill-estimated due to an insufficient amount of training samples. The second interesting method is *TrainTau()* which concerns the Ridge regularization. Given a set of possible regularization constant τ , this method selects the most efficient value in estimating classification rate with cross-validation.

In term of members, the *GMMMachineLearningModel* includes:

- 6 variables describing the model: the number of class, the number of features, the proportion of each class in the training set, the number of samples in each class, the mean vectors of each class and finally the covariance matrices of each class;
- 1 optional meta-parameter: the regularization value τ for Ridge regularization;
- 1 variable the result of tau selection;

- 4 precomputed terms used for prediction: eigenvalues of covariance matrices, eigenvectors of covariance matrices, 2 other terms.

To implement the selection algorithm, a new class *GMMSelectionMachineLearningModel* inheriting from *GMMMachineLearningModel* is defined. This class keeps the same six methods inherited from *MachineLearningModel* but reimplements *Predict()* to predict with a reduced set of features and *Save()/Load()* to use a second file to handle the results of the selection.

The other methods available in class *GMMSelectionMachineLearningModel* are used to perform the selection algorithms presented in Section 3. The method *Selection()* aims to set the various parameter of the selection and then calls either the *ForwardSelection()* method to use forward feature selection or the *FloatingForwardSelection()* method to use SFFS. Then, 3 different methods are used to evaluate criterion functions during selection. Given a pool of variables available for selection, these functions evaluate the criterion functions for all the possible augmented set. The *ComputeJM()* method is used for Jeffries-Matusita distance, *ComputeDivKL()* for Kullback–Leibler divergence and *ComputeClassifRate()* for the three good classification criteria. It is to notice that the good classification criteria are evaluated using cross-validation.

Several new members are defined in the *GMMSelectionMachineLearningModel* class:

- 2 meta-parameters: one to specify the number of variables to use for prediction and one to enable the process to automatically set the previous meta-parameter to the number maximizing the criterion function;
- 3 variables describing the results of selection: the set to use for prediction, the evolution of criterion function at each iteration and the best sets for each iteration;
- 2 variables used for cross-validation: the set of submodels and the folds of samples;
- 2 precomputed terms used for prediction.

All the code structure is summarized in the UML diagram displayed in Figure 12.

7.3. Applications

In order to provide an easy way to use the developed algorithm, we finally develop OTB application. These application allows the user to use and parametrize the algorithm from command line. Three applications are built.

The application *otbcli_TrainGMMApp* create and train a model from class *GMMMachineLearningModel* and the application *otbcli_TrainGMMSelectionApp* same for a model of class *GMMSelectionMachineLearningModel*. These two applications use a raster image and shapefile with groundtruth as input and create a model file.

The last application *otbcli_PredictGMMApp* takes a raster image as input and a model file and perform classification of the image. It generates a classification map in an image file and optionally a confidence map. An other application is available in the OTB to compute classification rate from the image and a groundtruth file.

8. Experimental results comparison to classical classifiers

9. Conclusion

	Training time (s)			Classification time (s)			Cohen's kappa		
# samples by class	250	500	1000	250	500	1000	250	500	1000
GMM SFS kappa	257	496	955	7.7	8.6	8.7	0.684 (0.025)	0.698 (0.027)	0.711 (0.028)
GMM SFS JM	8.6	8.9	9.1	9.7	9.8	9.6	0.680 (0.028)	0.700 (0.028)	0.710 (0.030)
GMM SFFS JM	8.8	9.0	9.3	9.7	9.8	9.8	0.680 (0.028)	0.700 (0.028)	0.710 (0.030)
GMM Ridge	71.7	105	167	530	530	530	0.611 (0.040)	0.620 (0.036)	0.642 (0.034)
KNN	8.9	19.6	59.7	387	639	887	0.551 (0.035)	0.563 (0.033)	0.574 (0.030)
Random Forest	24.5	49.3	105.1	33.0	41.7	45.9	0.645 (0.026)	0.673 (0.023)	0.693 (0.023)

Figure 13: table

A. Update for cross validation (calculation)

A.1. Mean update

$$\begin{aligned}
\mu_c^{n_c} &= \frac{1}{n_c} \sum_{j=1}^{n_c} \mathbf{x}_j \\
&= \frac{1}{n_c} \sum_{j=1}^{n_c - v_c} \mathbf{x}_j + \frac{1}{n_c} \sum_{j=n_c - v_c + 1}^{n_c} \mathbf{x}_j \\
&= \frac{n_c - v_c}{n_c} \mu_c^{n_c - v_c} + \frac{v_c}{n_c} \mu_c^{v_c} \\
&= \mu_c^{n_c - v_c} + \frac{v_c}{n_c} (\mu_c^{v_c} - \mu_c^{n_c - v_c})
\end{aligned}$$

$$\boxed{\mu_c^{n_c} = \mu_c^{n_c - v_c} + \frac{v_c}{n_c} (\mu_c^{v_c} - \mu_c^{n_c - v_c})} \quad (25)$$

$$\boxed{\mu_c^{n_c - v_c} = \frac{n_c \mu_c^{n_c} - v_c \mu_c^{v_c}}{n_c - v_c}} \quad (26)$$

A.2. Covariance matrix update

$$\begin{aligned}
\Sigma_c^{n_c} &= \frac{1}{n_c - 1} \sum_{j=1}^{n_c} (\mathbf{x}_j - \mu_c^{n_c})(\mathbf{x}_j - \mu_c^{n_c})^t \\
&= \frac{1}{n_c - 1} \sum_{j=1}^{n_c} (\mathbf{x}_j - \mu_c^{n_c - v_c} - \frac{v_c}{n_c} (\mu_c^{v_c} - \mu_c^{n_c - v_c})) (\mathbf{x}_j - \mu_c^{n_c - v_c} - \frac{v_c}{n_c} (\mu_c^{v_c} - \mu_c^{n_c - v_c}))^t \\
&= \frac{1}{n_c - 1} \sum_{j=1}^{n_c} (\mathbf{x}_j - \mu_c^{n_c - v_c})(\mathbf{x}_j - \mu_c^{n_c - v_c})^t \\
&\quad + \frac{v_c^2}{n_c^2} (\mu_c^{v_c} - \mu_c^{n_c - v_c})(\mu_c^{v_c} - \mu_c^{n_c - v_c})^t \\
&\quad - \frac{v_c}{n_c} (\mathbf{x}_j - \mu_c^{n_c - v_c})(\mu_c^{v_c} - \mu_c^{n_c - v_c})^t \\
&\quad - \frac{v_c}{n_c} (\mu_c^{v_c} - \mu_c^{n_c - v_c})(\mathbf{x}_j - \mu_c^{n_c - v_c})^t
\end{aligned}$$

$$\bullet (\mu_c^{v_c} - \mu_c^{n_c - v_c})(\mu_c^{v_c} - \mu_c^{n_c - v_c})^t = \frac{n_c^2}{(n_c - v_c)^2} (\mu_c^{v_c} - \mu_c^{n_c})(\mu_c^{v_c} - \mu_c^{n_c})^t$$

- $\frac{1}{n_c} \sum_{j=1}^{n_c} (\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c - v_c})^t = \frac{n_c v_c}{(n_c - v_c)^2} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t$
- $\frac{1}{n_c} \sum_{j=1}^{n_c} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c - v_c})(\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})^t = \frac{n_c v_c}{(n_c - v_c)^2} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t$

$$\begin{aligned}
\Sigma_c^{n_c} &= \frac{1}{n_c - 1} \sum_{j=1}^{n_c} (\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})(\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})^t - \frac{v_c^2}{(n_c - v_c)^2} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t \\
&= \frac{1}{n_c - 1} \sum_{j=1}^{n_c - v_c} (\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})(\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})^t + \frac{1}{n_c - 1} \sum_{j=n_c - v_c + 1}^{n_c} (\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})(\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})^t \\
&\quad - \frac{n_c v_c^2}{(n_c - 1)(n_c - v_c)^2} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t \\
&= \frac{n_c - v_c - 1}{n_c - 1} \Sigma_c^{n_c - v_c} + \frac{1}{n_c - 1} \sum_{j=n_c - v_c + 1}^{n_c} (\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})(\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})^t \\
&\quad - \frac{n_c v_c^2}{(n_c - 1)(n_c - v_c)^2} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t
\end{aligned}$$

$$\begin{aligned}
\sum_{j=n_c - v_c + 1}^{n_c} (\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})(\mathbf{x}_j - \boldsymbol{\mu}_c^{n_c - v_c})^t &= \sum_{j=n_c - v_c + 1}^{n_c} (\mathbf{x}_j - \boldsymbol{\mu}_c^{v_c} + \boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c - v_c})(\mathbf{x}_j - \boldsymbol{\mu}_c^{v_c} + \boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c - v_c})^t \\
&= (v_c - 1) \Sigma_c^{v_c} + v_c (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c - v_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c - v_c})^t \\
&\quad + \sum_{j=n_c - v_c + 1}^{n_c} (\mathbf{x}_j - \boldsymbol{\mu}_c^{v_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c - v_c})^t \\
&\quad + \sum_{j=n_c - v_c + 1}^{n_c} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c - v_c})(\mathbf{x}_j - \boldsymbol{\mu}_c^{v_c})^t \\
&= (v_c - 1) \Sigma_c^{v_c} + \frac{n_c^2 v_c}{(n_c - v_c)^2} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t
\end{aligned}$$

$$\boxed{\Sigma_c^{n_c} = \frac{n_c - v_c - 1}{n_c - 1} \Sigma_c^{n_c - v_c} + \frac{v_c - 1}{n_c - 1} \Sigma_c^{v_c} + \frac{n_c v_c}{(n_c - 1)(n_c - v_c)} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t} \quad (27)$$

$$\boxed{\Sigma_c^{n_c - v_c} = \frac{1}{n_c - v_c - 1} ((n_c - 1) \Sigma_c^{n_c} - (v_c - 1) \Sigma_c^{v_c} - \frac{n_c v_c}{(n_c - v_c)} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t)} \quad (28)$$

References

- [1] I. Guyon, S. Gunn, M. Nikravesh, and L. A. Zadeh, *Feature Extraction: Foundations and Applications (Studies in Fuzziness and Soft Computing)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [2] L. O. Jimenez and D. A. Landgrebe, "Supervised classification in high-dimensional space: geometrical, statistical, and asymptotical properties of multivariate data," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 28, no. 1, pp. 39–54, 1998.

- [3] M. Fauvel, C. Dechesne, A. Zullo, and F. Ferraty, "Fast forward feature selection of hyperspectral images for classification with gaussian mixture models," *Selected Topics in Applied Earth Observations and Remote Sensing, IEEE Journal of*, vol. 8, no. 6, pp. 2824–2831, 2015.
- [4] A. R. Webb, *Statistical pattern recognition*. John Wiley & Sons, 2003.
- [5] E. Christophe, J. Inglada, and A. Giros, "Orfeo toolbox: a complete solution for mapping from high resolution satellite images," *International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 37, pp. 1263–1268, 2008.
- [6] G. Camps-Valls, J. Mooij, and B. Scholkopf, "Remote sensing feature selection by kernel dependence measures," *IEEE Geoscience and Remote Sensing Letters*, vol. 7, no. 3, pp. 587–591, 2010.
- [7] R. G. Congalton and K. Green, *Assessing the accuracy of remotely sensed data: principles PCAd practices*. CRC press, 2008.
- [8] T. J. Hastie, R. J. Tibshirani, and J. H. Friedman, *The elements of statistical learning : data mining, inference, and prediction*. Springer series in statistics, Springer, 2009.
- [9] L. Bruzzone and C. Persello, "A novel approach to the selection of spatially invariant features for the classification of hyperspectral images with improved generalization capability," *Geoscience and Remote Sensing, IEEE Transactions on*, vol. 47, no. 9, pp. 3180–3191, 2009.
- [10] A. E. Hoerl and R. W. Kennard, "Ridge regression: Biased estimation for nonorthogonal problems," *Technometrics*, vol. 12, no. 1, pp. 55–67, 1970.
- [11] L. Bruzzone, F. Roli, and S. B. Serpico, "An extension of the jeffreys-matusita distance to multiclass cases for feature selection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 33, no. 6, pp. 1318–1321, 1995.
- [12] J. Biesiada and W. Duch, "Feature selection for high-dimensional data—a pearson redundancy based filter," in *Computer Recognition Systems 2*, pp. 242–249, Springer, 2007.
- [13] B. Demir and S. Ertürk, "Phase correlation based redundancy removal in feature weighting band selection for hyperspectral images," *International journal of Remote sensing*, vol. 29, no. 6, pp. 1801–1807, 2008.
- [14] P. M. Narendra and K. Fukunaga, "A branch and bound algorithm for feature subset selection," *IEEE Transactions on Computers*, vol. 100, no. 9, pp. 917–922, 1977.
- [15] A. W. Whitney, "A direct method of nonparametric measurement selection," *IEEE Transactions on Computers*, vol. 100, no. 9, pp. 1100–1103, 1971.
- [16] P. Somol, P. Pudil, J. Novovičová, and P. Pačlık, "Adaptive floating search methods in feature selection," *Pattern recognition letters*, vol. 20, no. 11, pp. 1157–1163, 1999.
- [17] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene selection for cancer classification using support vector machines," *Machine learning*, vol. 46, no. 1-3, pp. 389–422, 2002.
- [18] J. Weston, A. Elisseeff, B. Schölkopf, and M. Tipping, "Use of the zero-norm with linear models and kernel methods," *Journal of machine learning research*, vol. 3, no. Mar, pp. 1439–1461, 2003.
- [19] G. Hughes, "On the mean accuracy of statistical pattern recognizers," *IEEE transactions on information theory*, vol. 14, no. 1, pp. 55–63, 1968.