

# Lab 2

PHONEME dataset  
Classification with MAP criterion  
PCA and MDA feature selection



Escola Tècnica Superior d'Enginyeria  
de Telecomunicació de Barcelona

UNIVERSITAT POLITÈCNICA DE CATALUNYA

# Introduction. Phoneme dataset

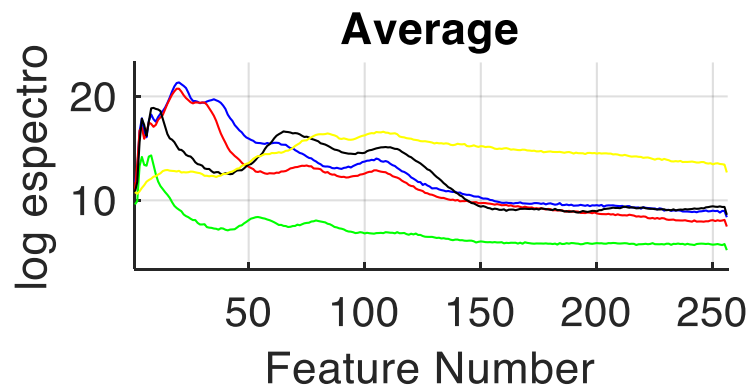
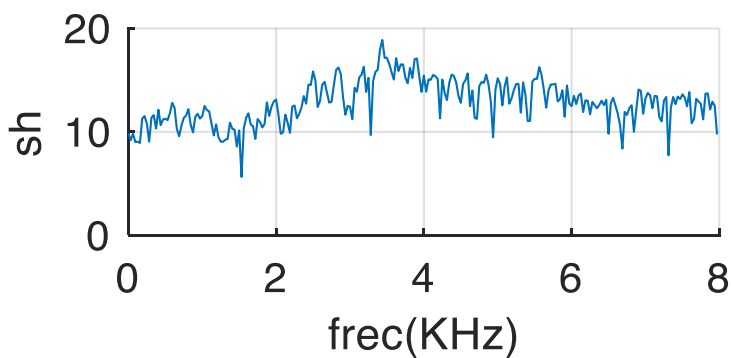
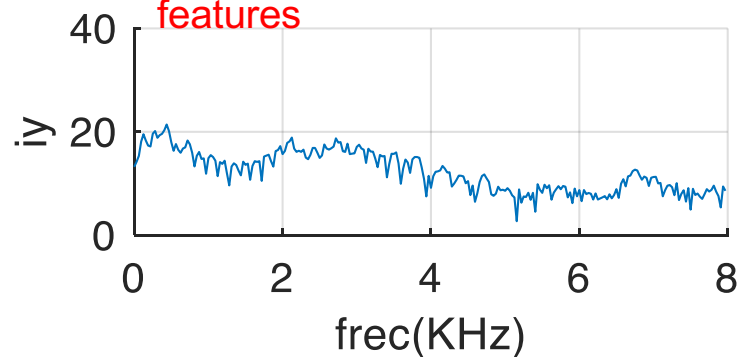
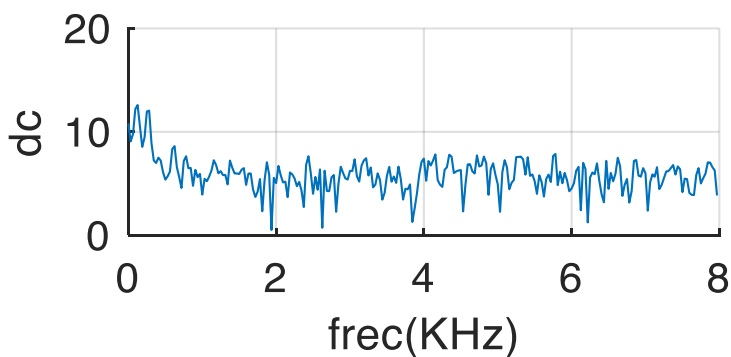
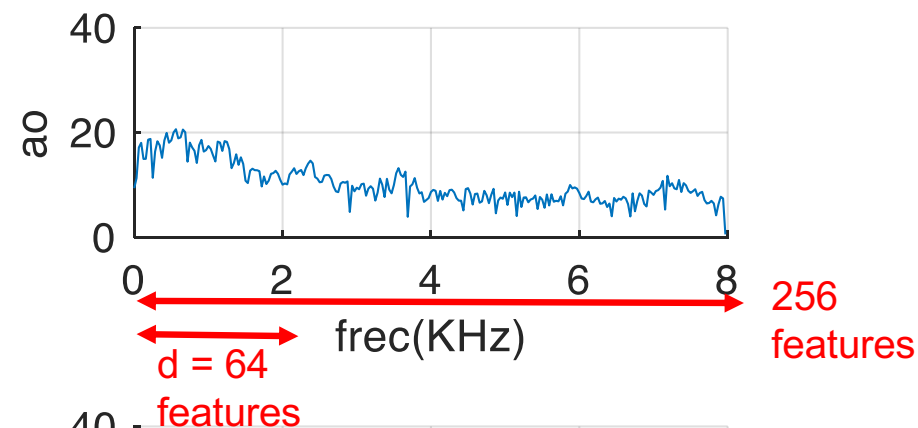
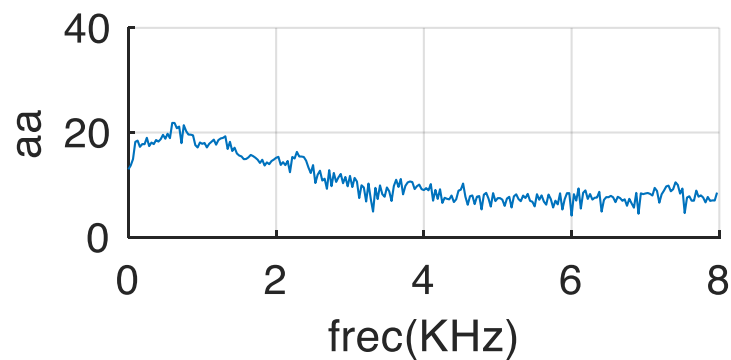
## Objectives:

- Work with a real dataset
- Split the dataset into training and test subsets
- Dimensionality reduction (feature selection) using PCA and MDA. Application to synthetic and real datasets

## PHONEME dataset:

- Each vector has been obtained computing  $\log \left( |TF(x(n))|^2 \right)$  where the sequence  $x(n)$  corresponds to part of a recording of a phoneme at a sampling rate of 16 kHz.
- Vectors correspond to 5 possible phonemes or classes: 'aa' (695) 'ao' (1022) 'dcl' (757) 'iy' (1163) 'sh' (872).
- For each vector we initially have 256 features, corresponding to the spectrum between 0 and 8 kHz.
- In the first part of Lab2 we will work just with the first 64 samples (frequencies 0 to 2 kHz).

# Example: one vector per class



# Dataset partition into training and test sets

Generated variables:

$$X = \begin{bmatrix} \quad \end{bmatrix} \begin{matrix} \updownarrow N \\ \leftarrow d = 256 \end{matrix}$$

$$\text{Labels} = \begin{bmatrix} \quad \end{bmatrix} \begin{matrix} \updownarrow N \end{matrix}$$

70%

30%

$X_{\text{train}} =$

$\text{Labels}_{\text{train}} =$

$$\begin{bmatrix} \quad \end{bmatrix} \begin{matrix} \updownarrow N_{\text{train}} \\ \leftarrow d = 256 \end{matrix}$$

$$\begin{bmatrix} \quad \end{bmatrix} \begin{matrix} \updownarrow N_{\text{train}} \end{matrix}$$

$X_{\text{test}} =$

$\text{Labels}_{\text{test}} =$

$$\begin{bmatrix} \quad \end{bmatrix} \begin{matrix} \updownarrow N_{\text{test}} \\ \leftarrow d = 256 \end{matrix}$$

$$\begin{bmatrix} \quad \end{bmatrix} \begin{matrix} \updownarrow N_{\text{test}} \end{matrix}$$

# Classifier design

## Design of a linear (LC) and a quadratic (QC) classifier

- Dataset with  $d=256$  (or 64) features
- Reduced dataset using manual selection of 2 features
- Reduced dataset (dimensión  $d'$ ) using PCA and MDA

# Dimensionality reduction

## Objective:

- Reduce the number of features (assuming column vectors) :

$$\mathbf{x}_k \text{ (} d \text{ características)} \quad \Rightarrow \quad \mathbf{z}_k = \underset{d' \times d}{\mathbf{W}^T} \mathbf{x}_k \text{ (} d' \text{ características)}$$

- **Be careful !!! In the lab we will work with row-vectors, so we need to compute the right product of  $\mathbf{W}$**

- **The dimensionality reduction helps to**
  - Simplify the classifier structure
  - Reduce the computational cost
  - Remove redundant information
- **Take into account...**
  - The reduction matrix  $\mathbf{W}$  must be created using the training dataset

# Scatter matrix

$$\mathbf{m}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in D_i} \mathbf{x} \quad \text{Mean of samples from class } i$$

$$\mathbf{m} = \frac{1}{N} \sum_{\mathbf{x} \in \{D_1, \dots, D_c\}} \mathbf{x} = \frac{1}{N} \sum_{i=1}^c N_i \mathbf{m}_i \quad \text{Mean of all samples}$$

$$\mathbf{S}_T = \sum_{\mathbf{x} \in \{D_1, \dots, D_c\}} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T \quad \text{Total data dispersion}$$

---

$$\mathbf{S}_T = \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T + \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T$$

$\mathbf{S}_C$



**Sum of intra-class  
scatter matrices**

$\mathbf{S}_B$



**Inter-class scatter  
matrix**

# PCA (Principal Component Analysis)

## Objective:

- Maximize:  $\text{trace}(\mathbf{W}^T \mathbf{S}_T \mathbf{W})$
- Constraints:  $\mathbf{w}_i^T \mathbf{w}_i = E$

## Solution (Matlab function `pca.m`):

- Columns of  $\mathbf{W}$ : eigenvalues associated to the largest eigenvalues  $d'$  of  $\mathbf{S}_T$ :

$$\mathbf{S}_T \mathbf{w}_i = \lambda_i \mathbf{w}_i$$

## Problem:

- PCA minimizes the approximation squared error but it does not guarantee the separability of the classes



# PCA (Principal Component Analysis)

## PCA Transformation

- Obtention of the PCA matrix from the training dataset

$$\mathbf{W\_pca} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_d] \text{ (pca.m)}$$

- Projection to a smaller dimension  $d'$  using  $\mathbf{W\_red} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{d'}]$

$$\mathbf{W\_red} = \mathbf{W\_pca}(:, 1..d')$$

- Transformation of training and test datasets:

```
X_train_pca=X_train*W_red;  
X_test_pca=X_test*W_red;
```

- Representation of training and test dataset for linear and quadratic classifiers, varying the dimensionality of the feature space, from  $d'=1$  to  $d'=d$ .

# MDA (Multiple Discriminant Analysis)

## Objective:

- Maximize intra-class separability while minimizing the inter-class scatter
- We measure the separability and scatter using the ellipsoid volumes, assuming data Gaussianity

## Formulation:

- Maximization: 
$$\mathbf{W} = \arg \max_{\mathbf{W}} \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_C \mathbf{W}|}$$

## Solution (Matlab function mda\_clp.m):

- $d' \leq \min(d, c-1)$  (c: number of classes)
- **W columns:** eigenvectors associated to the largest eigenvalues:

$$\mathbf{S}_B \mathbf{w}_j = \sigma_j \mathbf{S}_C \mathbf{w}_j \quad \Rightarrow \quad \mathbf{S}_C^{-1} \mathbf{S}_B \mathbf{w}_j = \sigma_j \mathbf{w}_j$$

# Lab2

## Part1:

- Use Phoneme dataset ( $c=5$  classes,  $d=256$  features)
- Train Lc and Qc classifiers using the first  $d=64$  features
- Train Lc and Qc classifiers using  $d'=2$  manually selected features

## Part2:

- Use synthetic Gaussian datasets ( $c=3$  classes,  $d=3$  features) for different SNR values
- Train Lc and Qc classifiers using all the features
- Train Lc and Qc classifiers after dimensionality reduction using PCA and MDA

## Part3:

- Use Phoneme dataset
- Train Lc and Qc classifiers using  $d'$  features selected with PCA and MDA
- Show Lc and Qc training/test error curves for varying number of features selected with PCA and MDA