# Pulsar Detection Problem

Mohamed Khaled Hassan Aly Motrash

S295805@studenti.polito.it

## Abstract: -

The goal of this paper is to analyze and study the pulsar Detection Data based on 8 different features which represent a signal .A potential signal detection known as a 'candidate'. As we have unbalanced data distribution between the classes, we will use some of the machine learning algorithms to analyze the data and learn the different relationships between the features and their distributions. After that we will analyze the different performances of different models to see the optimal one to use for that application.

## 1. Introduction:

Pulsars are rare type of Neutron star that produce radio emission while it rotates as it produces a pattern of broadband radio emission that is detectable on earth. Thus, pulsar search involves looking for periodic radio signals with large radio telescopes. The importance of analyzing this data is believed to have a link to the space-time, the inter-stellar medium, and states of matter.

Each pulsar produces slightly different emission pattern which varies through each rotation. in practice almost all detections are caused by radio frequency interference (RFI) and noise, making legitimate signals hard to find.

So, our goal is to be able to correctly classify the data taking into consideration the unbalanced distribution of data.

The dataset consists of different samples that are called 'Candidate' as each candidate(row) is described by 8 continuous variables and a single class variable (1 for Positive ,0 for Negative).

The first 4 features are simple statistics that describe a longitude-resolved version of the signal that has been averaged in both time and frequency.

The remaining 4 are obtained from the DM-SNR curve. A data mining stream lit application for astrophysical prediction using random forest classification.

Features interpretation:

1. Mean of the integrated profile.
2. Standard deviation of the integrated profile.
3. Excess kurtosis of the integrated profile.
4. Skewness of the integrated profile.
5. Mean of the DM-SNR curve.
6. Standard deviation of the DM-SNR curve.
7. Excess kurtosis of the DM-SNR curve.
8. Skewness of the DM-SNR curve.
9. Class.

# 2. Feature Analyze:

Each sample of the Pulsar detection data set consists of 8 features that represent different features grouped in the first 4 features that represent array of continuous variables obtained from the integrated pulse profile (folded profile). The last 4 features are obtained from the DM-SNR curve.

Each sample "Candidate" contains information about a pulsar radiation. These data set contains real pulsar examples embedded within RFI/ noise signals.
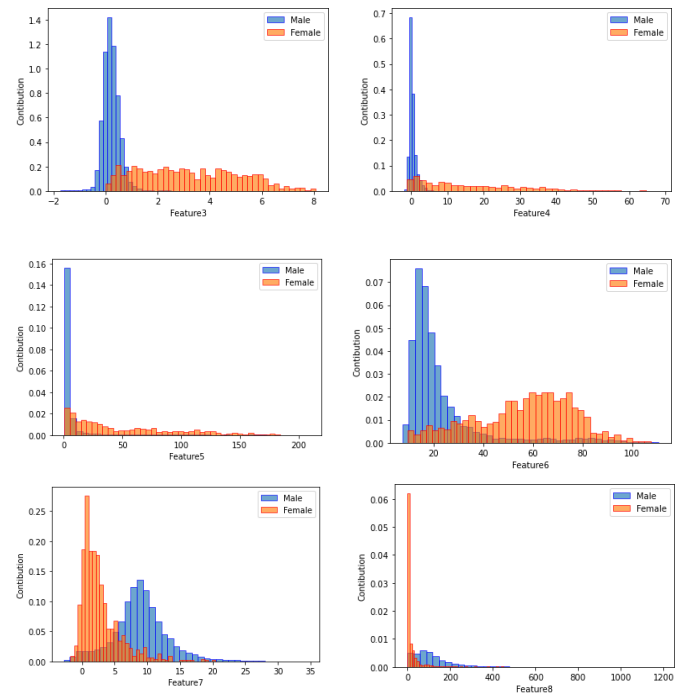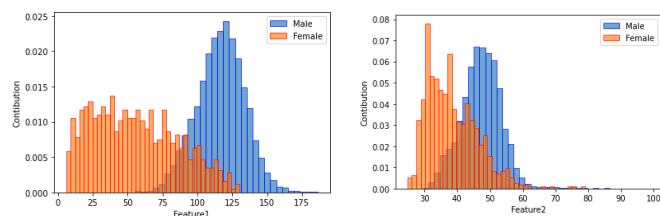
In our data set these data have been checked by human annotators represented by positive (class 1) in our data set

By analyzing the training Data, we find that we are dealing with a binary data with only class0 (Negative class) or class1(Positive) and we have an equivalent distribution of the testing samples such that we have:

- 8108 samples for class 0 (Negative)

- 821 samples for class 1 (Positive)

- 8969 samples for testing

On the following figure we could see the Raw feature distribution: -

Figure1: plots of Raw pulsar data set features Distribution
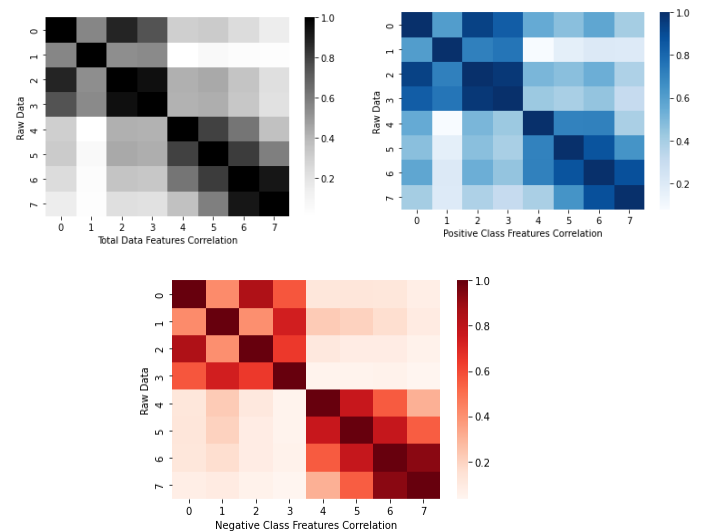




Now we will analyze the correlation between the features by plotting the heat map of the Pearson correlation coefficient according to the equation:

$$\frac{Cov(X,Y)}{\sqrt{Var(X)}\sqrt{Var(Y)}}$$

The heat map will display in a better way the correlation between different features and hence gives us a better point of view of our data.

Figure2: Heat maps of Pulsar detection for the RAW dataset features (Gray) & Positive dataset features (Blue) & Negative dataset features (Red)

From the plots shown in Figure 1, we can observe that features are not well distributed in Gaussian distribution with several outliers. And we can notice the high correlation among the first 4 features with each other and the last 4 features with each.

And that could tell us that a dimensionality reduction could be an effective choice if we decided to use it on number of dimensions of 4 but that will risk losing a lot of information from the other 4 features.

And that's what we will be analyzing in this paper with different classification techniques.

From the previous results it has become a fact that we must apply preprocessing techniques, such as Gaussianization of our dataset as we will work all our classification techniques and hypothesis based on the assumption that the data is distributed based on gaussian distribution.

# 3. Z-Normalization (Gaussian):

For x feature we will transform it by first computing the rank over the training dataset:

$$r(x) = \frac{\sum_{i=1}^{N} \mathbb{I}[x_i < x] + 1}{N + 2}$$

Then computing the transformed features using the inverse of the cumulative distribution function (percent point function, p.p.f)

$$y = \Phi^{-1}(r(x))$$

The following figure shows the feature distribution after applying Z-Normalization (Gaussianizing the Features): -
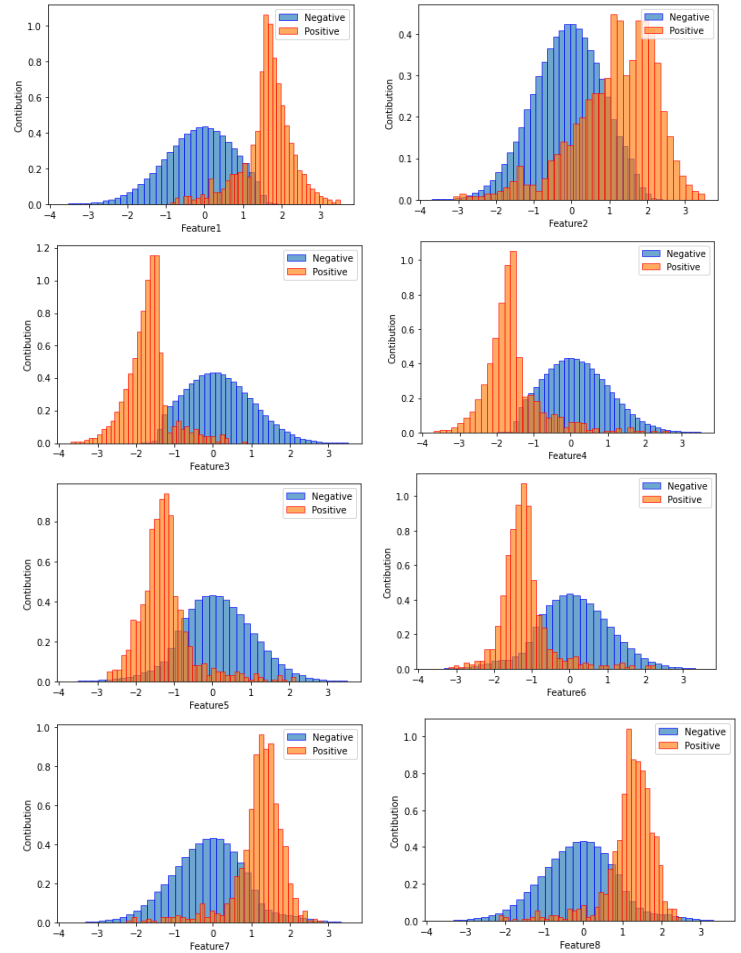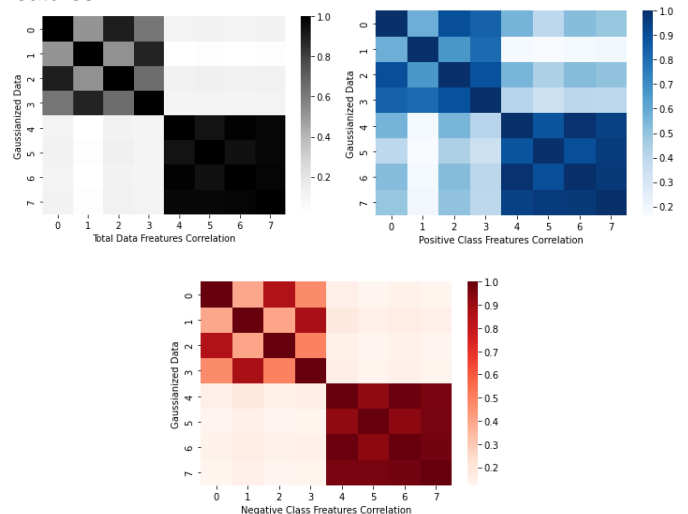


Figure3: plots of Gaussian (Z-Normalized ) pulsar data set features Distribution

From the previous figure the data is now well distributed, and the class Negative is better distributed giving the higher number of Candidates.

Now displaying the heat map again for the features: -

Figure4: Heat maps of Pulsar detection for the Z-Normalized features

From the previous figure it is now clear the correlation between the different features and there is a significant correlation between the first 4 features and the second 4 features in the positive class.

We could see a strong correlation coefficient between most of the features and that gives us an idea on the dimensionality reduction techniques if we will use any.

In such case we could see that possible using the original dataset without using any Dimensionality reduction techniques will allow the modules to learn more efficiently however, we could see that features 1 : 4 and 4 : 8 are highly correlated from other features so we could try a Dimensionality reduction technique PCA with m =4 but doing so will cause us to lose some important information so we will try using PCA with m =7 and with no PCA and compare the results.

# 4. Training methods:

## Single split:

In this technique we will use one single split of data with percentage of 80% for training and 20% for evaluation.

By doing so we will have the opportunity to estimate the model hyper-parameters only once and that will save us some time in the training phase and hence train more models.

## K-Fold cross validation:

Here we will have the chance to train the same model K times and evaluate the classifier.

And improving the hyper-parameters accordingly and this technique will give us better results than the previous technique but will consume more time in training each model.

We will use both techniques and compare the results. We expect better results using the K-fold approach.

We will use different applications with different priors since there is no information on the cost of different classes we will assume.

$$(\tilde{\pi}, Cfp, Cfn) = (0.5, 1, 1)$$

$$(\tilde{\pi}, Cfp, Cfn) = (0.1, 1, 1)$$

$$(\tilde{\pi}, Cfp, Cfn) = (0.9, 1, 1)$$

We will measure the performance of the model in terms of normalized minimum detection costs in all these applications.

# 5. Classifiers:

## 5.1. Gaussian Classifiers:

The first model of our analysis is Gaussian Classifiers: -

- MVG (Multivariate Gaussian Classifier)
- Naïve Bayes
- Tied Covariance
- Tied Naïve Bayes

Gaussian model is basically a model which depends on assign the class with the highest posterior probability C according to the following equation:

$$c_t^* = \arg\max_c P(C_t = c | X_t = x_t)$$

Gaussian Classifiers works under the assumption that data follows a gaussian distribution.

The calculation of the class posterior probability $P(Ct=c|Xt=xt)$ depends on the class-conditional distribution $X|C$

$$P(C_t = c | X_t = x_t) = \frac{f_{X,C}(x_t, c)}{\sum_{c' \in C} f_{X,C}(x_t, c')}$$

Applying the previous variants of the Gaussian Classifier on the data after applying the Z-Normalization with the different application priors we expect that for a data that is not equally distributed there would be an application with biased prior could give better results.

Table1: Gaussian Classifiers **without Z-Normalization** and without Dimensionality Reduction **NO PCA**

|  | Single split | | | 5-Folds | | |
| --- | --- | --- | --- | --- | --- | --- |
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| **Multivariate Gaussian** | 0.267 | 0.134 | 0.57 | 0.209 | 0.140 | 0.664 |
| **Naïve Bayes** | 0.392 | 0.163 | 0.618 | 0.272 | 0.157 | 0.366 |
| **Tied Covariance** | 0.274 | 0.184 | 0.664 | 0.204 | 0.144 | 0.638 |
| **TC & NB** | 0.359 | 0.232 | 0.644 | 0.269 | 0.138 | 0.428 |

Table2: Gaussian Classifiers **Z-Normalized** and without Dimensionality Reduction **NO PCA**

|  | Single split | | | 5-Folds | | |
| --- | --- | --- | --- | --- | --- | --- |
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| **Multivariate Gaussian** | 0.274 | 0.147 | 0.626 | 0.237 | 0.157 | 0.447 |
| **Naïve Bayes** | 0.433 | 0.129 | 0.667 | 0.365 | 0.142 | 0.480 |
| **Tied Covariance** | 0.264 | 0.126 | 0.589 | 0.209 | 0.124 | 0.495 |
| **TC & NB** | 0.480 | 0.145 | 0.665 | 0.407 | 0.148 | 0.428 |

Table3: Gaussian Classifiers **Z-Normalized** and with Dimensionality Reduction **PCA = 4**

|  | Single split | | | 5-Folds | | |
| --- | --- | --- | --- | --- | --- | --- |
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| **Multivariate Gaussian** | 0.253 | 0.133 | 0.604 | 0.230 | 0.140 | 0.488 |
| **Naïve Bayes** | 0.243 | 0.135 | 0.558 | 0.228 | 0.139 | 0.429 |
| **Tied Covariance** | 0.233 | 0.122 | 0.541 | 0.203 | 0.126 | 0.491 |
| **TC & NB** | 0.233 | 0.124 | 0.566 | 0.207 | 0.125 | 0.491 |

Table4: Gaussian Classifiers **Z-Normalized** and with Dimensionality Reduction **PCA = 7**

|  | Single split | | | 5-Folds | | |
| --- | --- | --- | --- | --- | --- | --- |
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| **Multivariate Gaussian** | 0.270 | 0.151 | 0.657 | 0.239 | 0.158 | 0.550 |
| **Naïve Bayes** | 0.261 | 0.141 | 0.587 | 0.262 | 0.155 | 0.479 |
| **Tied Covariance** | 0.247 | 0.127 | 0.601 | 0.215 | 0.125 | 0.507 |
| **TC & NB** | 0.250 | 0.121 | 0.562 | 0.230 | 0.132 | 0.482 |

By analyzing the previous results, we could reach to several points: -

- The application with prior 0.9 which increases the prior of the class Negative over the positive and given the fact that the biased dataset it is clear that this application is the worst as in almost all the results it gives the worst outcome.
- Unlike what is expected the application increases the prior of class Positive over the Negative (prior=0.1) gives slightly less results as compared to the application with prior = 0.5
- So we take conclusion that application with prior = 0.5 gives better results in nearly all of the applications.
- Without applying the Z-Normalization we couldn't get clear results as the data are considered sparse.
-  Comparing the different Gaussian classifiers we could see that in different prior applications [0.1 , 0.5 , 0.9] we could see a clear advantage to the **Tied Covariance Gaussian Classifier.**
- Applying the PCA =4 gave us the best results as there is higher correlation between 4 features and as expected it gave the best results despite the lost information it manages to outperform the PCA = 7 by a slight value or come closely

A Good Candidate to solve our problem is using **Tied Covariance Gaussian Classifier**

## 5.2. Logistic Regression:

We will now analyze the Logistic regression model we will perform the regularized model.

$$J(w,b) = \frac{\lambda}{2}\|w\|^2 + \frac{1}{n}\sum_{i=1}^{n}\log\left(1 + e^{-z_i(w^T x_i + b)}\right) \ , \quad z_i = \begin{cases} 1 & \text{if } c_i = 1 \\ -1 & \text{if } c_i = 0 \end{cases} \quad (\text{i.e. } z_i = 2c_i -$$

We will train our model to obtain our model parameters $(w, b)$. After which we will calculate the score of our data according to different values of $\lambda$.

$$s(x_t) = w^T x_t + b$$

To help selecting an optimal $\lambda$ we plot different values of $\lambda$ (regularization coefficient) against the minDCF and choose the value that gives us the min DCF and consider it for our application.

Large values of $\lambda$ gives us lower performances so we will choose a $\lambda$ which will be as close to 0

Figure 5: min DCF plotted against different values of regularization coefficient **without Z-Normalization** and without Dimensionality Reduction **NO PCA**
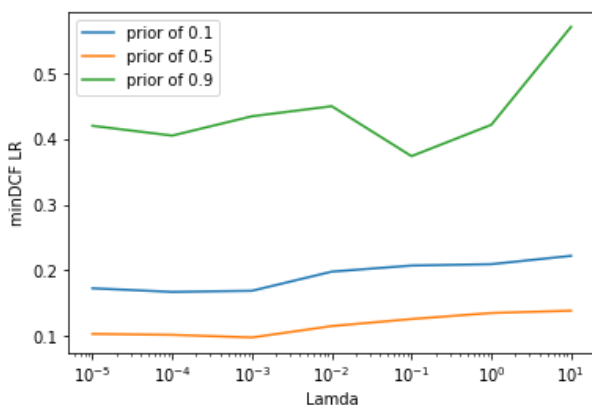
Figure 6: min DCF plotted against different values of regularization coefficient **Z-Normalization** and without Dimensionality Reduction **NO PCA**
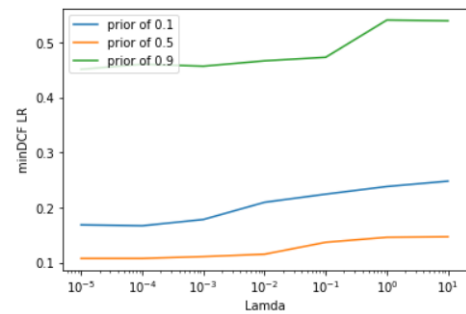
Figure 7: min DCF plotted against different values of regularization coefficient **Z-Normalization** and with Dimensionality Reduction **PCA =4**
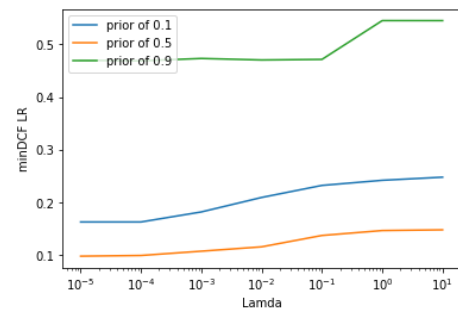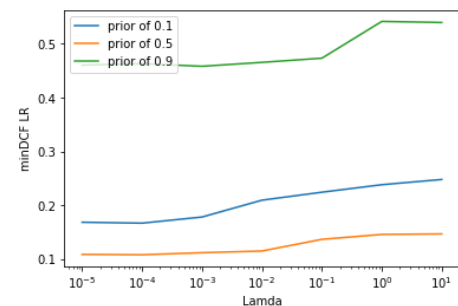
Figure 8: min DCF plotted against different values of regularization coefficient **Z-Normalization** and with Dimensionality Reduction **PCA =7**

As we can see from the plots a good value to be for Lambda $\lambda = 10^{-5}$ and $\lambda = 10^{-4}$

As we see the similarity between all the applications and the slight difference between using PCA and Raw data

The following table shows the numbers for results of the previous Plots: -

| | $\lambda = 10^{-4}$ | | | $\lambda = 10^{-5}$ | | |
|---|---|---|---|---|---|---|
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| LR (No PCA, No Z-Nom.) | 0.172 | 0.101 | 0.405 | 0.167 | 0.102 | 0.419 |
| LR (NO PCA) | 0.167 | 0.107 | 0.460 | 0.168 | 0.107 | 0.451 |
| LR (PCA = 4) | 0.163 | 0.099 | 0.468 | 0.163 | 0.098 | 0.469 |
| LR (PCA = 7) | 0.167 | 0.108 | 0.463 | 0.168 | 0.109 | 0.460 |

Also, here we could find that applying PCA =4 is performing much better than the other applications of PCA=7 and no PCA.

And another observation is that the application that uses prior =0.5 behaves the best and the application with prior=0.9 is the worst.

Overall, by choosing $\lambda = 10^{-5}$ we notice that it gives us a slightly better performance than that of $\lambda = 10^{-4}$ so we will be using $\lambda = 10^{-5}$.

So, our candidate is **LR ($\lambda = 10^{-5}$, ~π = 0.5).**

## 5.3. Linear Support Vector Machines:

Support Vector Machines are linear classifiers that look for maximum margin separation hyperplanes.

The (primal) SVM objective consists in minimizing: -

$$J(w,b) = \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{n} \max(0, 1 - z_i(w^T x_i + b))$$

By starting from the logistic regression, the choice of an optimal solution depends on the selection of a linear classifier which helps the best in separating the 2 classes with minimum errors as possible.

However, with these constraints we find an infinite number of hyper planes which could give us the same result.

So intuitively we can select a hyper place that separates the classes with the largest margin and that is the distance of the closest point with respect to the separation hyperplane.

We will obtain the separation hyperplane by solving the Dual problem.

Maximize w.r.t. $\alpha$

$$\alpha^T \mathbf{1} - \frac{1}{2}\alpha^T H \alpha$$

Where:

$$0 \leq \alpha_i \leq C, \; i = 1, \ldots, n$$

$$\sum_{i=1}^{n} \alpha_i z_i = 0$$

In order to select the optimal C value, we plot the minDCF graph for a balanced SVM to get the optimal value for the C.

Figure 9: min DCF plotted against different values of regularization coefficient **without Z-Normalization** and without Dimensionality Reduction **NO PCA**
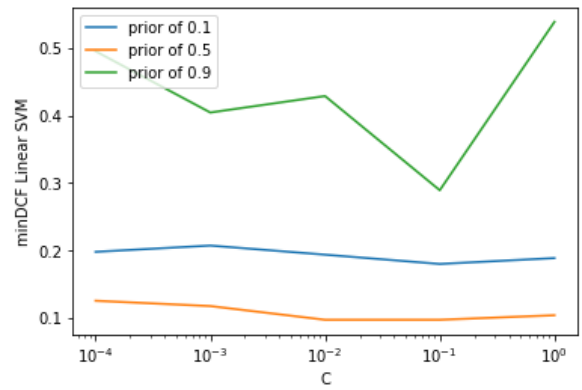


Figure 10: min DCF plotted against different values of regularization coefficient **Z-Normalization** and without Dimensionality Reduction **NO PCA**
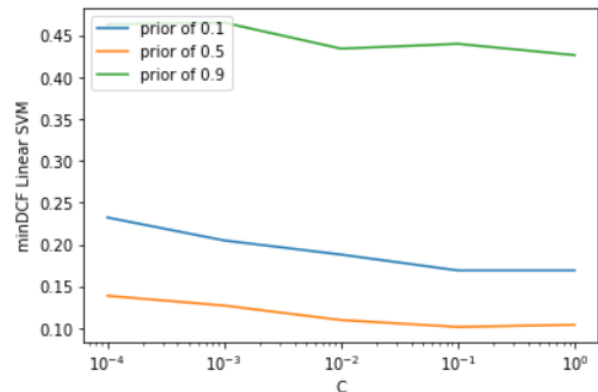
Figure 11: min DCF plotted against different values of regularization coefficient **Z-Normalization** and with Dimensionality Reduction **PCA =4**
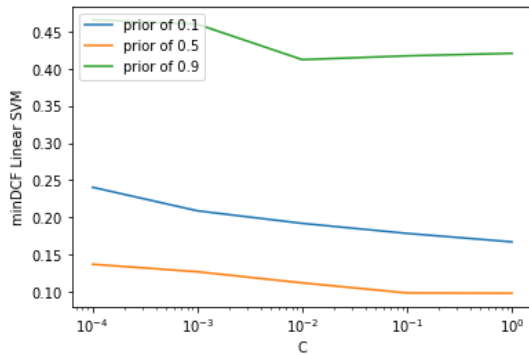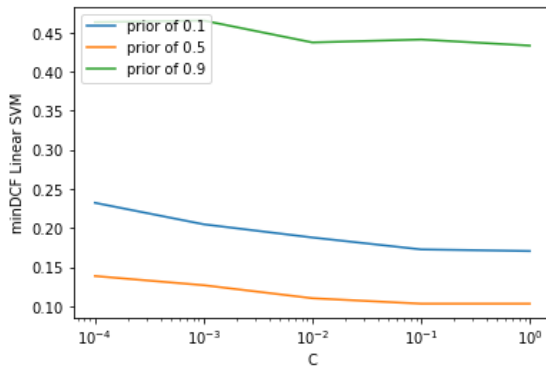


Figure 12: min DCF plotted against different values of regularization coefficient **Z-Normalization** and with Dimensionality Reduction **PCA =7**



From the graphs we could see some similarity in the results with best results for prior = 0.5 and also 2 close values of C parameter that is at C = $10^{-1}$ values and C = $10^0$

| | C = $10^{-1}$ | | | C = $10^0$ | | |
|---|---|---|---|---|---|---|
| **~π** | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| **LSVM (No PCA, No Z-Nom.)** | 0.180 | 0.097 | 0.289 | 0.188 | 0.097 | 0.539 |
| **LSVM (NO PCA)** | 0.169 | 0.101 | 0.440 | 0.169 | 0.104 | 0.426 |
| **LSVM (PCA = 4)** | 0.178 | 0.098 | 0.417 | 0.167 | 0.097 | 0.421 |
| **LSVM (PCA = 7)** | 0.173 | 0.103 | 0.441 | 0.171 | 0.103 | 0.433 |

From the previous results it is clear that choosing any of these values **C= 0.1,1** will give us a close result.

Also, by applying the Linear SVM it is also noticed the previous notations that is the application with 0.5 prior performs also better.

Linear SVM with PCA = 4 gives the best results with C =1

So, the Candidate here is **LSVM (PCA=4, C=1)**

## 5.4. Kernel Support Vector Machines:

SVMs allow for non-linear classification through an implicit expansion of the features in a higher dimensional space. The SVM dual objective depends on the training samples only through dot-products, and we can compute a classification score through scalar products between training and evaluation samples.

$$k(\boldsymbol{x}_1, \boldsymbol{x}_2) = \Phi(\boldsymbol{x}_1)^T \Phi(\boldsymbol{x}_2)$$

Where k is the kernel function as previously used.

$$\widehat{H}_{i,j} = z_i z_j k(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

Here we will try different kernels:

-   Polynomial kernel of degree d:

$$k(\boldsymbol{x}_1, \boldsymbol{x}_2) = (\boldsymbol{x}_1^T \boldsymbol{x}_2 + c)^d$$

-   Radial Basis Function kernel:

$$k(\boldsymbol{x}_1, \boldsymbol{x}_2) = e^{-\gamma \|\boldsymbol{x}_1 - \boldsymbol{x}_2\|^2}$$

The choice of the kernel and of its hyper-parameters (. c and $\gamma$) can also be made through cross validation.

The next Table shows the impact of different functions for the kernel support vector machine.

First RBF Kernel (K=0, C=1)

| | γ = 1 | | | γ = 10 | | |
|---|---|---|---|---|---|---|
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| KSVM (No PCA, No Z-Nom.) | 0.306 | 0.196 | 0.406 | 0.502 | 0.201 | 0.391 |
| KSVM (NO PCA) | 0.219 | 0.132 | 0.498 | 0.218 | 0.142 | 0.767 |
| KSVM (PCA = 4) | 0.236 | 0.138 | 0.517 | 0.208 | 0.128 | 0.672 |
| KSVM (PCA = 7) | 0.208 | 0.132 | 0.497 | 0.212 | 0.143 | 0.764 |

Second Poly Kernel (K=1, C=1, d=2)

| | c = 0 | | | c = 1 | | |
|---|---|---|---|---|---|---|
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| KSVM (No PCA, No Z-Nom.) | 1.0 | 0.935 | 1.0 | 1.0 | 0.934 | 1.0 |
| KSVM (NO PCA) | 1.0 | 0.99 | 0.999 | 1.0 | 0.99 | 0.998 |
| KSVM (PCA = 4) | 1.0 | 0.954 | 0.994 | 1.0 | 0.988 | 0.984 |
| KSVM (PCA = 7) | 1.0 | 0.993 | 0.995 | 1.0 | 0.998 | 0.999 |

After the previous results it is clear that using Poly kernel is not efficient at all .

However, by using RBF application we could see the best results is when using PCA = 4 and $\gamma$ = 10.

All results are poor but best candidate here: -

**KSVM_RBF ($\gamma$ = 10, PCA =4, K=0, C =1)**

## 5.5.  Gaussian Mixture Model:

GMM is a classifier that takes the normal Gaussian distributions  to a multiple distribution.

Gaussian mixture model could be considered basically as  just clustering technique that for a certain group it applies the gaussian model and models it.

The one class is divided to several components which are treated as summation of numbered Gaussian densities.

According to the eqution:

$$f_X(x) = \sum_{c=1}^{K} f_{X|C}(x|c)P(C=c) = \sum_{c=1}^{K} \pi_c \mathcal{N}(x|\mu_c, \Sigma_c)$$

Such that for each component is measured by means of weighted combination of Gaussians:-

$$f_X(x) = \sum_{c=1}^{K} w_c \mathcal{N}(x; \mu_c, \Sigma_c)$$

With parameters as follows:

$$M = [\mu_1 \dots \mu_K] \quad , \quad S = [\Sigma_1 \dots \Sigma_K]$$

$$w = [w_1 \dots w_K]$$

In order to estimate our initial values to start maximizing the likelihood we could use these techniques: -

- EM algorithm
- LBG algorithm that works as initial guess for the EM algorithm and works as duplicator for the components.

We could also try different covariance matrices: -

- Normal Covariance Matrix
- Diagonal Covariance Matrix
- Tied Covariance Matrix

By applying each of these combinations with the  applications with different priors [0.1,0.5,0.9] we could measure and spot the best performance

The following Table summarizes the different methods with **8 Components** for all: -

Table4: GMM Classifiers **No Normalization** and without= Dimensionality Reduction **No PCA**

|  | Single Fold | | | 5 - Fold | | |
|---|---|---|---|---|---|---|
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| **EM GMM** | 0.250 | 0.126 | 0.649 | 0.219 | 0.114 | 0.468 |
| **EM Diagonal GMM** | 0.257 | 0.145 | 0.536 | 0.202 | 0.141 | 0.539 |
| **EM Tied GMM** | 0.230 | 0.171 | 0.847 | 0.191 | 0.156 | 0.721 |

Table4: GMM Classifiers **Z-Normalized** and without Dimensionality Reduction **No PCA**

|  | Single Fold | | | 5 - Fold | | |
|---|---|---|---|---|---|---|
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| **EM GMM** | 0.219 | 0.111 | 0.493 | 0.196 | 0.122 | 0.451 |
| **EM Diagonal GMM** | 0.260 | 0.116 | 0.539 | 0.241 | 0.135 | 0.595 |
| **EM Tied GMM** | 0.449 | 0.243 | 0.753 | 0.444 | 0.239 | 0.628 |

Table4: GMM Classifiers **Z-Normalized** and without Dimensionality Reduction **PCA =4**

|  | Single Fold | | | 5 - Fold | | |
|---|---|---|---|---|---|---|
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| **EM GMM** | 0.206 | 0.120 | 0.535 | 0.191 | 0.102 | 0.133 |
| **EM Diagonal GMM** | 0.239 | 0.109 | 0.506 | 0.191 | 0.109 | 0.498 |
| **EM Tied GMM** | 0.365 | 0.227 | 0.682 | 0.365 | 0.208 | 0.619 |

Table4: GMM Classifiers **Z-Normalized** and without Dimensionality Reduction **PCA =7**

|  | Single Fold | | | 5 - Fold | | |
|---|---|---|---|---|---|---|
| ~π | 0.1 | 0.5 | 0.9 | 0.1 | 0.5 | 0.9 |
| **EM GMM** | 0.196 | 0.111 | 0.580 | 0.193 | 0.101 | 0.403 |
| **EM Diagonal GMM** | 0.209 | 0.108 | 0.568 | 0.197 | 0.120 | 0.420 |
| **EM Tied GMM** | 0.432 | 0.238 | 0.723 | 0.400 | 0.232 | 0.638 |

As seen from the tables it is clear that by using the original covariance matrix the results are much better than the diagonal and tied Covariances.

And it is also appear a close similarity between using the PCA =4 and PCA =7 in some applications however in our model we will use the PCA =4 as it gives s slightly better results.

Given the fact that the application with 0.5 prior is still gives us a better result.

So, the best candidate that we could consider is **GMM (PCA = 4) EM Full Covariance GMM**
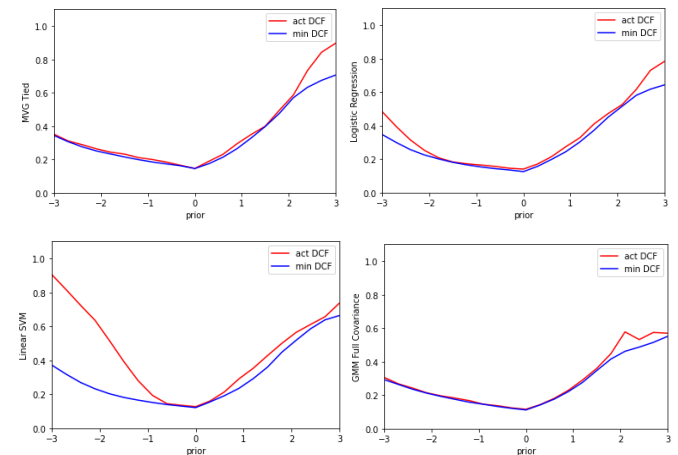
# 6. Scores Calibration:

Now we move our attention to the calibration of the results after we determined several applications that we will consider as a solution to our problem.

- **MVG Tied Covariance (PCA =4, ~π = 0.5)**
- **Logistic Regression ($\lambda = 10^{-5}$, ~π = 0.5, PCA=4).**
- **Linear SVM (PCA=4, C=1)**
- **GMM (PCA = 4) EM Full Covariance GMM**

We then try all of these algorithms on the hole data set to see their performance for **real DCF** and the **min DCF.**

Figure13: All algorithms performance with actual DCF and min DCF



Results are different in terms of min DCF and act DCF and for that we will try to make score calibration technique that will be computing the transformation function and this will help mapping the actual DCF to be as close to the min DCF as possible as this will be later

applied to the Testing samples to get the best possible results.

By assuming the f function to be in the form:

$$f(s) = \alpha s + \beta$$

We will need to estimate the values of 'a and 'b'. and f here could be interpreted as the log-likelihood ration for the 2 class hypotheses.
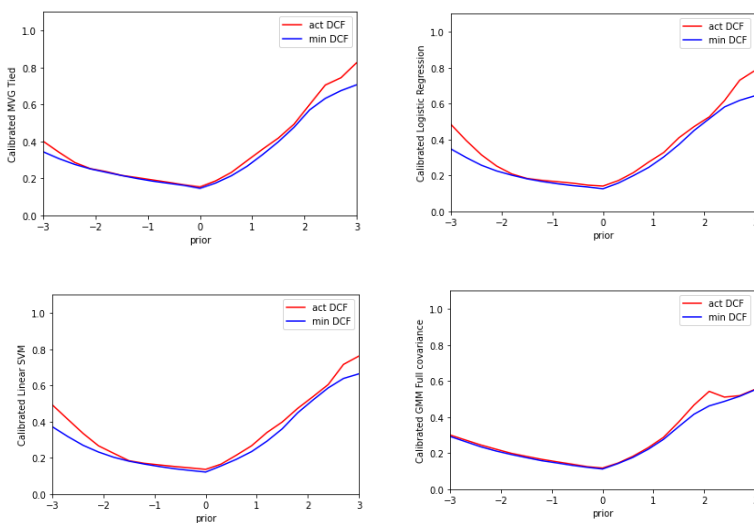
$$f(s) = \log \frac{f_{S|C}(s|\mathcal{H}_T)}{f_{S|C}(s|\mathcal{H}_F)} = \alpha s + \beta$$

Where the class posterior probability for prior is:

$$\log \frac{P(C = \mathcal{H}_T|s)}{P(C = \mathcal{H}_F|s)} = \alpha s + \beta + \log \frac{\tilde{\pi}}{1 - \tilde{\pi}}$$

And by applying these techniques to our data we could get the results after the Calibration as follows: -

Figure14: All **Calibrated** algorithms performance with actual DCF and min DCF



As we see from the plots the actual DCF is now becoming more like that of the min DCF and after doing that the modules are now ready to be tested on the real data.

We managed to successfully map the miss calibrated scores to become better calibrated.

# 7. **Testing Data Evaluation:**

Now we apply all of the 4 modules to the testing data (unseen data) and we then see the performance on reality.

In the previous modules we have seen that prior = 0.5 gave us the best performance overall so we will apply it for all the modules.

The following table has the overall performance comparing the min DCF with the act DCF: -

Table4: GMM Classifiers **Z-Normalized** and without Dimensionality Reduction **PCA =4**

|  | minDCF | ActualDCF |
|---|---|---|
| **~π** | **0.5** | **0.5** |
| **MVG Tied Full Cov** | 0.143 | 0.150 |
| **LR (λ = 10⁻⁵)** | 0.120 | 0.133 |
| **L SVM (C =1)** | 0.115 | 0.130 |
| **GMM EM Full Covariance GMM** | 0.107 | 0.111 |

The results are consistent with the ones obtained during the previous analysis over the training set and it is clear that the model performs very good and similar to what we have expected and that confirms our choices.

We can also see that Gaussian mixture model with the Full Covariance has managed to get the best performance on the unseen data.

## 8. Conclusions:

During the analysis we discovered that for the pulsar dataset linear classification is better that the quadratic ones. Also, it is clear that our model managed to get close enough to the testing data .and that makes our choices correct in selecting these modules.