

Machine Learning and Pattern Recognition: final report

Claudio Macaluso, s317149

February 5, 2025

1 Preliminary hypothesis: a visual inspection of data

All the analyses in the initial sections are based on graphical evaluations made by observing the scatter plots and histograms of the data (Figure 1). Therefore, here we rely on purely qualitative assessments.

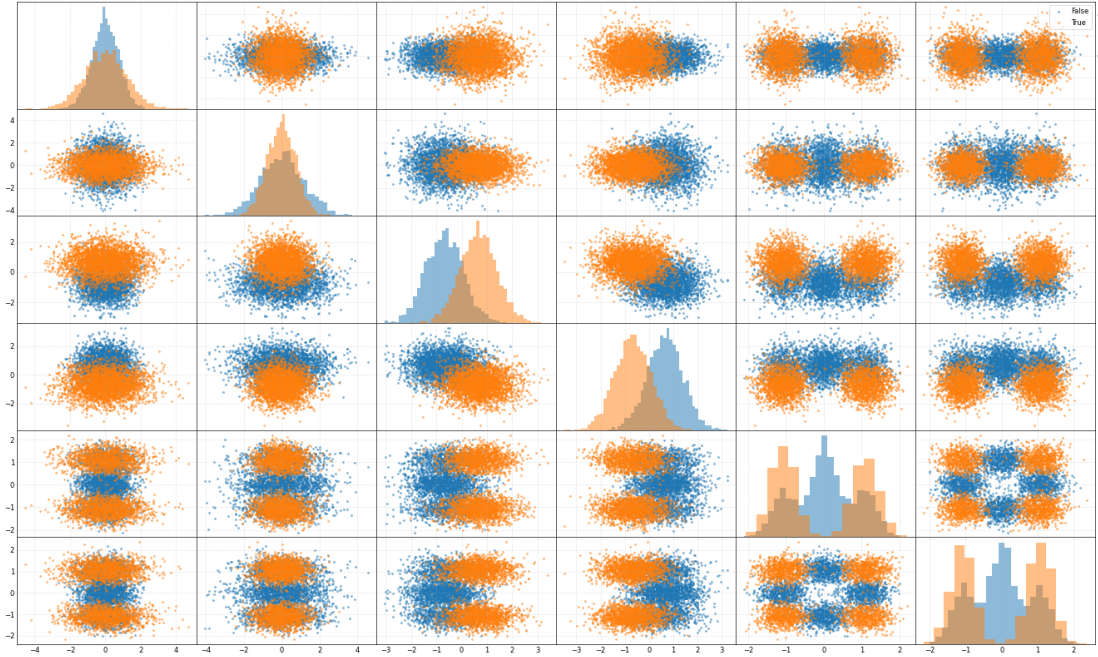


Figure 1: Spoofing dataset scatter matrix

Looking at the plots, we can notice that features can be grouped given their common characteristics.

Beginning with the first two features, they are both unimodal: this characteristic can be useful if we use some model that assume an underlying distribution that is inherently unimodal (e.g. Gaussian models). Unfortunately, as we can see they show a significant overlap, which may reduce their discriminative power compared to other features. Despite this, these two features combined with the other ones can anyway contribute positively to the classification task. Another thing that have to be taken into account, is that the variances are different for the two classes and this can lead to a decrease in performance if we use some model that assume the same variance for all the classes (e.g. Tied Covariance Gaussian Model).

Moving on the last two, they seem less overlapped compared to the previous ones, but in this case they're multimodal: multimodality can be challenging, specifically if, as we said, the model we use assumes unimodal distribution (e.g. Gaussian Models). In the case of these two features, an extremely simple model could poorly perform. Since the multimodality of the features, taking into account mean and variance of the classes might not be significative anymore: mean doesn't represent any actual cluster of data points and, since variance represent spread around mean, we can't rely on the variance as a good parameter of spread of data neither. In the end, for this two features, we can expect that models assuming unimodal distribution (e.g. Gaussian Models) perform poorly

compared to others. Moreover, looking at the scatter plots, we can imagine that a quadratic separation surface (e.g. quadratic Logistic Regression, Support Vector Machine using non-linear kernels) could lead to better results than a linear one (e.g. Tied Covariance Gaussian Model, linear Logistic Regression or Support Vector Machine using linear kernel).

Proceeding to the third and the fourth features, they are unimodal and have a lower overlap, so they are well suited for models that assumes unimodal distribution. Besides, they are the most "discriminative" features overall, since they show less overlap with respect to the other features. We can confirm this statement analysing LDA directions: we can notice, from (1), that the found vectors highly correspond to the directions of the third and the fourth features, since LDA tend to select linear combination of the features in a way that make classes more distinguishable in the projected data, confirming the hypothesis for which these two features lead to a more discriminative separation of classes.

$$\mathbf{v} = \begin{bmatrix} -0.0106 \\ 0.0134 \\ -0.9657 \\ 0.9677 \\ 0.0218 \\ -0.0329 \end{bmatrix} \quad (1)$$

Furthermore, since these features have similar class variance, it could be suitable to apply some models that make the assumption of same covariance among classes (e.g. Tied Covariance Gaussian Model). Last thing about the third and fourth features, as we can see from scatter plots, the two are correlated, so using models assuming non correlated features (e.g. Naive Bayes Gaussian Models) can lead to a decrease in performance for these two features, so we will need to take care about this in the further analyses.

2 Dimensionality reduction techniques

Now, let's evaluate the scatter plots and histograms of data projected using dimensionality reduction techniques.

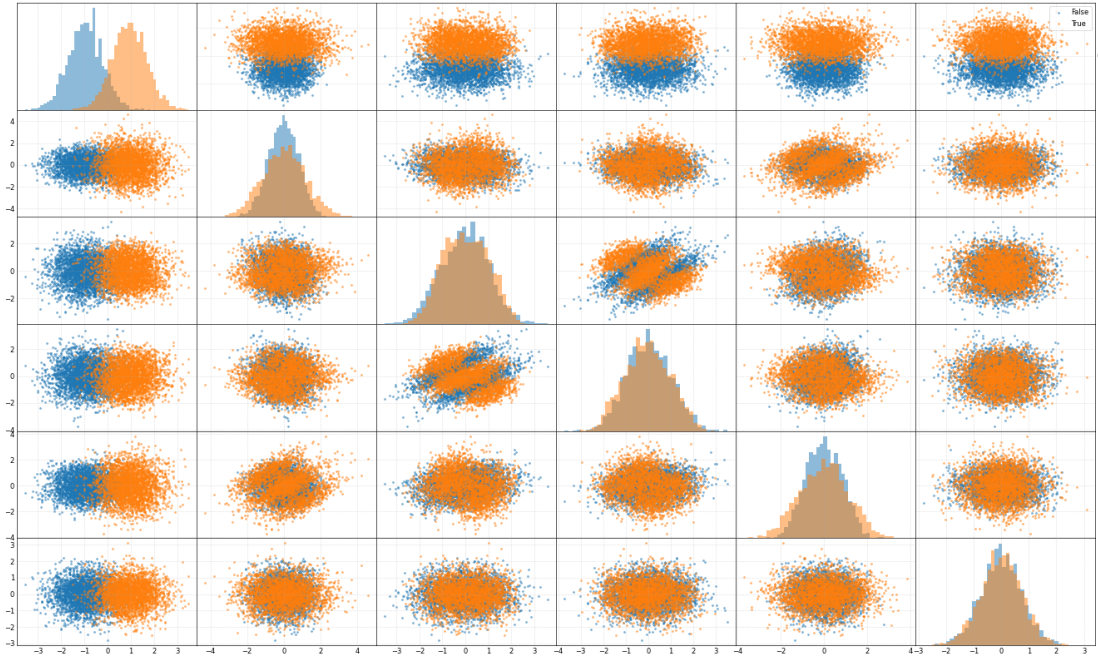


Figure 2: Spoofing dataset scatter matrix projected according to PCA directions

As we expected, since PCA is a non supervised dimensionality reduction technique (i.e. class agnostic), the found directions seem to not separate well classes, except for the first feature. From a first analysis to the scatter plots of the PCA projected data, we can see that, except for the first component, the others seem to show high

overlap between classes. Calculating the means and variances of classes in the projected space, we can see that the first component is the only one with a consistent difference in the mean of the classes. For the variances of the classes instead, we can notice that second and fifth features seem those with the biggest difference, so, disregarding the fact that visually they appear to have a high degree of overlap, they could still be useful for the classification task. A possible benefit of PCA projection in this case could be that, from the histograms, we can notice that all the new features have unimodal distribution, so some models could improve their performance due to these new characteristics of the projected data. This last one is just a coincidence, since PCA objective isn't inherently searching for a projection space in which features are unimodal, but still can be taken into account as an improvement in the new features space. Finally, as we expect, we can notice, calculating Pearson correlation matrix, that features are decorrelated in this new space (considering whole dataset, without class division), since projecting using PCA has, as effect, the decorrelation of features.

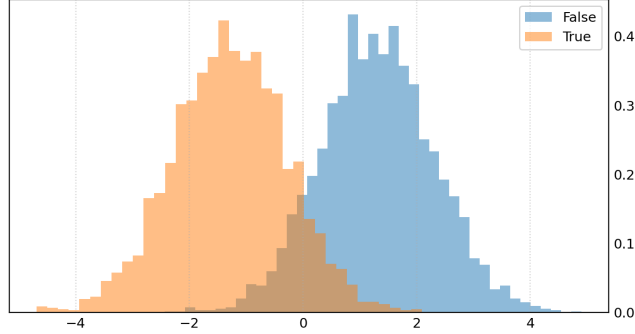


Figure 3: Spoofing dataset scatter matrix projected according to LDA directions

Clearly, a better separation can be achieved by applying LDA dimensionality reduction technique, that, unlike the previous one, is a supervised technique that aims to find the direction in which to project data for which the separation between classes is maximized.

Setting up a simple classification rule for 1-dimensional data that uses the mean of the means of the classes as a threshold to classify each sample, we can begin to check the impact of dimensionality reduction techniques on classification on Table 1.

Table 1: Error rate with decreasing number Principal Components using a simple LDA mean based classifier

Dimensions	Error Rate	Explained Variance
RAW	9.30%	-
PC6	9.30%	100.00%
PC5	9.25%	90.68%
PC4	9.40%	74.29%
PC3	9.25%	57.69%
PC2	9.05%	40.94%
PC1	9.25%	23.99%

With this configuration, we can conclude that, looking just at error rate metrics as an indicator of the goodness of the dimensionality reduction hyperparameter m used in the context of this simple mean based classifier, the optimal value is 2. Using this configuration, retained variance is equal to 40.94%, and it is enough to use this 2-dimensional representation to capture the most significant patterns in the data while discarding less relevant variability.

3 Models evaluation

Now we can begin to analyse various models and their performance for our specific dataset. For each model, we'll write first of all, what do we expect, and then evaluate their effective performance using some chosen metrics.

3.1 Gaussian Models

3.1.1 Early analysis

Gaussian models are a group of probabilistic generative models that try to estimate the probability distribution of data starting from a representative dataset provided at training time. In particular, they assume that data are distributed following a Gaussian distribution, using as parameters the class-specific mean and variance. Decision boundaries are formed by comparing the likelihood of data points under these found Gaussian distributions, combined with prior probabilities to apply Bayes' theorem to find scores used in classification.

The strong assumption of Gaussian distributed data is a simplification used to deal with a well known distribution probability density function and all of its properties: due to this approximation, we expect these model to perform well in cases in which the underlying real data distribution is actually near to the Gaussian one. We can try to visualize the goodness of the hypothesis plotting, above data histograms, the Gaussian distribution obtained by using the mean and the variance of each feature for each class.

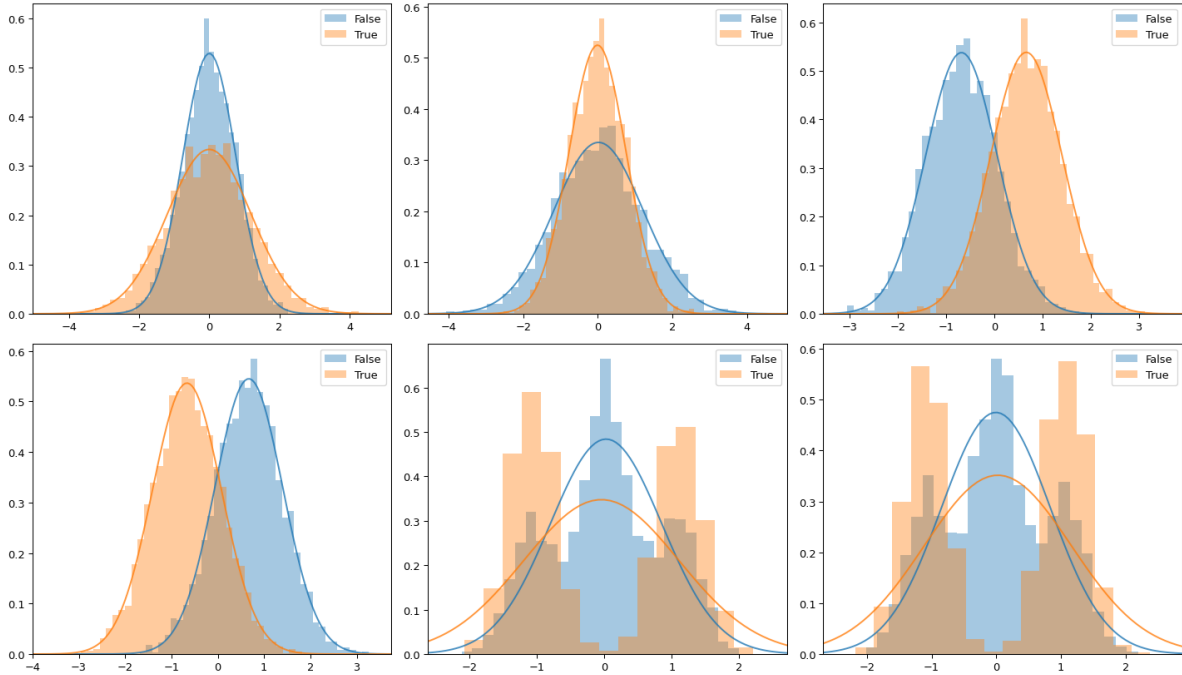


Figure 4: Per class features histograms of the spoofing dataset and their respective gaussian distribution, found using the mean and the variance of the per class features data

From Figure 4, we can see that for most of the features, at least visually, the Gaussian assumption is largely satisfied, with the exception of the last two features: in fact, as we said in the data visual inspection section, those feature are multimodal and bad fits with models that assume an unimodal probability distribution (i.e. Gaussian models).

But, as we stated during the dimensionality reduction techniques analysis, we noticed that, applying PCA to data make them turn into unimodal distributions, so let's try it and see if the assumption is better satisfied with projected data.

Indeed, as we can see from Figure 5 the features appear now all unimodal, and the class Gaussian curves seem to fit better in this case, even if, as we said, the price to pay is an higher degree of overlap of the features overall: it's likely that, even if the Gaussian assumption is now satisfied, the PCA don't provide any relevant benefit, since, for the PCA projected features excluding the first one, classes are highly overlapped and, for this reason, quite not distinguishable each other. We will evaluate the consequences of all these reasoning quantitatively in further analyses to check if PCA is really convenient to apply or not and how it impact on Gaussian models.

3.1.2 Quantitative evaluation

As we can see from Table 2 by now, using uniform prior and uniform cost matrix, Multivariate Gaussian Models outperform the others.

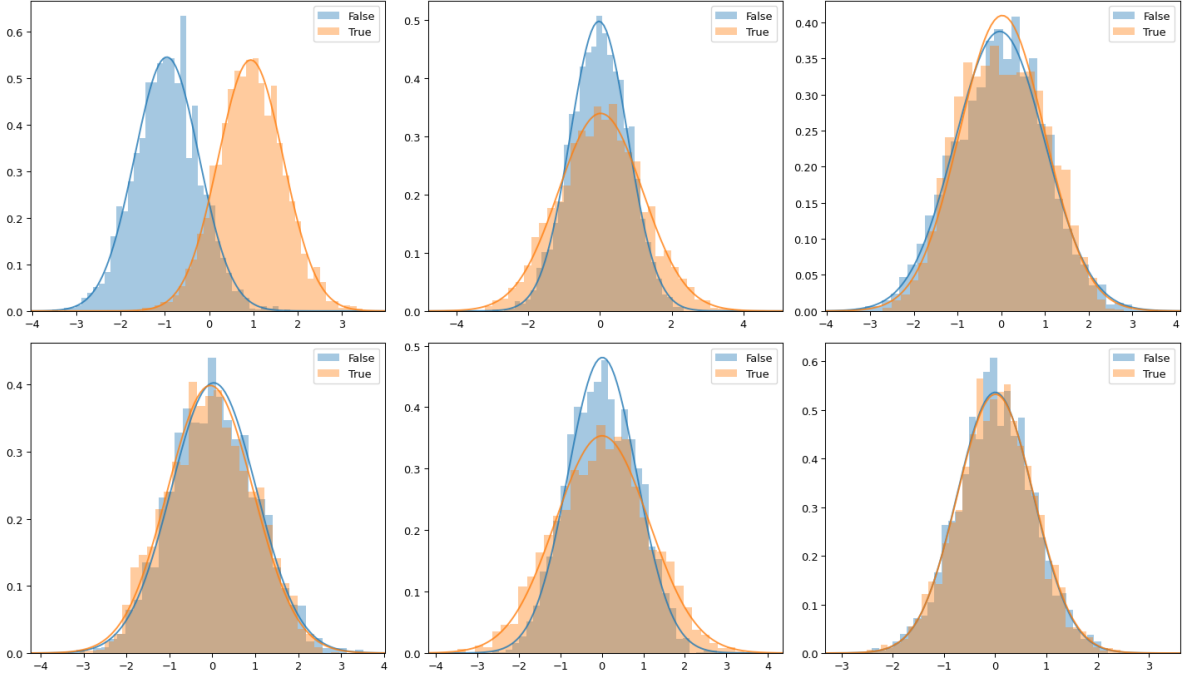


Figure 5: Per class features histograms of the PCA projected spoofing dataset and their respective gaussian distribution, found using the mean and the variance of the per class features data

Table 2: Table showing the error rate for Gaussian models

Model	Error Rate
Multivariate	7.00%
Naive	7.20%
Tied	9.30%
NaiveTied	9.30%
LDA	9.30%

We can try to evaluate error rate using subsets of features, and, in light of the results, draw conclusions about data and models behaviour.

Table 3: Error rate for Gaussian models, using subsets of features

Model	Error Rate / Features subset		
	1-2	3-4	1-4
Multivariate	36.50%	9.45%	7.95%
Naive	36.30%	9.45%	7.65%
Tied	49.45%	9.40%	9.50%
NaiveTied	49.40%	9.40%	9.50%
LDA	49.45%	9.40%	9.50%

The results in Table 3 shows that, as we analysed in the previous section, discarding features that badly fits Gaussian assumption lead to a slightly worse, but still comparable to the results obtained with the full set of features: this is due probably to the fact that, while still providing some useful information, the Gaussian assumption is a strong assumption for the last two features, since, as we said, their distribution is far from a Gaussian one.

Moreover, we can notice that, while using first two features lead to bad results overall, Naive model outperform the others: infact, if we analyse the correlation matrix, we can notice that first two features are the less correlated pair of all. Since Naive Bayes hypothesis assumes that features are not correlated, probably for this subset the hypothesis lead to a better understanding of underlying pattern, and not considering correlation allows to simplify the built data representation in a way that help the classification task. For Tied models and LDA, we can observe that they perform poorly for this subset, probably due to its poor separability.

For the second and the third feature, we obtain overall good results: this confirm us, as we stated in the first

section and proved analysing LDA direction results, that these two features are the most discriminant and have enough information alone to be considered the most useful in the classification task. Here the tied assumption works better than other subsets probably due to the separability of the features.

Table 4: Error rate for Gaussian models, using a decreasing number Principal Components

Model	Error Rate / PCA					
	m=6	m=5	m=4	m=3	m=2	m=1
Multivariate	7.00%	7.10%	8.05%	8.80%	8.80%	9.25%
Naive	8.90%	8.75%	8.85%	9.00%	8.85%	9.25%
Tied	9.30%	9.30%	9.25%	9.25%	9.25%	9.35%
NaiveTied	9.20%	9.20%	9.15%	9.10%	9.35%	9.35%
LDA	9.30%	9.30%	9.25%	9.25%	9.25%	9.35%

Finally, using PCA to preprocess data, we can see from Table 4 that, for Multivariate Gaussian Model, its performance degrades as the components taken into account decrease. Naive Bayes Model shows consistent performance across all PCA components, but still doesn't perform better than MVG. Tied covariance model exhibit minimal sensitivity to the reduction in dimensions, maintaining error rates around 9.25-9.35%. Their performance does not improve significantly even with all features.

From this results we can conclude that PCA captures the principal directions of variance in the data, but it does not necessarily align these directions with the features most useful for classification. This can explain why reducing the dimensionality through PCA degrades performance for the Multivariate model, which can effectively leverage the original feature structure to obtain optimal performance.

Now, as suggested in the provided text, we consider 5 different application for our analysis: $(\pi = 0.5, C_{fn} = 1.0, C_{fp} = 1.0)$, $(\pi =, C_{fn} = 1.0, C_{fp} = 1.0)$, $(\pi = 0.1, C_{fn} = 1.0, C_{fp} = 1.0)$, $(\pi = 0.5, C_{fn} = 1.0, C_{fp} = 9.0)$ and $(\pi = 0.5, C_{fn} = 9.0, C_{fp} = 1.0)$. In the analysis of the five applications provided, the effective prior $\tilde{\pi}$ was computed using the formula:

$$\tilde{\pi} = \frac{\pi_T C_{fn}}{\pi_T C_{fn} + (1 - \pi_T) C_{fp}}$$

Notably, the last two applications are equivalent to other two in terms of effective prior. This implies that, while the raw priors and costs differ, their impact on decision-making thresholds is identical, so we can discard them.

As we can see from Table 5, comparing minDCF of different models for different PCA configuration, we can see that the best model performance are obtained for the first two application by Multivariate Gaussian model without PCA and with PCA using 6 Principal Components, while for the last application the best value for minDCF is obtained by Naive Gaussian model without any preprocessing technique. Comparing minDCF values we can conclude that the first application is the one in which we achieve best performance, since values vary little from each other and the lowest value found for this application is even the lowest overall. The models performance for this application are also comparable each other. Overall, the performance of the first application are the best from the point of view of minDCF: values vary little from each other, and the best model have the smallest minDCF of all. For the other two, values for minDCF vary more than the first one, and they are relatively higher than the first one, with the third one having slightly lower values than the third one.

Comparing actDCF with minDCF values to check calibration, we notice that for the first application the calibration appear to be comparable between the different models, with Naive Bayes obtaining the best result using PCA with 5 Principal Components. For the second and the third one, the calibration loss appear to be comparable among models, with a single exception for the second application in which the only configuration that outperform the others is for the Multivariate Gaussian model trained using PCA with 2 Principal Components. For the third application, the order of magnitude of the calibration loss is higher than the second one: since the application that we're going to use will be this, we'll probably need to calibrate scores to obtain optimal performance from this set of models.

In the following analysis, we're going to see that, among all, the Gaussians, as well as Gaussian Mixture, are almost already well calibrated compared to the other models, since their scores are probabilistically interpretable without any post-processing.

Table 5: Table showing the DCF and minDCF using Gaussian models for the effective priors $\tilde{\pi} = 0.5$, $\tilde{\pi} = 0.9$ and $\tilde{\pi} = 0.1$ using a decreasing number Principal Components

Model	$\tilde{\pi} = 0.5$				$\tilde{\pi} = 0.9$				$\tilde{\pi} = 0.1$			
	ER	mDCF	DCF	MLoss	ER	mDCF	DCF	MLoss	ER	mDCF	DCF	MLoss
MVG	7.00%	0.1302	0.1399	0.0098	13.95%	0.3423	0.4001	0.0577	13.75%	0.2629	0.3051	0.0422
MVG[PC6]	7.00%	0.1302	0.1399	0.0098	13.95%	0.3423	0.4001	0.0577	13.75%	0.2629	0.3051	0.0422
MVG[PC5]	7.10%	0.1331	0.1419	0.0088	13.85%	0.3512	0.3980	0.0468	13.70%	0.2738	0.3041	0.0303
MVG[PC4]	8.05%	0.1537	0.1609	0.0072	14.95%	0.4150	0.4598	0.0448	15.75%	0.3012	0.3529	0.0517
MVG[PC3]	8.80%	0.1734	0.1759	0.0025	15.75%	0.4392	0.4680	0.0288	16.70%	0.3563	0.3879	0.0316
MVG[PC2]	8.80%	0.1731	0.1760	0.0029	16.10%	0.4384	0.4434	0.0050	16.70%	0.3526	0.3879	0.0353
MVG[PC1]	9.25%	0.1769	0.1850	0.0081	17.40%	0.4342	0.4775	0.0433	16.75%	0.3686	0.3970	0.0283
NBG	7.20%	0.1311	0.1439	0.0128	14.20%	0.3510	0.3893	0.0383	13.60%	0.2570	0.3022	0.0452
NBG[PC6]	8.90%	0.1727	0.1780	0.0053	15.70%	0.4359	0.4512	0.0153	16.50%	0.3535	0.3920	0.0385
NBG[PC5]	8.75%	0.1737	0.1750	0.0013	15.65%	0.4340	0.4660	0.0320	16.55%	0.3545	0.3930	0.0385
NBG[PC4]	8.85%	0.1717	0.1770	0.0053	15.50%	0.4313	0.4630	0.0317	16.75%	0.3614	0.3970	0.0356
NBG[PC3]	9.00%	0.1746	0.1799	0.0053	15.70%	0.4343	0.4591	0.0248	16.65%	0.3645	0.3950	0.0305
NBG[PC2]	8.85%	0.1710	0.1770	0.0060	16.05%	0.4323	0.4424	0.0101	16.65%	0.3562	0.3869	0.0307
NBG[PC1]	9.25%	0.1769	0.1850	0.0081	17.40%	0.4342	0.4775	0.0433	16.75%	0.3686	0.3970	0.0283
TCG	9.30%	0.1812	0.1860	0.0048	17.05%	0.4421	0.4626	0.0204	16.80%	0.3628	0.4061	0.0432
TCG[PC6]	9.30%	0.1812	0.1860	0.0048	17.05%	0.4421	0.4626	0.0204	16.80%	0.3628	0.4061	0.0432
TCG[PC5]	9.30%	0.1812	0.1860	0.0049	17.05%	0.4451	0.4626	0.0174	16.75%	0.3648	0.4051	0.0402
TCG[PC4]	9.25%	0.1821	0.1850	0.0029	17.00%	0.4441	0.4615	0.0174	16.65%	0.3610	0.4031	0.0421
TCG[PC3]	9.25%	0.1830	0.1850	0.0020	16.75%	0.4342	0.4565	0.0223	16.50%	0.3681	0.4082	0.0401
TCG[PC2]	9.25%	0.1789	0.1850	0.0062	17.45%	0.4352	0.4785	0.0433	16.70%	0.3630	0.3960	0.0330
TCG[PC1]	9.35%	0.1769	0.1870	0.0101	17.55%	0.4342	0.4806	0.0464	16.60%	0.3686	0.4021	0.0334
NTG	9.30%	0.1803	0.1860	0.0057	17.50%	0.4434	0.4716	0.0282	16.90%	0.3631	0.4080	0.0449
NTG[PC6]	9.20%	0.1800	0.1840	0.0041	17.40%	0.4300	0.4696	0.0396	16.75%	0.3611	0.4131	0.0520
NTG[PC5]	9.20%	0.1800	0.1840	0.0040	17.40%	0.4300	0.4696	0.0396	16.65%	0.3611	0.4112	0.0500
NTG[PC4]	9.15%	0.1789	0.1830	0.0041	17.40%	0.4361	0.4696	0.0336	16.70%	0.3641	0.4122	0.0480
NTG[PC3]	9.10%	0.1790	0.1820	0.0031	17.45%	0.4361	0.4706	0.0346	16.65%	0.3610	0.4112	0.0502
NTG[PC2]	9.35%	0.1769	0.1870	0.0101	17.50%	0.4312	0.4796	0.0484	16.45%	0.3590	0.3910	0.0320
NTG[PC1]	9.35%	0.1769	0.1870	0.0101	17.55%	0.4342	0.4806	0.0464	16.60%	0.3686	0.4021	0.0334

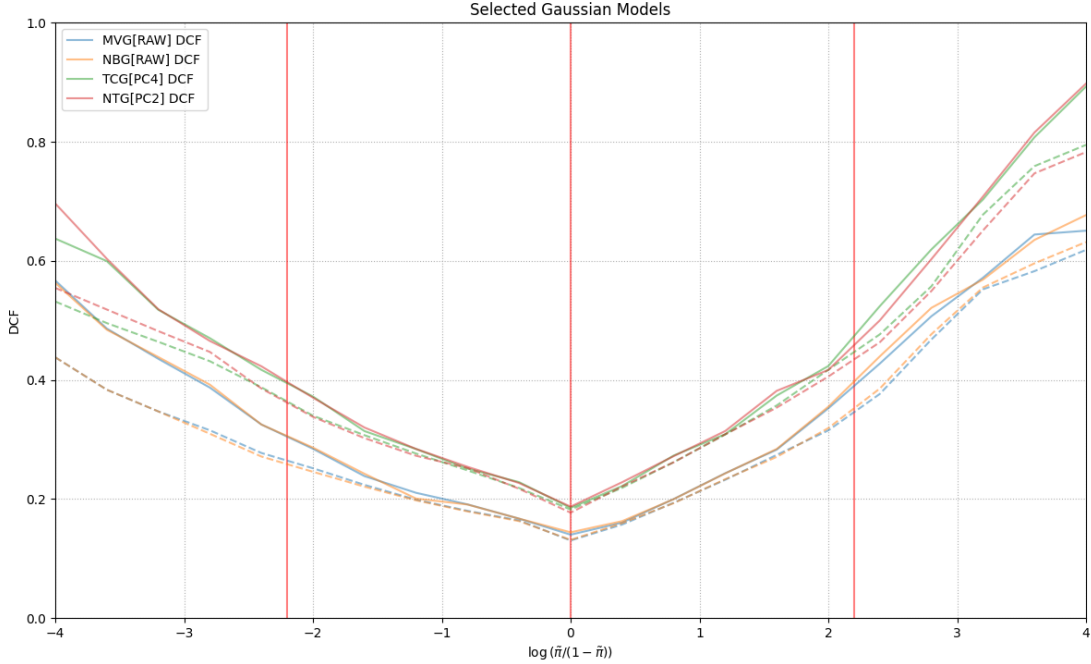


Figure 6: Bayes error plots of selected Gaussian models to compare

From Figure 6 we can notice that the minDCF for the Multivariate Gaussian and Naive Bayes models remain relatively consistent and comparable across different applications, indicating robust performance using different working points. The rankings of these models were also consistent, with these two models, with NB on top, outperforming Tied and Tied Naive Gaussian models in terms of minDCF. All models exhibited good calibration, as evidenced by the low miscalibration loss values. However, Tied Covariance Naive Bayes model had the lowest miscalibration loss, which indicates that it aligns well with the given application’s effective prior, even though its overall performance, evaluated through minDCF, was not the best.

3.2 Logistic Regression

3.2.1 Early analysis

Logistic Regression Model is a discriminative model that aims to discover a linear separation surface used to classify our data. The model is discriminative since, in the learning process, the found probability distribution is the posterior one, specializing the model in discrimination of data rather than generation. This model, in its original formulation, has a fundamental problem: when data are linearly separable, the minimization problem that the model aims to solve is not defined. A regularization term must be introduced: this solve the above issue, but on the other hand make the gathered scores lose their probabilistic interpretation. This loss is more consistent when dealing with a lot of samples: this happens because the model is already robust to overfitting with ample data, making regularization counterproductive. However, with small datasets, regularization mitigates overfitting, resulting in better performance.

We can check this property by plotting the DCF and minDCF for the whole dataset and just a slice of it (50 samples): we can see from Figure 7 that, using just few samples, make the DCF increase for higher lambda values and minDCF get lower when the regularization term lambda increase. So, since we have to use all the dataset, lambda value has to be chosen carefully taking into account this property.

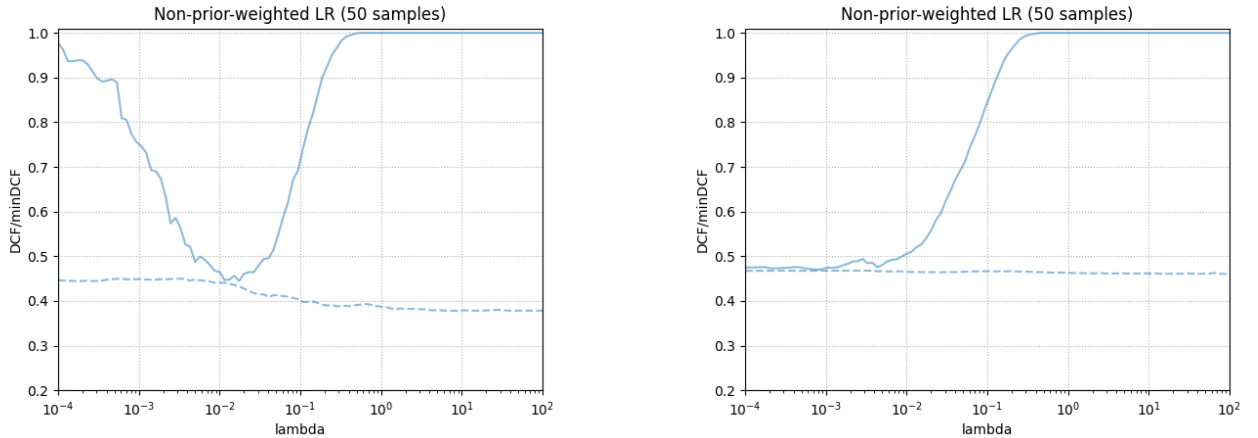


Figure 7: Comparison between DCF degradation for large and small datasets

3.2.2 Quantitative evaluation

Let’s now plot various model configuration using different values of the regularization coefficient λ .

Using non-prior-weighted and prior-weighted Logistic Regression model, we can see that, in term of minDCF, the results are comparable, with the non-prior-weighted being slightly better in term of DCF, a problem that can be solved by calibration in further analysis.

Since Logistic Regression aims to find a linear separation surface, it can be not enough to represent the optimal boundaries to use for classification. To use Logistic Regression building up non-linear decision boundaries in original features space, we can exploit Quadratic Feature Mapping. Using this configuration, we reach the lowest minDCF values overall, and even in actDCF it shows the best results compared to the other models: this is due to model’s ability to effectively capture non-linear relationships between features, maintaining a balance between complexity and generalization.

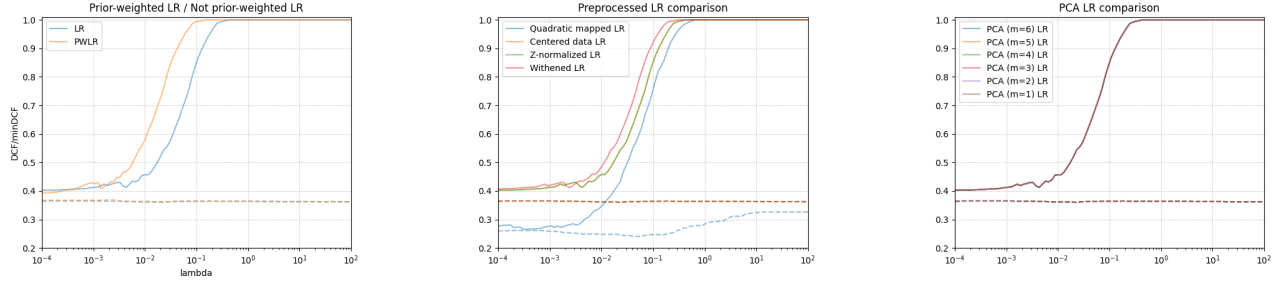


Figure 8: Comparison between minDCF and DCF of considered Logistic Regression configuration

Since the data are already mostly normalized, applying preprocessing techniques such as centering, Z-normalization and withening isn't beneficial, so they're not useful. DCF values seems to be slightly better, but for the minDCF the values are absolutely comparable with the non-preprocessed Logistic Regression.

The same reasoning can be applied to PCA.

Table 6: Performance of chosen model

Model	λ	Error Rate	minDCF	DCF
LogReg [QFM]	0.0132194115	6.20%	0.2969	0.3729

Finally, selecting the model based on performance achieved in minDCF, the final decision is for the non-prior-weighted Logistic Regression with Quadratic Feature mapping (its performance are reported in Table 6).

3.3 Support Vector Machine

3.3.1 Early analysis

Support Vector Machine is another discriminative machine learning model that, in its base form, aims to find a separation hyperplane that maximize the distance between the hyperplane itself and the point that is located the nearest to it (i.e. maximize the margin). It provide a way to interpret the regularization term of the Logistic Regression and can achieve good performance, expecially using non-linear kernels to train and predict: these methods offer a way to find non-linear hyperplane in the original space, operating in high dimensional, implicit feature space without ever computing the coordinates in that space, exploiting the fact that the objective in the dual form relies just on the inner product of the considered vectors. The biggest downside of SVM models is that the found scores are not probabilistically interpretable: we'll almost certainly need to calibrate the scores to obtain good generalization capabilities and good performance overall on evaluation set.

As a preliminary analysis, we can check if a non-linear surface will be determinant from the dataset characteristics. We know, from previous analysis, that the last two features are those that shows the most complex distribution overall. We take in analysis those features to check if, eventually, a quadratic surface would provide better separation of data points. In this context, we can use a degree 2 kernel that maps the last two features in a 1D vector, using the formula:

$$\mathbf{y}_i = \mathbf{x}_{i,[4:6]}$$

$$z_i = \mathbf{y}_{i,[0:1]}\mathbf{y}_{i,[1:2]}$$

We can now plot these one dimensional samples, taking into consideration that a linear rule correspond to classify samples based on a threshold, while a quadratic rule correspond to inequalities that involve the sample being inside an interval, and so on. As we can see from Figure 9, for the considered transformation, a simple threshold wouldn't be enough to correctly classify points, since, even if their distribution are clustered in two different points of the axis, the overlapping between the two classes is consistent. From this evaluation we can conclude that, if non-linear separability of the last two features isn't balanced by the others, we expect non-linear models to outperform linear ones.

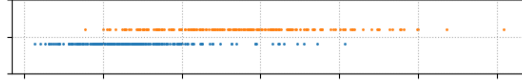


Figure 9: Transformed samples using last two features (1/10 samples plot)

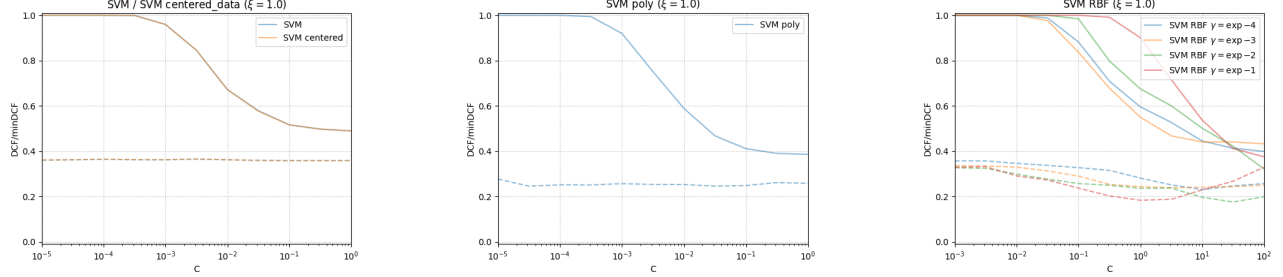


Figure 10: Comparison between minDCF and DCF of considered Support Vector Machine configuration

3.3.2 Quantitative evaluation

From the calibration point of view, mostly all trained models are comparable, with the only exception of Radial Basis Function kernel SVM, that shows a calibration loss that decrease proportionally with γ values. This is due to the excessive complexity that a high value of γ give to the model, suggesting that a lower value is beneficial to the generalization capabilities of the model itself (for relatively low value of regularization).

Running linear SVM, we can see that regularization term significantly affect performance in DCF and minDCF, but still the performance are better for non-linear model.

As we expected, trying with Polynomial kernel we can see that it achieve better results in term of minDCF and DCF than the linear one. We can notice even that the Polunomial kernel make the SVM performance stable across different regularization configuration. However, best performance overall are obtained using Radial Basis Function kernel, that shows better calibration with respect to the others and better performance in general. It means that the probability distribution of the original data is better captured by the shape given by the last kernel in the original space (Figure 11).

Table 7: Performance of chosen model

Model	K	C	Kernel	Dual loss	Dual gap	Error Rate	minDCF	DCF
SVM	1	31.62278	RBF ($\gamma = e^{-2}$)	1.209865e+04	4.825311e+05	4.50%	0.1755	0.4216

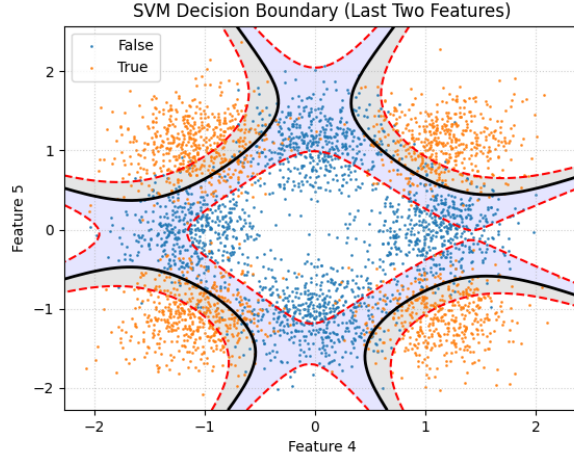


Figure 11: Decision boundaries of chosen model for the last two features

Finally, the selection process, even this time made by selecting the model achieving the lowest minDCF value, end with the chosen model to be the SVM using Radial Basis Function Kernel with the parameters showed in Table 7. A possible issue for this solution can be found in the value of the duality gap: since really high, this means that the optimization problem do not converge properly, due maybe to numerical issues or not optimal hyperparameter tuning. The model will still be selected as the final one, since, even if the problem do not converge properly, the performance achieved are still better than the other, and this is the selection criterion that has been decided to prioritize.

3.4 Gaussian Mixture Model

3.4.1 Early analysis

The Gaussian Mixture Models are generative model that assume that the underlying distribution from which data are generated is a Gaussian Mixture, i.e. a complex shape of a distribution can be approximated by the sum of a variable number of Gaussian distributions. The number of Gaussian distribution to use is an hyperparameter to tune to obtain a Mixture distribution that is the most similar to the real one.

3.4.2 Quantitative evaluation

From the results obtained by the different configuration of Gaussian Mixture Model, we can notice that those that works better are those that use more components for the True label and less for the False label. This result may indicate that the distribution of the True label is more complex and need to use more Gaussians to be well represented by a GMM. Conversely, using an excessive number of Gaussians for the False label may overfit to training data, reducing its generalization ability.

The Diagonal Covariance models outperform the others, suggesting that classification improves when using diagonal covariance matrices. This could be because full covariance matrices introduce too much flexibility, leading to overfitting, whereas the Diagonal Covariance assumption enhances generalization by reducing model complexity. Moreover, the Tied Covariance model performs well only when using 32 Gaussian components for the True label. This further supports the idea that the True label has a more complex distribution, requiring more components for proper representation. Additionally, the Tied Covariance constraint may be too restrictive, preventing the model from capturing the underlying distribution effectively when using fewer than 32 components.

Table 8: Performance of chosen models

Model	Error Rate	minDCF	DCF
GMM[DC] (nT=16,nF=8)	4.95%	0.1253	0.1548

Decision boundary of the selected model (Table 8) confirm what we saw in previous analysis: the last two features distribution are fitted the best using some complex distribution, as in this case with Gaussian Mixture Model.

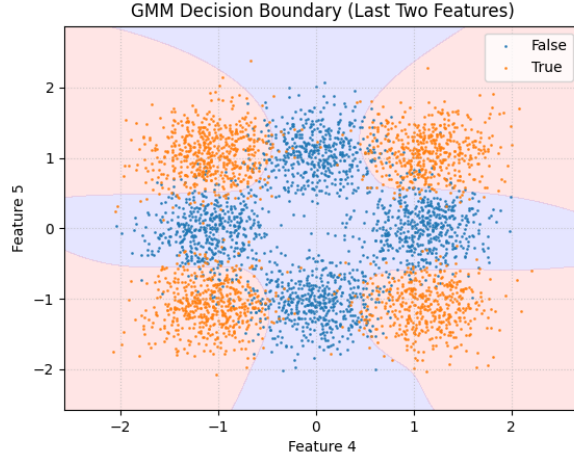


Figure 12: Diagonal Covariance Gaussian Mixture model ($nT=16$, $nF=8$) decision boundaries

4 Final Selection

Table 9: Performance of chosen model

Model	Error Rate	minDCF	DCF
GMM[DC] ($nT=16, nF=8$)	4.95%	0.1253	0.1548
GM [NB-RAW]	13.60%	0.2570	0.3022
LR [QFM] ($\lambda = 0.0132$)	6.20%	0.2969	0.3729
SVM[RBF] ($K = 1, C = 31.62, \gamma = e^{-2}$)	4.50%	0.1755	0.4216

For the application in exam, as we can see from Table 9, the model that perform best in term of minDCF is Diagonal Covariance Gaussian Mixture model, followed by Support Vector Machine using RBF kernel, Naive Bayes Gaussian Model and finally Logistic Regression using Quadratic Feature Mapping. Let's now check how these selected models perform in a various range of applications.

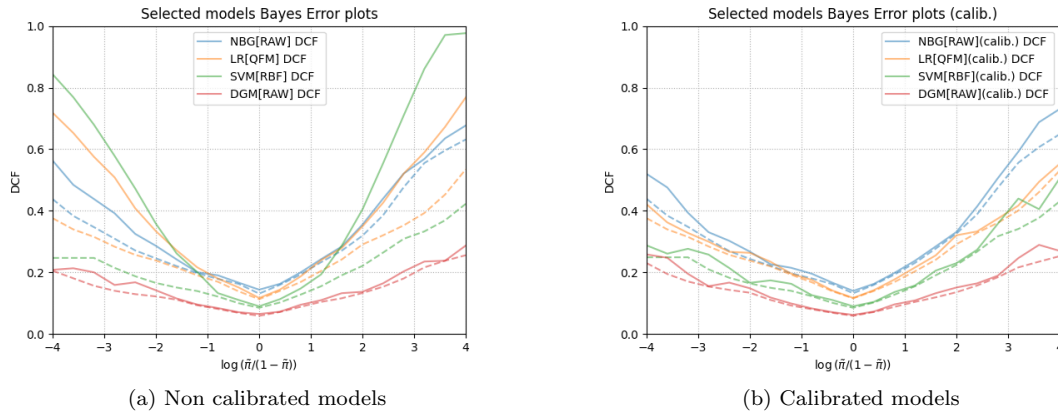


Figure 13

As we can see, the model showing the most stable performance across different applications is Gaussian Mixture model. While, as we could expect, the model that have scores that are not probabilistically interpretable have a worse calibration, such as Support Vector Machine and Logistic Regression.

Performing calibration on selected models, we obtain performance reported in Figure 13, in which we can notice that the models DCFs are now lower, with still GMM outperforming all the others models. Results obtained with fused scores are reported in Figure 14. The fusion doesn't provide any significant advantage with respect to

the performance achieved using Gaussian Mixture Model only, and so the final choice is for the Gaussian Mixture Model alone to be used.

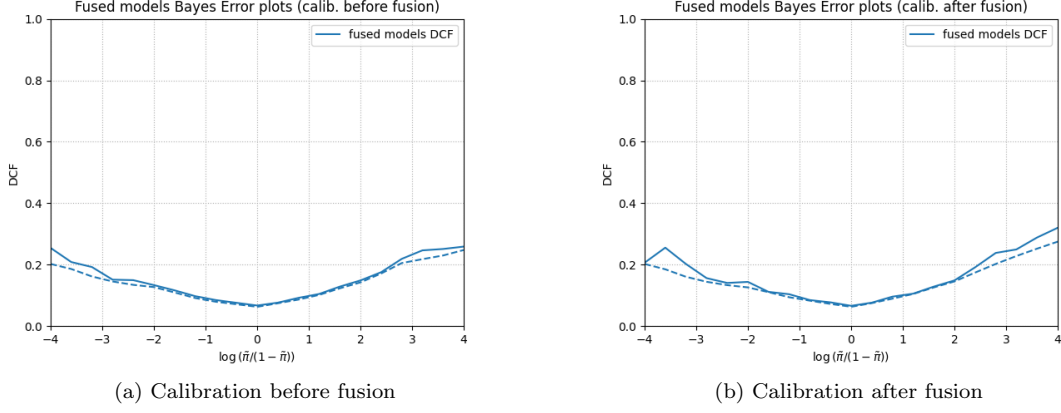


Figure 14

5 Delivered System

Finally, the evaluation of the model on the evaluation set is provided below.

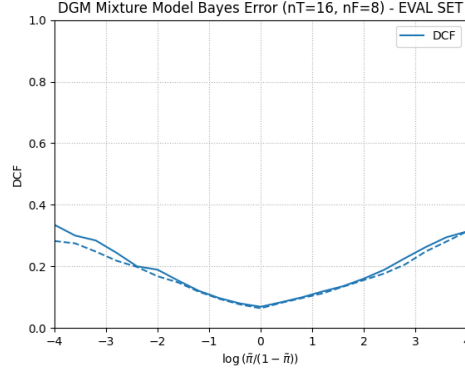


Figure 15: Delivered System Bayes Error plot

Table 10: Delivered System performance on evaluation set

Model	Error Rate	minDCF	DCF
GMM[DC] (nT=16,nF=8)	5.50%	0.1814	0.1927

The performance observed on the evaluation set is consistent with the results obtained during validation, confirming the reliability of the models even with data that have never been used to make choices about the model. This suggests that the chosen configurations generalize well and that no significant overfitting occurred.