

Report : Operational Feature Selection in Gaussian Mixture Models

Adrien Lagrange

Dynafor

July 17, 2016

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 [?]. 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 [?]. 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 [1]. Nevertheless, methods based on Gaussian Mixture Models (GMM) have several interesting properties described in [2] 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. Non linear parsimonious features selection

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 .

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

2.1. Gaussian Mixture Models

The hypothesis of Gaussian mixture models (GMM) is that the distribution of a given sample 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 . The Gaussian model assumes that each f_c is, conditionally to c , a Gaussian distribution of parameters μ_c and Σ_c ($f_c(\mathbf{x}|\theta) = \mathcal{N}_c(\mu_c, \Sigma_c)$) 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 following 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 using the Bayes' law, the maximum a posteriori (MAP) estimate is used and thus 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

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

It is interesting to notice the covariance and its inverse 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 (means, 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 available.

A lack of samples induces ill-conditioned covariance matrices and so unstable inversion. There are two major solutions to this problem. The first option is to use a regularization method to stabilize the inversion of the covariance matrices. The 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 3 is implemented but we mainly focus on a feature selection method named sequential forward features selection presented in Section 4.1 and an amelioration of this algorithm the sequential floating forward features selection introduced in Section 4.2.

In order to evaluate the benefits of such methods for classification, we introduce in the next section several functions called criterion functions. These functions aim either to evaluate a criterion of good classification or either to evaluate the separability/similarity of class distributions. The choice of one of these criterion functions is up to the user and can depend on the specificity of the task. We purposely choose to present only some of possible functions among the most used. It is also to be noticed that the criterion function is not only use for evaluation in the case of the features selection algorithm presented hereafter but is indeed a parameter of the algorithm.

2.2. Criterion function

2.2.1. Measures of good classification

All this 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 an other class and the False Positive (FP) corresponding to the samples wrongly classified as part of of this class. Figure 1 illustrates the definition of this values.

		Prediction outcome	
		p	n
Actual value	p'	True Positive	False Negative
	n'	False Positive	True Negative

Figure 1: Confusion matrix

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 measure the probabilities of agreement between predictions and ground-truth. This metric tends to give an equal importance to each class but is hard to interpret.

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).

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

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

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

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 $pr = \sum_{c=1}^C \frac{TP_c}{FP_c} \frac{TP_c}{FN_c}$.

In order to compute these metrics, the predicted labels are needed and so we need to compute the decision function (Equation 5). Moreover to estimate this measures of good classification, a cross-validation method is used which reduces both bias and variance of the estimation.

2.2.2. Measures of similarity between distributions

The Kullback–Leibler divergence is a measure of distance between two distributions. It actually 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 k the dimension of the distribution.

It can be noticed that the KL divergence is not symmetric 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 an other measure of similarity between two distributions. In fact, we do not discuss the efficiency of this measure in our work and simply used it to compute the next measure. The Bhattacharyya distance is computed 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 [3], 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.

3. Ridge Regularization

3.1. 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 ill-conditioned matrix. 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, to perform the regularization, the covariance matrix is decomposed with a diagonalization algorithm for symmetric matrices s.t. $\boldsymbol{\Sigma} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^t$ where \mathbf{Q} is the matrix of eigenvectors and $\boldsymbol{\Lambda}$ is the diagonal matrix of eigenvalues. When the decomposition is available, the regularization almost finish and the decision function is simply rewritten as follow

$$Q_c(\mathbf{x}) = -(\mathbf{x} - \boldsymbol{\mu}_c)^t \mathbf{Q}_c (\boldsymbol{\Lambda}_c + \tau \mathbf{I})^{-1} \mathbf{Q}_c^t (\mathbf{x} - \boldsymbol{\mu}_c) - \log((2\pi)^d |\boldsymbol{\Sigma}_c + \tau \mathbf{I}|) + 2 \log(\pi_c) \quad (17)$$

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 element of the matrix.

3.2. 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. Moreover, to have a good estimation of the quality of the classification with a

given τ a cross-validation method is used. 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.

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 used to test all τ using Equation 17.

Algorithm 1 Ridge Regularization training steps

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

- 1: Randomly cut \mathcal{S} in k subsets such as $\mathcal{S}_1 \cup \dots \cup \mathcal{S}_k = \mathcal{S}$ and $\mathcal{S}_i \cap \mathcal{S}_j = \emptyset$
 - 2: **for all** \mathcal{S}_i **do**
 - 3: Learn GMM with $\mathcal{S} \setminus \mathcal{S}_i$
 - 4: Compute decomposition of covariance matrices of each classes s.t. $\Sigma_c = \mathbf{Q}_c \Lambda_c \mathbf{Q}_c^t$
 - 5: **for all** $\tau \in \mathcal{P}$ **do**
 - 6: Estimate classification rate on \mathcal{S}_i using Equation 17
 - 7: **end for**
 - 8: **end for**
 - 9: Average the classification rate over the k folds
 - 10: Compute $\arg \max_{\tau}$ of the classification rates to get the best parameter τ^*
 - 11: Learn GMM with the whole training set \mathcal{S} and parametrize the GMM with τ^*
-

Finally, another trick could be used to optimize the gridsearch. As explained in the next section, 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.

4. Features selection

We recall that the aim of feature selection method is to select a subset of variables to described each samples and so to get rid of high dimensionality and its curse. It is also important to underline that the proposed method is a *selection* method and not an *extraction* method. It means the subset of variables obtained at the end 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 4.1) and the sequential floating forward feature selection method (4.2) the second being a more complex variation of the first.

4.1. 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 by the user among the ones presented in Section 2.2 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 $maxVar$ has been selected.

To summarize, the features are selected one by one and the set of selected features are never put into questions. It results in a reasonable computational time but a suboptimal solution to the selection problem.

Algorithm 2 Sequential forward features selection

Require: $S, J, maxVar$

```

1:  $S = \emptyset$ 
2:  $F = \{\text{all variables}\}$ 
3: while Repeat till  $S$  contained  $maxVar$  variables 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 2 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.

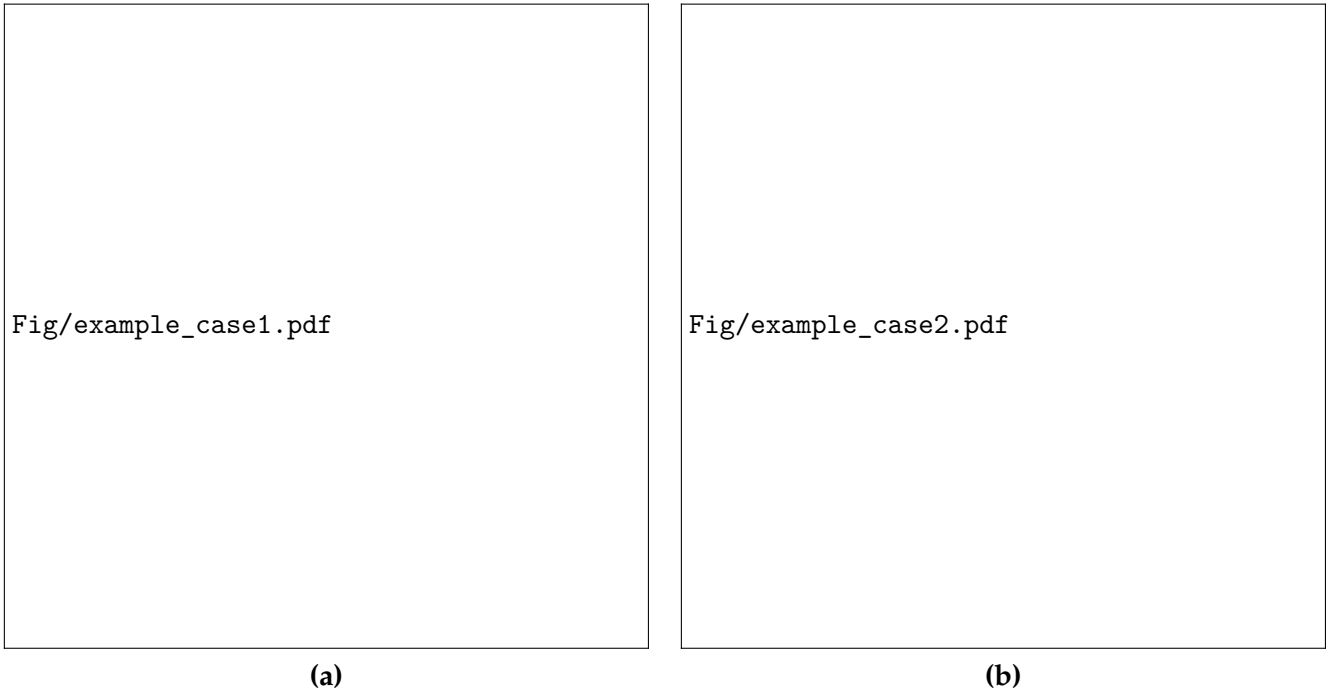


Figure 2: Example datasets: **(a)** 2 classes and 2 features but only one informative feature, **(b)** same dataset with a reduced gap between means of informative feature.

In the case of the example presented in Figure 2, we compute the various criterion functions using a single feature and summarize it in Table 1. It corresponds to first step of the SFS. And we can see that all the criterion functions are maximum with the feature 1 which makes sense because the projection of the learned Gaussian distribution effectively seems separable only for feature 1. So the

algorithm will has expected select the feature 1 as the best features to perform classification.

	Case a		Case b	
	Feat 1	Feat 2	Feat 1	Feat 2
Overall Accuracy	0.995	0.445	0.935	0.445
Cohen's kappa	0.99	-0.11	0.87	-0.11
F1-score mean	0.995	0.43	0.935	0.43
KL divergence	13.04	0.25	3.48	0.25
Jeffrey-Matusita distance	0.35	0.19	0.31	0.19

Table 1: Criterion functions values corresponding to Figure 2 dataset (50% for training/50% for testing).

4.2. 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 of the complexity of the dataset.

4.3. Implementation

In term of computational efficiency, the main issue is to find an efficient way to compute the mean vector and the covariance matrix and its inverse and to do it only when it is absolutely necessary. More specifically, we develop an optimized way to derive the submodel for cross-validation from the full model and we minimize the number of costly operations during a step of the selection algorithm.

4.3.1. Update for cross validation

Based on [1], 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)

$$\mu_c^{n_c - v_c} = \frac{n_c \mu_c^{n_c} - v_c \mu_c^{v_c}}{n_c - v_c} \quad (18)$$

$$\Sigma_c^{n_c - v_c} = \frac{1}{n_c - v_c - 1} ((n_c - 1) \Sigma_c^{n_c} - (v_c - 1) \Sigma_c^{v_c} - \frac{n_c v_c}{(n_c - v_c)} (\mu_c^{v_c} - \mu_c^{n_c})(\mu_c^{v_c} - \mu_c^{n_c})^t) \quad (19)$$

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  $S$ 
```

where ν_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.

4.3.2. 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 $\Sigma^{(k-1)}$ the covariance matrix of the $k-1^{th}$ iteration, i.e., the covariance of the selected features and $\Sigma^{(k)}$ the covariance matrix of the k^{th} iteration, i.e., the covariance of the augmented set. Then, the inverse of the covariance matrix $(\Sigma^{(k)})^{-1}$, the quadratical term $(\mathbf{x}^{(k)})^t (\Sigma^{(k)})^{-1} \mathbf{x}^{(k)}$ and the determinant $\log |\Sigma^{(k)}|$ 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 [2] (chapter 9.2).

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

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

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 $(\Sigma^{(k)})^{-1}$ in function of $(\Sigma^{(k-1)})^{-1}$ is obtained

$$(\Sigma^{(k)})^{-1} = \begin{bmatrix} (\Sigma^{(k-1)})^{-1} + \frac{1}{\alpha} (\Sigma^{(k-1)})^{-1} \mathbf{u} \mathbf{u}^t (\Sigma^{(k-1)})^{-1} & -\frac{1}{\alpha} (\Sigma^{(k-1)})^{-1} \mathbf{u} \\ -\frac{1}{\alpha} \mathbf{u}^t (\Sigma^{(k-1)})^{-1} & \frac{1}{\alpha} \end{bmatrix} \quad (20)$$

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

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 $(\Sigma^{(k)})^{-1}$. Then, if we use the following notation for the previous matrix

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

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

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

$$(\Sigma^{(k-1)})^{-1} = \mathbf{A} - \alpha \mathbf{v} \mathbf{v}^t \quad (21)$$

A formula can also be deduced for the quadratical term of the decision function. If we note

$$(\mathbf{x}^{(k)})^t = \begin{bmatrix} (\mathbf{x}^{(k-1)})^t & x_k \end{bmatrix},$$

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

$(\mathbf{x}^{(k)})^t (\boldsymbol{\Sigma}^{(k)})^{-1} \mathbf{x}^{(k)} = (\mathbf{x}^{(k-1)})^t (\boldsymbol{\Sigma}^{(k-1)})^{-1} \mathbf{x}^{(k-1)} + \alpha \left(\begin{bmatrix} \mathbf{v}^t & \frac{1}{\alpha} \end{bmatrix} \mathbf{x}^{(k)} \right)^2 \quad (22)$

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

$$\log \left(|\boldsymbol{\Sigma}^{(k)}| \right) = \log \left(|\boldsymbol{\Sigma}^{(k-1)}| \right) + \log \alpha \quad (23)$$

To give an example, if a measure of good classification is used, we need to compute the decision rule and with the update rules, it can be rewritten as follow

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

It can be noticed that this new formula 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 Figure 3 illustrates the optimization of the selection algorithm.

Nevertheless, 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 $\mathbf{\Lambda}$ the diagonal matrix of eigenvalues and \mathbf{Q} the matrix of eigenvectors s.t. $\boldsymbol{\Sigma} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^t$.

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 ill-estimated 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.

5. 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

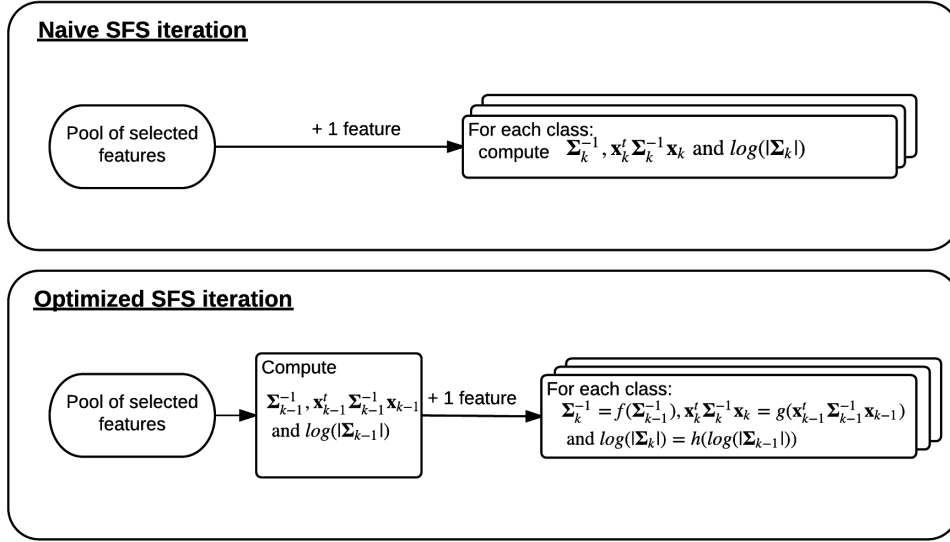


Figure 3: Iteration of sequential forward features selection.

also discussed. All the experimentation has been conducted with a single dataset presented in the next subsection.

5.1. Aisa dataset

The Aisa dataset is composed of an hyperspectral image with 256 spectral band between 42 nm and 42 nm and a shapefile containing groundtruth. The groundtruth is made with 16 classes representing landcover of the area. A representation of the dataset is displayed in Figure 4.

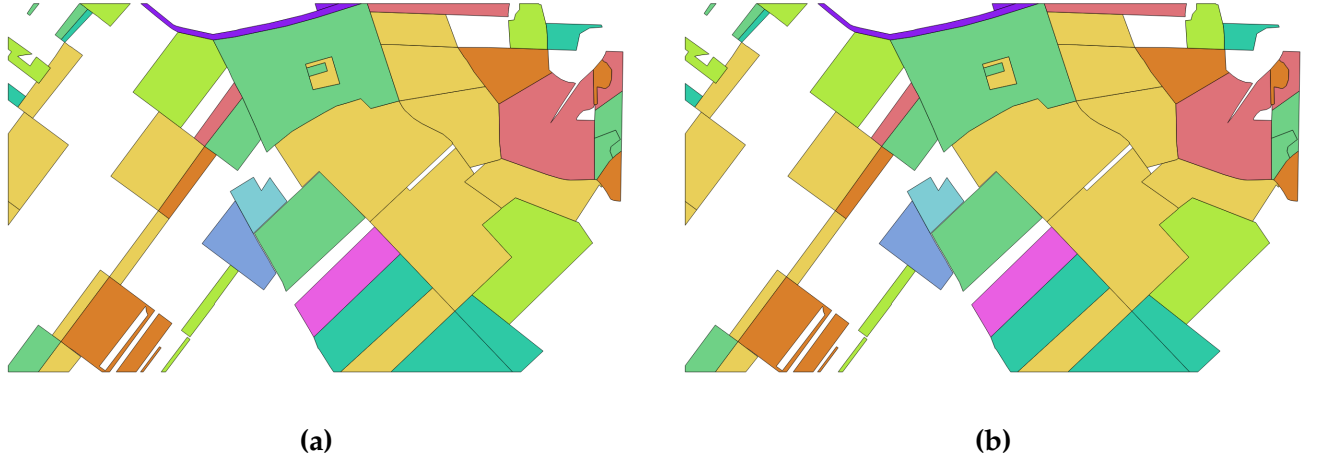


Figure 4: Aisa dataset : (a) color composition and (b) groundtruth.

5.2. 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 6 corresponds to the evolution of the classification rate in function of the number of selected variables with the three metrics introduced before (overall accuracy, Cohen's kappa and mean of F1-score). 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 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.

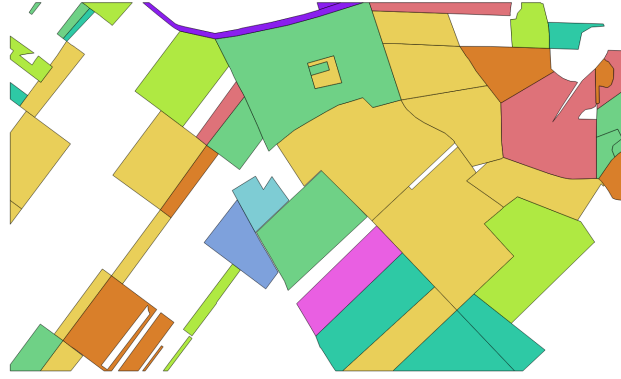


Figure 5: Classification results (overall accuracy, Cohen's kappa and mean of F1-score) averaged on 20 trials using forward selection and Jeffries-Matusita distance as criterion.

It is also interesting to see the influence of the size of the training set. From Figure ??, we see that, either with an equalized training set or with a proportional training set, the classification rate is 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.

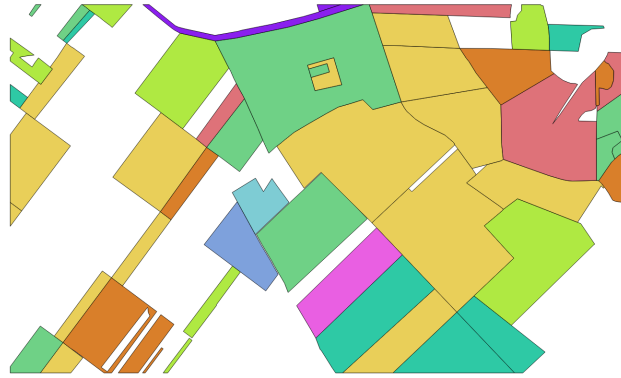


Figure 6: Classification results averaged on 20 trials using forward selection and Jeffries-Matusita distance as criterion: **(a)** with 50, 100, 200 samples by class in training set and **(b)** with 0.5%, 1%, 2.5% of dataset in training set keeping proportion.

Using Figure 7, 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 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.

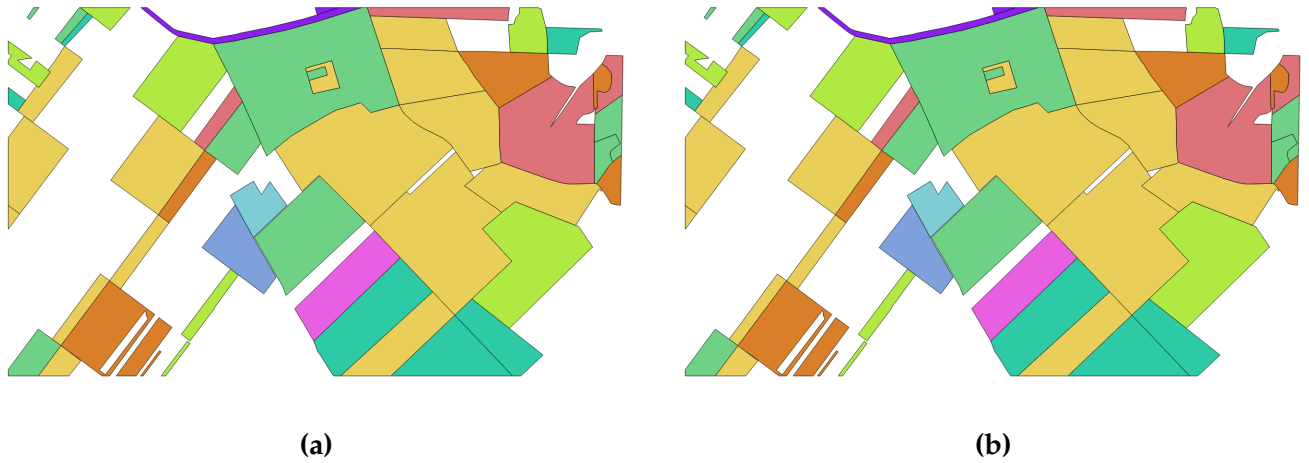


Figure 7: Classification results with 20 trials : **(a)** using Jeffries-Matusita distance as criterion and **(b)** using Cohen's kappa as criterion.

A. Update for cross validation (calculation)

A.1. Mean update

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

$$\boxed{\mu_c^{n_c} = \mu_c^{n_c - v_c} + \frac{v_c}{n_c} (\mu_c^{v_c} - \mu_c^{n_c - v_c})} \quad (24)$$

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

A.2. Covariance matrix update

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

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

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

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

$$\boxed{\Sigma_c^{n_c} = \frac{n_c - v_c - 1}{n_c - 1} \Sigma_c^{n_c - v_c} + \frac{v_c - 1}{n_c - 1} \Sigma_c^{v_c} + \frac{n_c v_c}{(n_c - 1)(n_c - v_c)} (\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})(\boldsymbol{\mu}_c^{v_c} - \boldsymbol{\mu}_c^{n_c})^t} \quad (26)$$

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

References

- [1] 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.
- [2] A. R. Webb, *Statistical pattern recognition*. John Wiley & Sons, 2003.
- [3] 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.