

---

# REPORT MACHINE LEARNING AND PATTERN RECOGNITION

---

## Fingerprint spoofing detection

Alessia Intini  
s309895  
Politecnico di Torino  
2023/24

**MLPR Project:  
Fingerprint spoofing detection**

## Contents

<b>1</b>	<b>Dataset Analysis</b>	<b>3</b>
1.1	Training and evaluation sets . . . . .	3
1.2	Features analysis . . . . .	3
<b>2</b>	<b>Dimensionality reduction</b>	<b>4</b>
2.1	PCA . . . . .	5
2.2	LDA . . . . .	6
<b>3</b>	<b>Multivariate Gaussian Density</b>	<b>7</b>
<b>4</b>	<b>Classification model analysis</b>	<b>7</b>
4.1	Gaussian Models . . . . .	7
4.1.1	MVG Gaussian Classifier . . . . .	7
4.1.2	Naive Bayes Classifier . . . . .	7
4.1.3	Tied Bayes Classifier . . . . .	7
4.1.4	Gaussian Models Comparison . . . . .	7
4.2	Logistic Regression Classifier . . . . .	7
4.2.1	Quadratic Logistic Regression(QLR) . . . . .	7
4.3	SVM Classifier . . . . .	7
4.3.1	Linear SVM . . . . .	7
4.3.2	Kernel SVM . . . . .	7
4.4	GMM Classifier . . . . .	7
<b>5</b>	<b>Score Calibration</b>	<b>7</b>
5.1	Calibration Analysis on Selected Models . . . . .	7
5.2	Calibrating Scores for Selected Models . . . . .	7
<b>6</b>	<b>Experimental Results</b>	<b>7</b>
6.1	Calibration on evaluation score . . . . .	7
6.2	Considerations . . . . .	7
<b>7</b>	<b>Conclusions</b>	<b>7</b>

# MLPR Project: Fingerprint spoofing detection

The goal of project is to perform a binary classification on fingerprint spoofing detection, that is to identify genuine vs counterfeit fingerprint images. The dataset consists of labeled samples corresponding to the genuine (True, label 1) class and the fake (False, label 0) class, the data is 6-dimensional.

## 1 Dataset Analysis

### 1.1 Training and evaluation sets

The datasets provided contain 6000 samples; the first 6 values in each row represent the features, while the last value is the label. Specifically they are:

- Training Set: 2990 samples belonging to the Fake class (label 0) and 3010 samples belonging to the Genuine class (label 1)
- Evaluation Set: 3010 samples belonging to the Fake class (label 0) and 2990 samples belonging to the Genuine class (label 1)

We will use the Training Set to perform all the analysis, during this phase the dataset is divided to use about 60% of it as the training dataset and the remaining 40% for validation. After this step, the most promising models were chosen and the evaluation dataset was used to evaluate them and make the final considerations.

### 1.2 Features analysis

We can start to analyze all features, here is the graph of each feature, through these histograms it is shown how each feature is distributed:

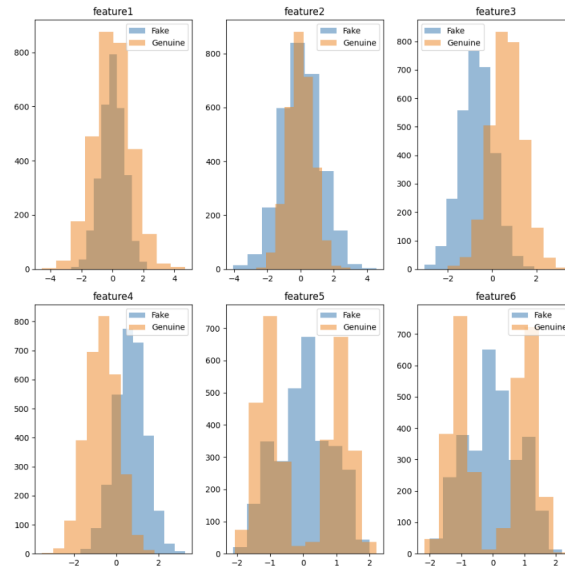


Figure 1: All features

By plotting a histogram for each feature, separately for the Fake and Genuine classes, it is possible

## MLPR Project: Fingerprint spoofing detection

to show whether the samples follow a Gaussian distribution and to what extent. Thus, we can say that some features follow a Gaussian distribution. Furthermore, it can be seen that features 4 and 5 may be the most discriminating features based on their ability to separate the data of the two classes.

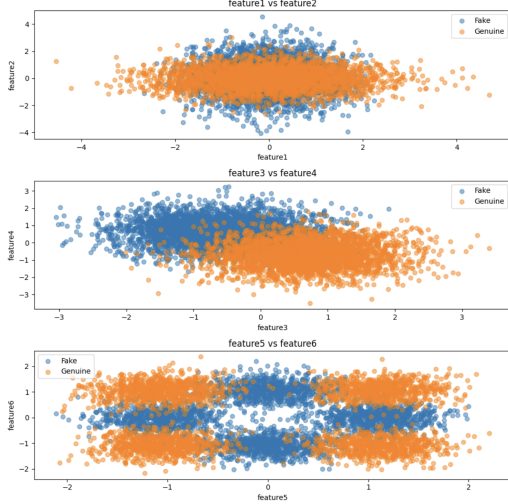


Figure 2: Distribution of feature pairs

Feature	Media	Varianza
1	0.00170711	1.00134304
2	0.00503903	0.9983527
3	-0.00560753	1.0024818
4	0.00109537	0.99029389
5	-0.00700025	1.00119747
6	0.00910515	0.99722374

Figure 3: Mean and variance for each features

After making these calculations we can conclude that the data are already centered since all features have mean close to zero.

Analyzing the pairwise averages and variances of the features, we can make some observations:

- Feature 1 and Feature 2: These two features have averages very close to zero, indicating that their values are similarly distributed around zero. However, the variance of Feature 2 is slightly lower than that of Feature 1, indicating that the values of Feature 2 are slightly less dispersed than those of Feature 1.
- Feature 3 and Feature 4: Feature 3 has a negative mean, while Feature 4 has a positive mean. This might indicate that the values of Feature 3 tend to be lower than those of Feature 4. Also, the variance of Feature 3 is slightly higher than that of Feature 4, indicating that the values of Feature 3 are more dispersed than those of Feature 4.
- Feature 5 and Feature 6: Feature 5 has a negative mean, while Feature 6 has a positive mean. This might indicate that the values of Feature 5 tend to be lower than those of Feature 6. Also, the variance of Feature 5 is slightly higher than that of Feature 6, indicating that the values of Feature 5 are more dispersed than those of Feature 6.

## 2 Dimensionality reduction

Dimensionality reduction is useful in this context to compute a mapping from an n-dimensional feature space to an m-dimensional space; it is applied because it can compress information, remove noise and simplify classification.

## MLPR Project: Fingerprint spoofing detection

### 2.1 PCA

PCA is an unsupervised dimensionality reduction technique that, given a centered dataset  $X = \{x_1, \dots, x_k\}$ , it aims to find the subspace of  $R^n$  that allows to preserve most of the information, i.e. the directions with the highest variance. We start by calculating the covariance matrix:

$$\mathbf{C} = \frac{1}{K} \sum_i (x_i - \bar{x})(x_i - \bar{x})^T \quad (1)$$

After that we can compute the eigen-decomposition of  $\mathbf{C} = \mathbf{U}\mathbf{\Sigma}\mathbf{U}^T$ , where  $\mathbf{U}$  is the matrix of eigenvectors and  $\mathbf{\Sigma}$  is the diagonal matrix of eigenvalues. Now we can project the data into the new subspace, it is spanned by the  $m$  columns of  $\mathbf{U}$  corresponding to the  $m$  values of the highest eigenvalues.

$$\mathbf{y}_i = \mathbf{P}^T(x_i - \bar{x}) \quad (2)$$

$\mathbf{P}$  is the matrix corresponding to the  $m$  columns of  $\mathbf{U}$  associated with the  $m$  highest eigenvalues of  $\mathbf{C}$ . At this point to choose which is the best value of  $m$  we need to check how much total variance in the data we are able to retain using different values of  $m$ . To do this evaluation, it was decided to represent on a graph the variance of the data as  $m$  increases, i.e., we exploited that each eigenvalue corresponds to the variance along the corresponding axis, and the eigenvalues are the elements of the diagonal of the matrix  $\mathbf{\Sigma}$ . The percentage will be calculated as the ratio of the sum of the first  $m$  eigenvalues to the sum of all eigenvalues.

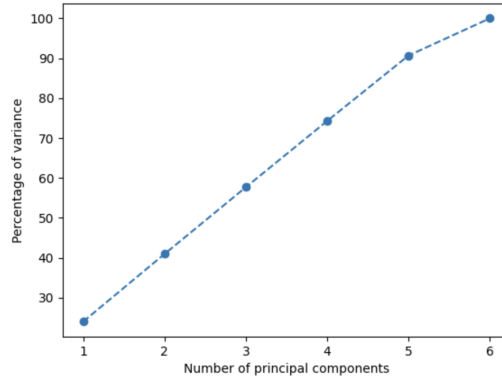


Figure 4: Variance of data for each  $m$  values

It is clear from Figure 4 that values of  $m < 5$  rapidly decrease the amount of variance retained in the data. Removing a feature is probably the best choice, since this way the percentage of variance remains fairly high.

## MLPR Project: Fingerprint spoofing detection

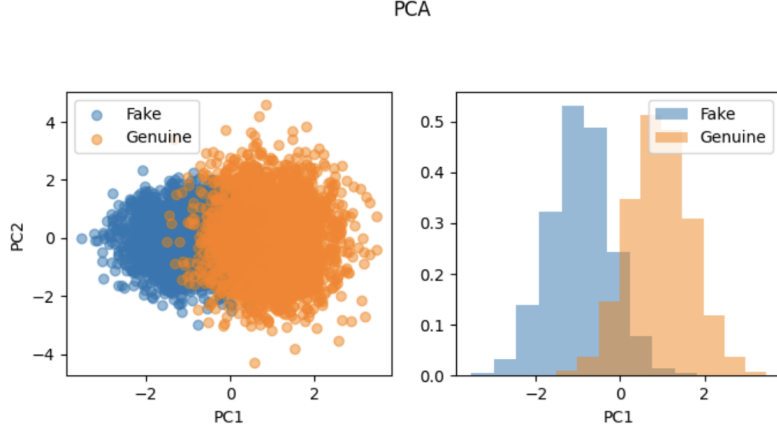


Figure 5: PCA for the two first components

After the analysis done to choose which  $m$  is best to use, we can represent the first two principal components found by calculating the PCA with  $m=6$ , visualizing these components can help to understand which features contribute most to the variance in the data. It can actually be observed that there is better separation between classes by applying PCA and choosing an  $m$  that has good variance of the data.

### 2.2 LDA

LDA is a supervised dimensionality reduction technique that aims to find the subspace that maximizes the separation between classes.

To compute the transformation matrix  $W$  we need to compute between-class and within-class covariance matrices.

$$\mathbf{S}_b = \frac{1}{N} \sum_{c=1}^K n_c (\mu_c - \mu)(\mu_c - \mu)^T \quad (3)$$

$$\mathbf{S}_w = \frac{1}{N} \sum_{c=1}^K \sum_{i=1}^{n_c} (x_{c,i} - \mu_i)(x_{c,i} - \mu_i)^T \quad (4)$$

where  $\mu_c$  is the mean of the  $c$ -th class,  $\mu$  is the mean of all the data, and  $n_c$  is the number of samples in the  $c$ -th class.

The LDA directions can be computed solving the generalized eigenvalue problem  $S_b w = \lambda S_w w$ , usually to solve the generalized eigenvalue problem we transform the  $S_w$  matrix into an identity matrix and then solve the eigenvalue problem for the other  $S_b$  matrix. This process allows us to find the discriminant vectors  $w$  that maximize the separation between the classes.

### 3 Multivariate Gaussian Density

## 4 Classification model analysis

### 4.1 Gaussian Models

#### 4.1.1 MVG Gaussian Classifier

#### 4.1.2 Naive Bayes Classifier

#### 4.1.3 Tied Bayes Classifier

#### 4.1.4 Gaussian Models Comparison

### 4.2 Logistic Regression Classifier

#### 4.2.1 Quadratic Logistic Regression(QLR)

### 4.3 SVM Classifier

#### 4.3.1 Linear SVM

#### 4.3.2 Kernel SVM

### 4.4 GMM Classifier

## 5 Score Calibration

### 5.1 Calibration Analysis on Selected Models

### 5.2 Calibrating Scores for Selected Models

## 6 Experimental Results

### 6.1 Calibration on evaluation score

### 6.2 Considerations

## 7 Conclusions