

# Wine Project Report

Ruggero Nocera (SXXXXX1)

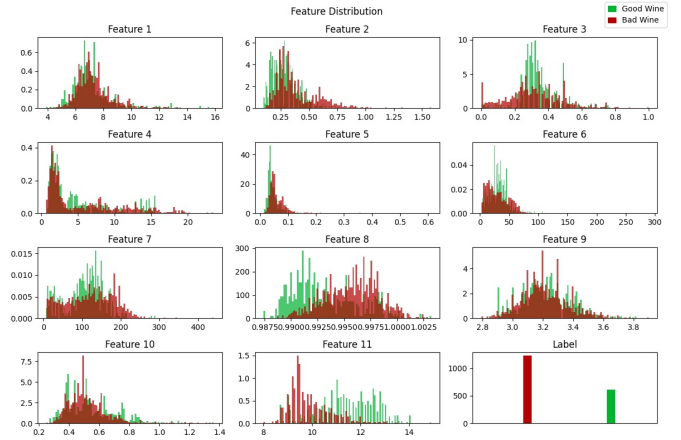
Quarta Matteo (SXXXXXX)

**Abstract** This paper is an analysis of the effectiveness of the various models studied during the course applied to a binary classification model. The dataset is a transformed version of the one provided in *Modeling wine preferences by data mining from physicochemical properties* from *Decision Support Systems, Elsevier* (P. Cortez, A. Cerdeira, F. Almeida, T. Matos and J. Reis.), where grades span from 0 to 10. Grade 6 has been removed, grades above 6 have been mapped to class 1 and grades below 6 to class 0. Due to limited time and computational power, results may be just close-to or far-from optimal, depending on the time required to train a model.

## Contents

<b>1 Preliminary Data Analysis</b>	<b>1</b>
1.1 Feature Distribution . . . . .	1
<b>2 Pre-Processing Analysis</b>	<b>3</b>
2.1 Pre-Processing MVG Classifiers . .	3
2.2 Pre-Processing for GMMs . . . . .	3
2.3 Pre-Processing for Polynomial Kernel SVMs . . . . .	3
2.4 Pre-Processing for RBF Kernel SVMs . . . . .	4
<b>3 Optimizing Models</b>	<b>4</b>
3.1 Optimizing MVG Classifiers . . . .	4

Figure 1: Histogram of Class' Features



## 1 Preliminary Data Analysis

### 1.1 Feature Distribution

Before discussing the model, their implementation and their effectiveness we briefly take a look at how the features are distributed. For convenience we shall now report the legend just one, but keep in mind that in all pictures red color is associated to class 0 (which we will be referring to as class *Bad*) and green color is associated to class 1 (which will be class *Good*).

First of all, our training dataset is unbalanced. In the next pages we will be classifying samples obtained from a K-Fold<sup>1</sup> Validation approach, using a theoretical threshold given by:

$$t = -\log \frac{\pi}{1 - \pi}$$

For the threshold to be optimal, we should use the empirical prior  $\pi \approx .33$  based on a frequentist approach; Instead we will be using a non-optimal prior  $\tilde{\pi} = .5$  as it is the application we are going to be targeting .

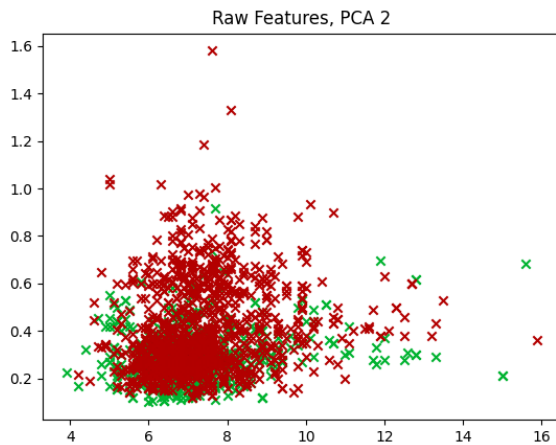
<sup>1</sup>K varies through models. For fast ones, 5 or 10 is used. For slower ones, 3 is used.

Coming back to features distributions, some things are to be noticed. While some features are similar between class Good and class Bad (see feature 1) others differ substantially and can be very helpful in discriminating samples (see feature 8).

Moreover, the features are distributed in various ways: while some do look pretty Gaussian (see feature 9) and others are instead fairly regular (see feature 5) and could thus be well estimated by Gaussian models<sup>2</sup>, other act in a more irregular way.

To take a closer look we now project the data on the two dimensional plane. We will be using 3 methods to do so: PCA<sup>3</sup>, Normalization + PCA, Normalization + Whitening.

Figure 2: 2D-PCA Projection

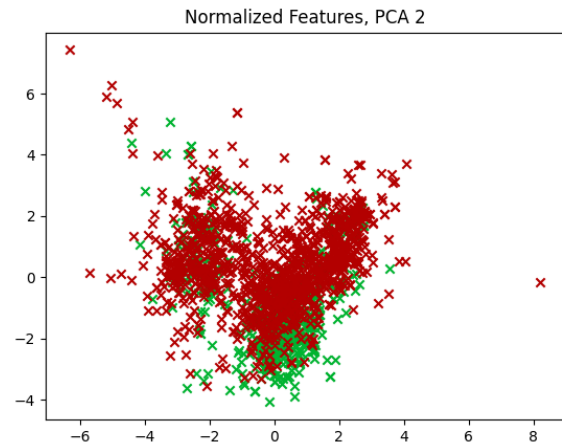


The points are very close one another, we thus expect linear models to be not so effective compared to others. The data is pretty *circularly* distributed, so correlation may not play an important role in discriminating samples.

<sup>2</sup>Meaning both Gaussian Classifiers and Gaussian Mixture Models

<sup>3</sup>Without data centering to appreciate the correlation

Figure 3: 2D-PCA Projection, Normalized Data

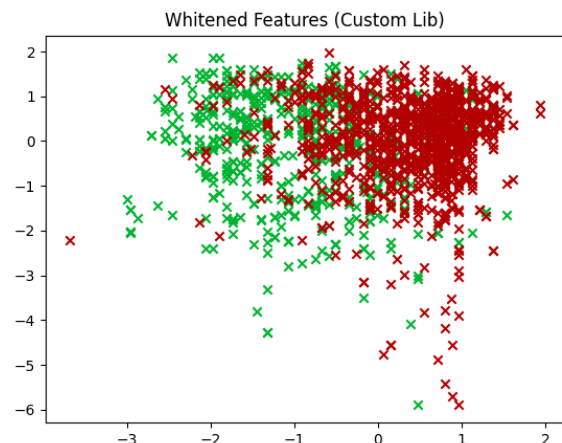


The normalized projection seems to split data in two clusters, each containing some samples of either class, but some points of different class seem to get far apart from the other. So normalization is a technique worth trying.

Note that for Normalization we refer to *Z-Normalization*. From some early tests we found that *Min-Max* normalization was not very effective, and with *Gaussianization* centering and scaling data in the same way as Z while being slower, we decided to use this method.

Lastly we take a look at normalized and whitened data.

Figure 4: 2D-PCA Projection, Whitened Data



Results do look interesting: points to get much

apart, even if we can spot many outliers. With the distribution being this way, we could expect even linear models to achieve decent results.

## 2 Pre-Processing Analysis

In this section we will run some dummy<sup>4</sup>models to infer wheter pre-processing, by the means of PCA, normalization, etc. can be useful or worth trying.

For now we will run all possible combinations of preprocessing (Raw, Normalized, Whitened) with PCA reductions (2-PCA, 3-PCA, ...).

### 2.1 Pre-Processing MVG Classifiers

Type	PCA	DCF	minDCF
Raw	9	0.401	<b>0.362</b>
Normalized	7	0.464	0.420
Whitened	2	0.430	0.410

Table 1: Full-Covariance MVG - Best Results

Type	PCA	DCF	minDCF
Raw	7	<b>0.375</b>	0.366
Normalized	9	<b>0.375</b>	0.371
Whitened	11	0.402	0.401

Table 2: Tied-Covariance MVG - Best Results

Type	PCA	DCF	minDCF
Raw	7	0.383	0.366
Normalized	10	0.429	0.391
Whitened	5	0.402	0.386

Table 3: Naive-Bayes MVG - Best Results

The result validate some of our assumptions and invalidate others. By looking at the DCF<sup>5</sup>values, our best model seems to be the Tied-Covariance ones. We are not very surprised to se that the

<sup>4</sup>Referring to models requiring parameter tuning or score recalibration

<sup>5</sup>By DCF we actually mean the *normalized* DCF, considering a prior  $\pi = .5$  and equal costs

Naive Bayes does sometimes outperform the Full-Covariance model as we noticed in Figure 2 that some features are not very correlated. The same can said about the the Tied covariance by looking at the same projection or feature 8 of Figure 1.

Normalization and Whitening do not help, they harm, sometimes even significantly the classification, while it looks like PCA can be of modest help.

An interesting fact that is that while the Full-Covariance model looks like the worse performing one, it is the model with the lowest minDCF values. While it is not an important difference compared with other models, it is the one with the biggest difference between the DCF and minDCF values, hinting that score calibration could be required.

### 2.2 Pre-Processing for GMMs

For GMMs we decide the number of components (for this dummy execution 4 was picked arbitrarily). We start with identity covariance matrices and by placing the means near the dataset mean.<sup>6</sup>

Type	PCA	DCF	minDCF
Raw	No	0.410	<b>0.394</b>
Normalized	9	<b>0.407</b>	0.402
Whitened	4	0.429	0.420

Table 4: GMM - Best Results

As we would expect after the results of the Gaussian Classifiers, also here whitening does not look like a good choice and normalization looks ineffective. High dimensionality still looks preferred.

The DCF values are close to the minDCF values, score calibration may be not necessary.

### 2.3 Pre-Processing for Polynomial Kernel SVMs

The polynomial kernel mapping  $\phi(\cdot)$  we will use is so defined:

<sup>6</sup>The starting points are used in the same way also in the optimization part.

$$\phi(\mathbf{x}) = \begin{bmatrix} \text{vec}(\mathbf{x}\mathbf{x}^T) \\ \sqrt{2}\mathbf{x} \\ 1 \end{bmatrix} \quad (1)$$

And we will be using the kernel function so described:

$$\begin{aligned} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) &= k(\mathbf{x}_1, \mathbf{x}_2) = \\ &= (\mathbf{x}_1^T \mathbf{x}_2 + c)^d + b = (\mathbf{x}_1^T \mathbf{x}_2 + 1)^2 + 0.5 \end{aligned}$$

The formulation of 2.3 would require us to strictly use  $c = 1$ , as we do in our dummy model, but later on this will be an hyperparameter to be optimized.

The results are interesting:

Type	PCA	DCF	minDCF
Raw	7	0.554	0.545
Normalized	6	<b>0.393</b>	<b>0.381</b>
Whitened	11	0.399	0.393

Table 5: Polynomail Kernel SVM - Best Results

While competitiveness with other models has to be made later, when each one will be fully used, here we see an interesting change: normalization greatly improves our classification. While a lower PCA has obtained the lowest DCF, it is not significantly lower than the one obtained with higher dimensionality, so we can say that here it's normalization helping us out. Whitening, again, does not make any difference, as the results are same or close to the ones obtained with normalization only.

Notice also how scores seems well calibration with our theoretical threshold even if they do not have a probabilistic interpretation.

## 2.4 Pre-Processing for RBF Kernel SVMs

The dummy RBF function used here is the following:

$$k(\mathbf{x}_1, \mathbf{x}_2) = e^{-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2} + b = e^{-0.05 \|\mathbf{x}_1 - \mathbf{x}_2\|^2} + 0.05$$

Small values of  $\gamma$  and  $b$  were picked to avoid numerical issues when using non-normalized models.

Type	PCA	DCF	minDCF
Raw	10	0.588	0.579
Normalized	11	<b>0.356</b>	<b>0.352</b>
Whitened	11	<b>0.356</b>	<b>0.352</b>

Table 6: RBF Kernel SVM - Best Results

Just like polynomial kernel, RBF prefers high dimensionality, so further testing in this direction is required. Again, normalization plays an important role while whitening seems to have no effect whatsoever and scores look extremely well calibrated.

Just like polynomial kernel SVMs, here normalization has no effect whatsoever while normalization greatly improves our performance. The RBF Kernel works significantly better when working with all or most components.

## 3 Optimizing Models

### 3.1 Optimizing MVG Classifiers

Since MVG Classifiers do not have parameters to be tuned, we are going to discuss only the possibility of recalibrating scores. In particular we will not discuss the Tied Covariance model, as based on its DCF values, the scores look already well calibrated.