

Machine Learning and Pattern Recognition Project Report

Paolo Magliano s314867

July 10, 2024

Abstract

This report describes the work done on Machine Learning and Pattern Recognition project. The report follows the structure of course laboratories, which studies and analyses the dataset and the models:

- Multivariate Gaussian Model
- Logistic Regression
- Support Vector Machines
- Gaussian Mixture Model

1 Lab 2: Dataset

The samples are computed by a feature extractor that summarizes high-level characteristics of a fingerprint image. The data is 6-dimensional and it consists of labeled samples corresponding to the genuine (True, label 1) class and the fake (False, label 0) class.

The dataset can be summarized by the image in Figures 1. It shows the data distribution of features and their correlation with each other.

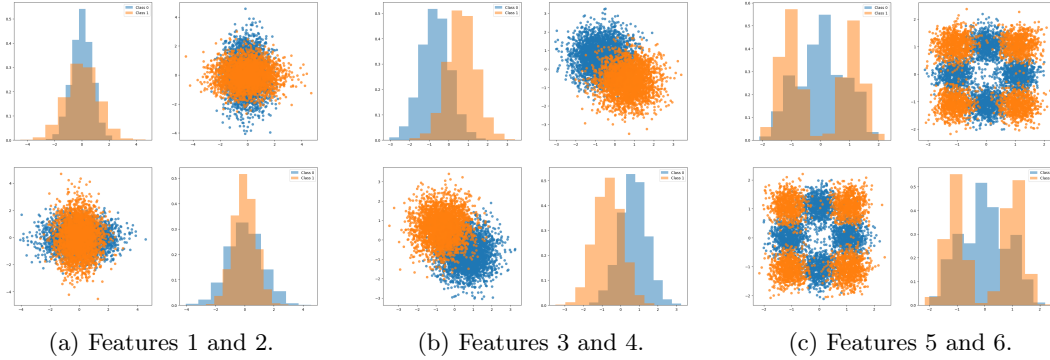


Figure 1: The figure shows histograms of the projected data on each feature. The scatter plot shows the correlation between the two features.

The first two features are mostly overlapping, so it is difficult to separate the two classes. The histograms show that they can be probably well approximated by a Gaussian density function. Further more means are very close to each other, but the variances are different.

The third and fourth features are more separated, but still overlapping. Also these features can be approximated by a Gaussian density function but the means are different and variance very similar. These features are more useful in a classification task as they are more discriminative.

The last two features are the most discriminative, they can't be approximated by a simple unimodal Gaussian density function because they seems to be bimodal or trimodal. Thanks to this behavior they form 4 distinct clusters (two for each feature) that can be easily separated. They are probably the most useful features if the classifier is able to capture the complex structure.

2 Lab 3: PCA and LDA

2.1 Dimensionality Reduction

Principal Component Analysis (PCA) is a dimensionality reduction technique that finds the directions of maximum variance in the data. It projects the data onto a lower-dimensional subspace while preserving as much variance as possible. The Figure 2 shows the results of PCA on the dataset. The first principal component captures most of the variance and it separates decently the two classes. The other principal components overlap, so they are not very useful for classification.

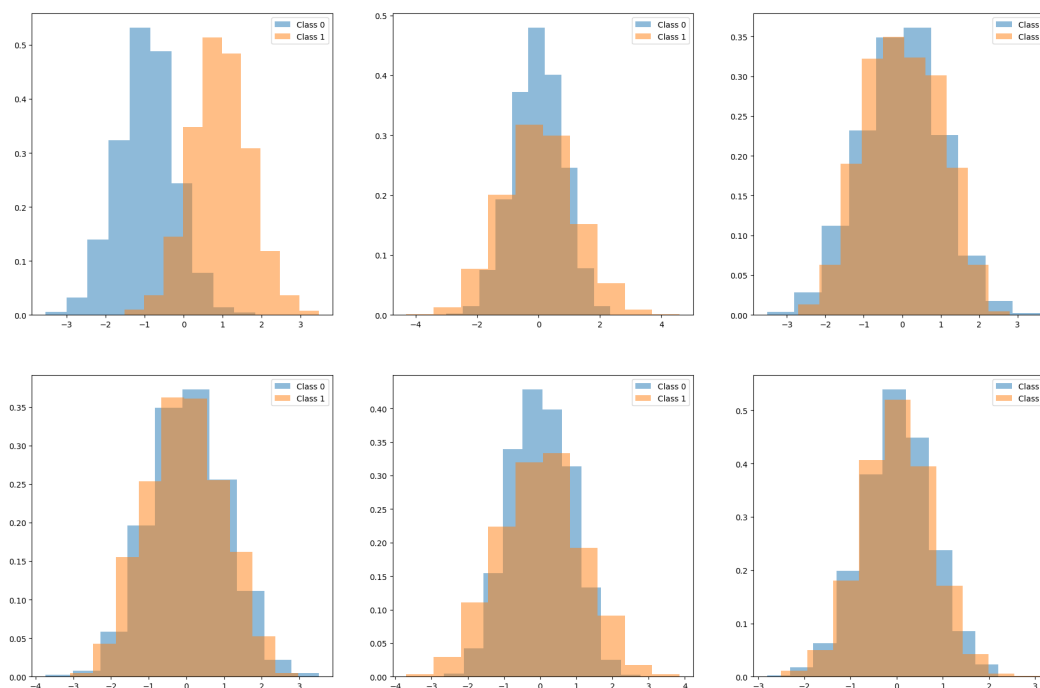


Figure 2: The figure shows the results of PCA on the dataset. The histograms show data distribution on each principal component. The top left one represents the first principal component, the bottom right is the last one.

Linear Discriminant Analysis (LDA) is a supervised dimensionality reduction technique that finds the directions that maximize the separation between classes. The Figure 3 shows the results of LDA on the dataset. Only one linear discriminant is found because the number of classes is 2. The linear discriminant separates the two classes partially. It is comparable to the first principal component of PCA.

Both cases simplify the dataset as much that the well-distinctive cluster of the last two features is lost. So, the preprocessing techniques could be useful to speed up the computation but they could also lose important information. In very simple dataset and this one, with only 6 dimensions, presumably the best choice is to use the original dataset.

2.2 Classification

The project analyses the performance of LDA method as a classifier. In order to get reliable results, the experiments' performance is evaluated with k-fold cross-validation with k equal to 10. It is also inspect the improvement of the classifier by using PCA as a preprocessing technique.

The Table 1 shows the performance results as error rate, Detection Cost Function (DCF) and Minimum DCF. The last one is useful to understand the gains of changing the threshold to optimal value.

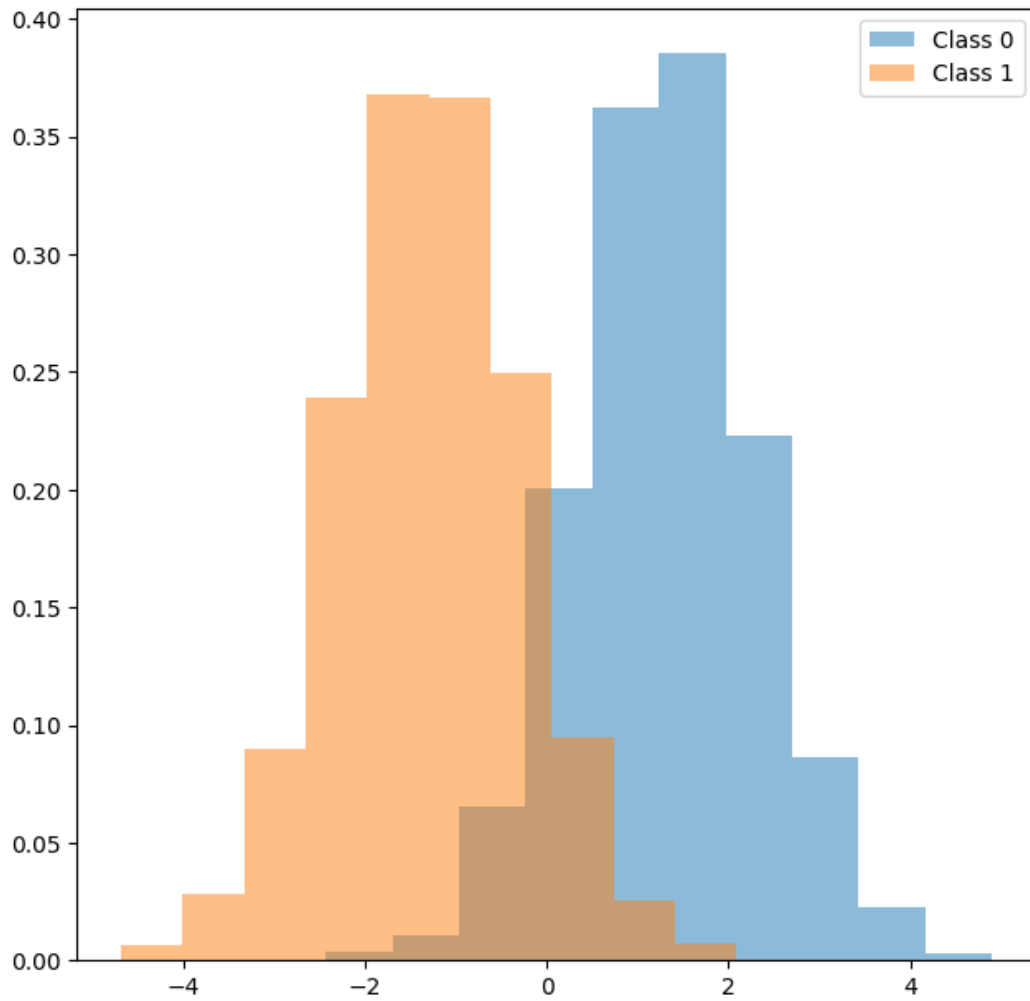


Figure 3: The figure shows the results of LDA on the dataset. The histogram shows data distribution on the only linear discriminant found.

PCA	Error Rate (%)	DCF	Min DCF
None	9.58	0.191	0.177
1 PC	9.68	0.193	0.179
2 PC	9.63	0.193	0.179
3 PC	9.43	0.189	0.177
4 PC	9.57	0.191	0.177
5 PC	9.58	0.192	0.177
6 PC	9.62	0.192	0.177

Table 1: Performance metrics for LDA models with various dimensionality reduction techniques.

As intuitively retrieved from PCA and LDA graphs, the capabilities of both methods are similar. In fact, the performance of the classifier is not improved by using PCA as a preprocessing technique.

In general the performance are terrible considering the simplicity of the dataset. Further more the threshold selection is not optimal, so the performance could be improved by tuning it a little bit.

3 Lab 4: Multivariate Gaussian density

As preliminary step, the project analyses the correspondence between the features distribution and the Gaussian density function. This is helpful to understand the potential efficacy of model based on Gaussian density. The figure 4 shows the histograms of the dataset and the Gaussian density function that best fits the data.

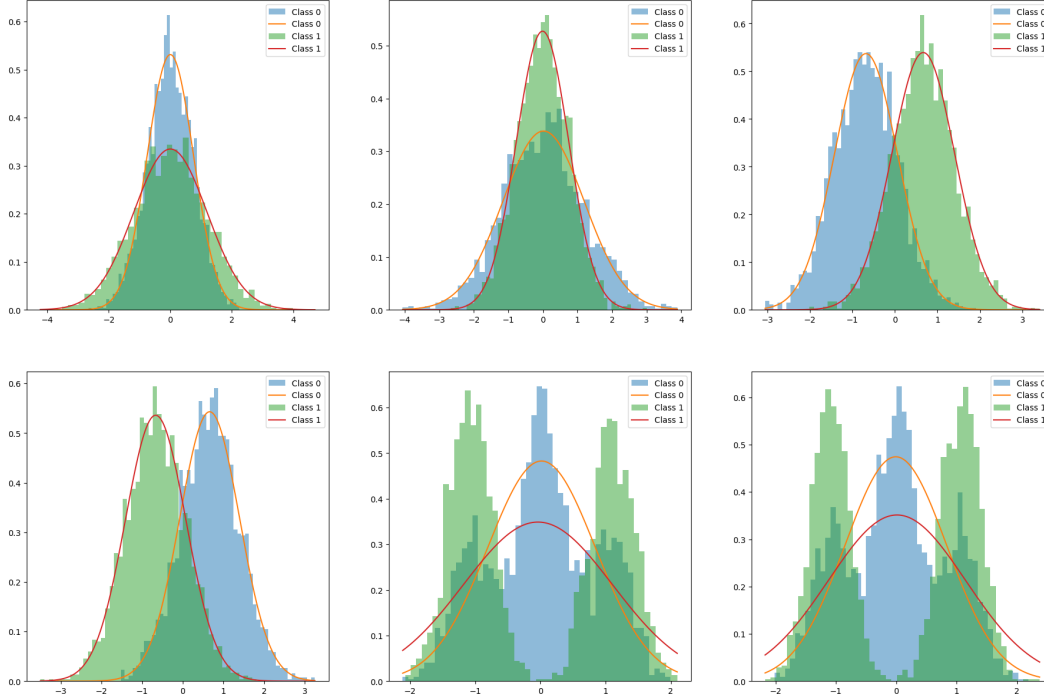


Figure 4: The figure shows the histograms of the dataset and the Gaussian density function that best fits each feature.

The results are consistent with the previous analysis. The first four features can be approximated by a unimodal Gaussian density function. Instead, the last two can not be well represented thanks to their multimodality.

4 Lab 5: Multivariate Gaussian Model

The project analyses the performance of Multivariate Gaussian (MVG) models with different covariance matrices: full, tied, and diagonal. The experiments are conducted with k-fold cross-validation with k equal to 10. The results are shown in Table 2.

Model	Error Rate (%)	DCF	Min DCF
MVG	7.45	0.148	0.133
TiedMVG	9.62	0.192	0.177
NaiveMVG	7.48	0.149	0.132

Table 2: Performance metrics for various MVG models without dimensionality reduction.

The results show that the tied MVG performs worse than the others. This is due to the fact that the covariance matrix change considerably along the features and so it can't be approximated by a single matrix. To visualize this better, the Figure 5 shows the estimated Gaussian density function for the tied MVG model.

The graphs show that the tied MVG model can practically use only features 3 and 4 as discriminative ones. The other features are approximated with the same gaussian function

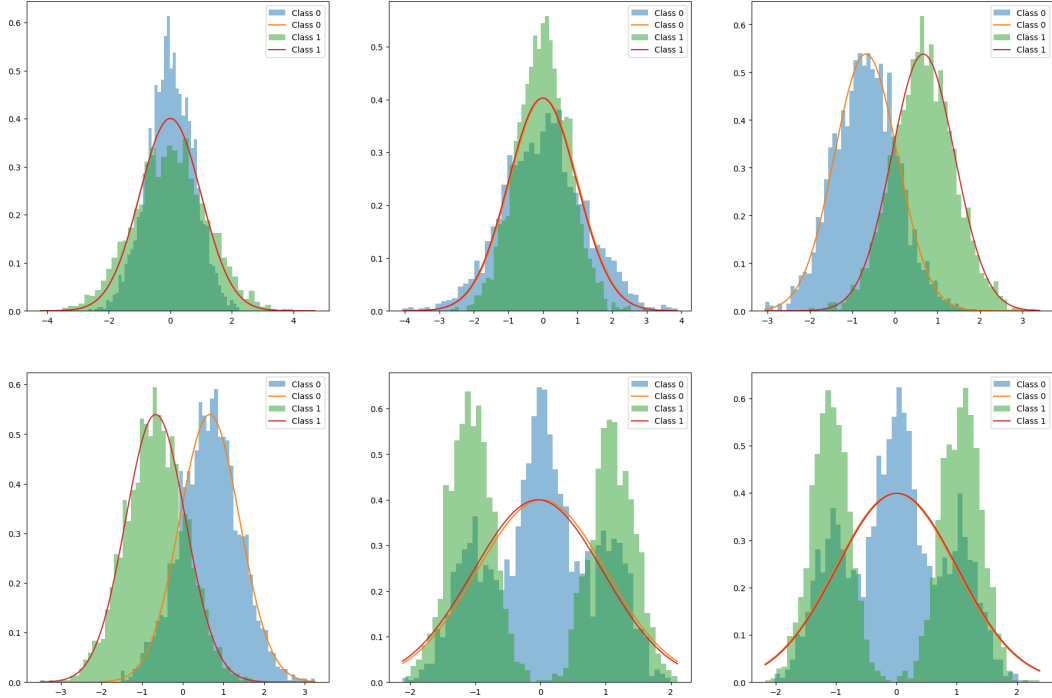


Figure 5: The figure shows the histograms of the dataset and the Gaussian density function with tied covariance matrix for each feature.

from both classes.

4.1 Correlation Analysis

In order to analyze further the performance of Naive MVG model respect to the full MVG model, the project computes the pearson correlation coefficient between the features for each class. The results are shown in Figure 6.

The graphs show the reason why the Naive MVG model performs as well as the full MVG model. The features are not correlated at all, so the assumption of independence is not violated.

4.2 Assumption Analysis

The Gaussian model assumes that features can be jointly modeled by Gaussian distributions. The good-ness of the model is therefore strongly affected by the accuracy of this assumption. As already discussed, the last two features do not satisfy this assumption. To analyze if indeed the last set of features negatively affects our classifier, we can try repeating the classification using only feature 1 to 4. The results are shown in Table 3.

Model	Error Rate (%)	DCF	Min DCF
MVG	8.35	0.167	0.149
TiedMVG	9.73	0.194	0.176
NaiveMVG	8.35	0.167	0.149

Table 3: Performance metrics for various MVG models with only the first 4 features.

For both MVG and Naive MVG models, the performance is slightly worse than the full dataset. This implies that the last two features are useful anyway, even if they are not well approximated by a Gaussian.

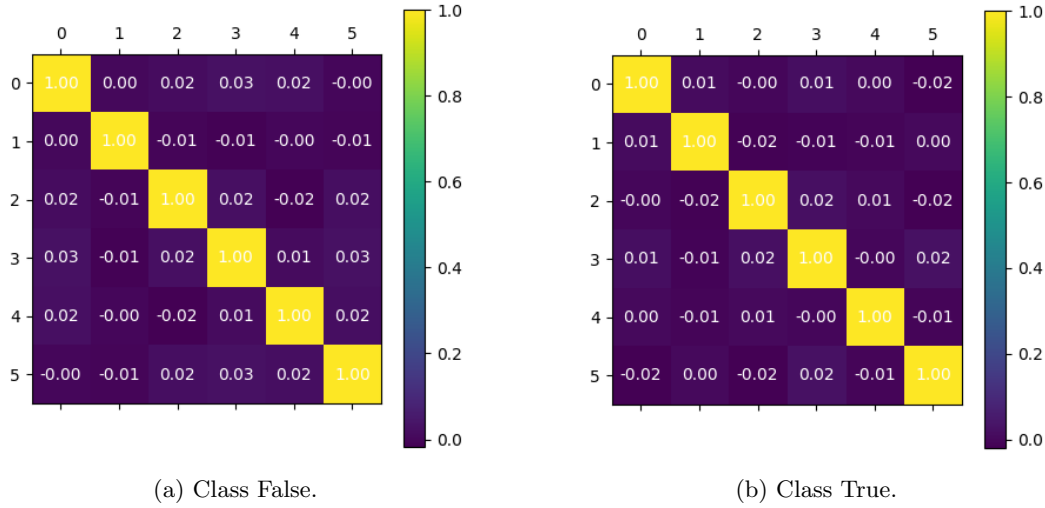


Figure 6: The figure shows the Pearson correlation coefficient between the features for each class. Correlated features are close to 1 or -1.

4.3 Mean and Variance of Features Analysis

To further explore how means and variance of the features affect the performance of different approaches, the project repeats the classification using only features 1 and 2 and then only features 3 and 4. This is done because features 1 and 2 means are similar but variances are not, whereas for features 3 and 4 the two classes mainly differ for the feature mean, but show similar variance. The results are shown in Table 4 and Table 5.

Model	Error Rate (%)	DCF	Min DCF
MVG	36.83	0.737	0.711
TiedMVG	49.55	0.99	0.887
NaiveMVG	36.80	0.737	0.711

Table 4: Performance metrics for various MVG models with only the first 2 features.

Model	Error Rate (%)	DCF	Min DCF
MVG	9.67	0.193	0.176
TiedMVG	9.68	0.194	0.177
NaiveMVG	9.68	0.194	0.177

Table 5: Performance metrics for various MVG models with only the 3 and 4 features.

Lets first analyze the generale performance between the two cases. The performance of the first two features is terrible, this is due to the fact that the two classes are not separable at all and this shows how important is the mean of class distribution. Instead, the performance of features 3 and 4 is very good, this is due to the fact that the two classes mean are very different and this remark how much these features are discriminative.

More over, in the first case Tied MVG model performs even worse because the features' variances are unsimilar. Instead in the second case, the performances are the same for all models because the features' variances are pretty the same.

4.4 Dimensionality reduction

Finally the project investigates the performance of MVG models with PCA as a preprocessing technique. The results are shown in Table 6.

Model	Metrics	PCA 1	PCA 2	PCA 3	PCA 4	PCA 5	PCA 6
MVG	Error	9.7%	9.23%	8.78%	8.33%	7.38%	7.45%
	DCF	0.194	0.185	0.175	0.166	0.147	0.148
	Min DCF	0.179	0.173	0.161	0.151	0.133	0.133
TiedMVG	Error	9.68%	9.63%	9.45%	9.58%	9.57%	9.62%
	DCF	0.193	0.193	0.189	0.191	0.191	0.192
	Min DCF	0.179	0.177	0.177	0.177	0.177	0.177
NaiveMVG	Error	9.7%	9.32%	9.3%	9.02%	8.58%	8.6%
	DCF	0.194	0.186	0.186	0.18	0.172	0.172
	Min DCF	0.179	0.174	0.166	0.162	0.158	0.158

Table 6: Performance metrics for various MVG models with PCA as a preprocessing technique.

The results confirm the hypothesis done in Dimensionality Reduction section. The performance of the classifier is not improved by using PCA as a preprocessing technique due to the lostness of the information. Indeed the best results are obtained with higher number of principal components. Tied MVG model is the only one that performs the same with and without PCA, this is due to PCA that provides features with similar variances.

5 Lab 7: Model Evaluation

The project proposes various metrics and tools to evaluate the performance of the models. The Detection Cost Function (DCF) is a metric that combines the error rate and the cost of false alarms and misses. The Minimum DCF is the DCF value at the optimal threshold.

5.1 Different Application

In this section the project investigates the performance of the models in different application setups where the prior, the cost of false positives and the cost of false negatives are different. Specifically:

- $(\pi_T, C_{FP}, C_{FN}) = (0.5, 1, 1)$
- $(\pi_T, C_{FP}, C_{FN}) = (0.1, 1, 1)$
- $(\pi_T, C_{FP}, C_{FN}) = (0.9, 1, 1)$
- $(\pi_T, C_{FP}, C_{FN}) = (0.5, 1, 9)$
- $(\pi_T, C_{FP}, C_{FN}) = (0.5, 9, 1)$

The image 7 shows the confusion matrix for the different application setups for the MVG model.

The confusion matrixes emphasize how the model choices depend on the prior and the cost of false positives and false negatives. For example, in the case $(0.1, 1, 1)$ the model is biased towards the negative class, so the false negatives are very low. Instead, in the case $(0.9, 1, 1)$ the model is biased towards the positive class, so the false positives are very low.

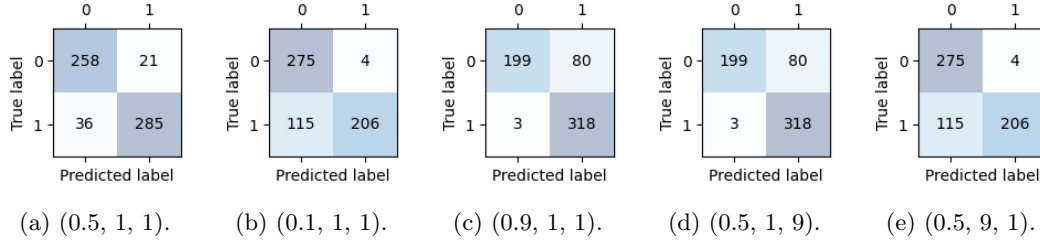


Figure 7: The figure shows the confusion matrix for the different application setups for the MVG model.

5.2 Application Evaluation

The project evaluates the performance of the models in the different application setups:

- $(\pi_T, C_{FP}, C_{FN}) = (0.5, 1, 1)$
- $(\pi_T, C_{FP}, C_{FN}) = (0.1, 1, 1)$
- $(\pi_T, C_{FP}, C_{FN}) = (0.9, 1, 1)$

The results for (0.5, 1, 1) are already presented in Table 2 and Table 6. The results for the other two setups are shown in Table 7. and Table 8.

Model	Metrics	None	1 PC	2 PC	3 PC	4 PC	5 PC	6 PC
MVG	Error	13.68%	17.77%	16.82%	16.12%	15.15%	13.72%	13.68%
	DCF	0.37	0.481	0.457	0.433	0.419	0.378	0.37
	Min DCF	0.308	0.413	0.409	0.381	0.355	0.312	0.308
TiedMVG	Error	17.42%	17.85%	17.55%	17.5%	17.4%	17.42%	17.42%
	DCF	0.471	0.483	0.474	0.471	0.472	0.475	0.471
	Min DCF	0.408	0.413	0.408	0.402	0.408	0.409	0.408
NaiveMVG	Error	13.97%	17.77%	16.98%	16.3%	15.9%	15.42%	15.38%
	DCF	0.373	0.481	0.466	0.439	0.431	0.41	0.409
	Min DCF	0.31	0.413	0.41	0.393	0.369	0.363	0.362

Table 7: Performance metrics for various MVG models with (0.1, 1, 1) application setup.

Model	Metrics	None	1 PC	2 PC	3 PC	4 PC	5 PC	6 PC
MVG	Error	14.05%	17.6%	16.25%	15.8%	14.87%	14.03%	14.05%
	DCF	0.367	0.464	0.45	0.439	0.399	0.364	0.367
	Min DCF	0.319	0.426	0.404	0.381	0.366	0.323	0.319
TiedMVG	Error	17.05%	17.5%	17.4%	17.2%	17.13%	17.03%	17.05%
	DCF	0.46	0.467	0.468	0.461	0.463	0.458	0.46
	Min DCF	0.423	0.426	0.426	0.417	0.423	0.425	0.423
NaiveMVG	Error	14.15%	17.6%	16.47%	16.05%	15.9%	15.57%	15.57%
	DCF	0.364	0.464	0.451	0.439	0.433	0.411	0.411
	Min DCF	0.32	0.426	0.403	0.397	0.384	0.359	0.36

Table 8: Performance metrics for various MVG models with (0.9, 1, 1) application setup.

The results show that the performance of the models is strongly affected by the application setup and all the models suffer of big miscalibration error. Over all setups MVG model with full covariance matrix performs better than the others. As can be seen from min DCF

values, all models can improve their performance by changing the threshold, despite they will never reach the value of (0.5, 1, 1) setup. The only setup well calibrated is (0.5, 1, 1) because the dataset is almost balanced according to the prior.

5.3 Bayes Error Plot

The project proposes a Bayes error plot to visualize the performance of the best MVG models: with PCA 6 for (0.1, 1, 1) application setup. The plot in Figure 8 compare the actual DCF and the min DCF for different application priors, understanding how well the model behaves in different scenarios.

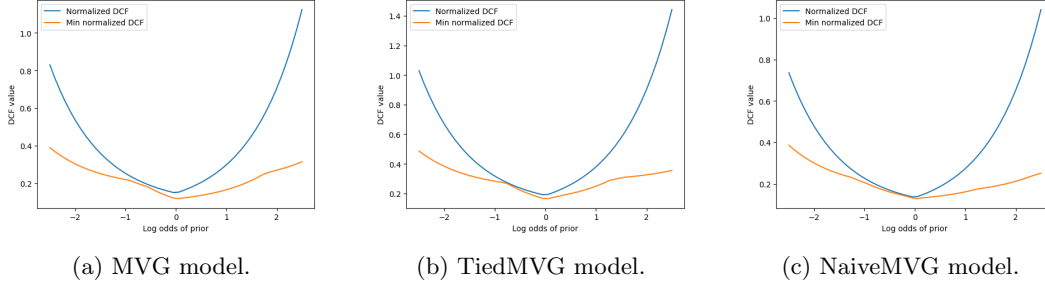


Figure 8: The figure shows the Bayes error plot for the MVG models with PCA 6 for (0.1, 1, 1) application setup.

The Bayes error plot confirm the previous hypothesis on miscalibrated models. In fact, the actual DCF is often far from the min DCF and only in few cases the values are close to each other.

6 Lab 8: Logistic Regression

6.1 regularization

The project firstly investigates the performance of Logistic Regression (LR) model as it varies the regularization parameter λ . The experiments are conducted with k-fold cross-validation with k equal to 10. The Figure 9a shows DCF and min DCF for different values of λ .

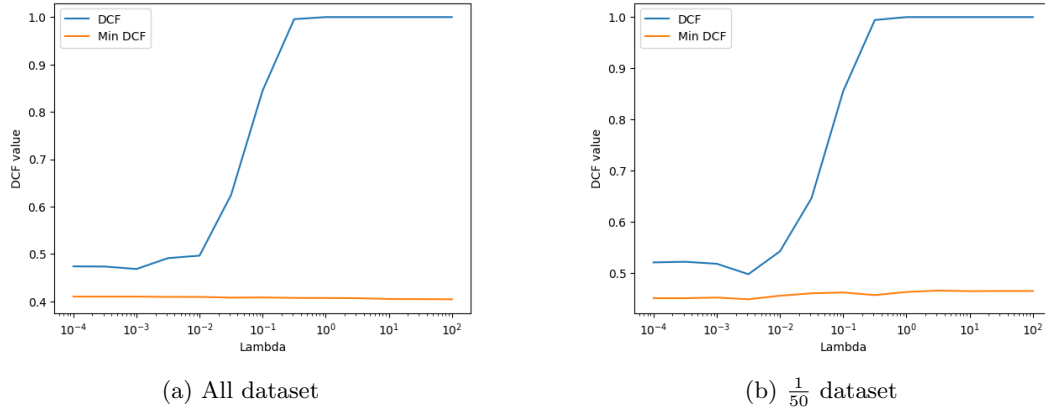


Figure 9: The figure shows the performance of LR models with different λ .

Since we have a large number of samples, regularization seems ineffective, and actually degrades actual DCF since the regularized models tend to lose the probabilistic interpretation of the scores. To better understand the role of regularization, the project analyses the results that we would obtain if we had fewer training samples. So, it repeats the previous

analysis, but keep only 1 out of 50 model training sample and the results are shown in Figure 9b.

7 Lab 9: Support Vector Machines