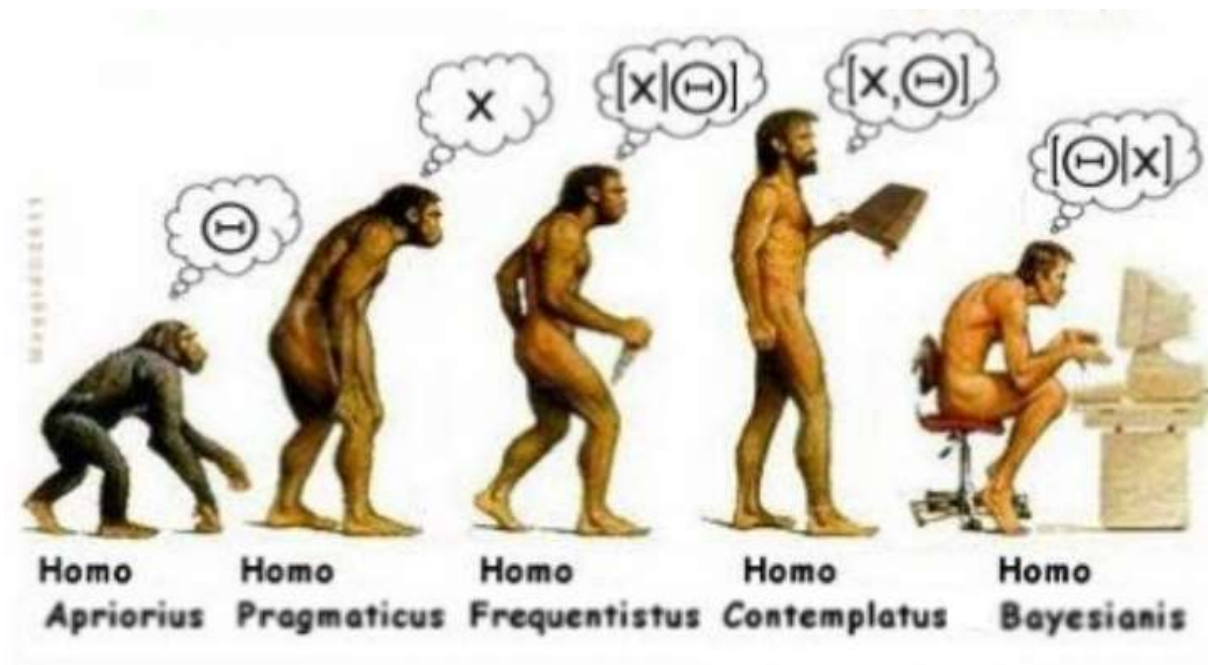


# Chapter 2

## Decision theory



**Recommended bibliography:** *Pattern Classification (2nd ed)* by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000, Chapters 2 & 3

**Credits:** Some figures are taken from *Pattern Classification (2nd ed)* by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000 with the permission of the authors

# CONTENTS

## **2.1 Bayesian decision**

2.1.1 Introduction

2.1.2 Maximum A Posteriory (MAP) decision rule

2.1.3 Minimum risk classifier

2.1.4 Discriminants and decision regions

2.1.5 Gaussian density function

2.1.6 Discriminants for Gaussian classes

2.1.7 Performance indicators

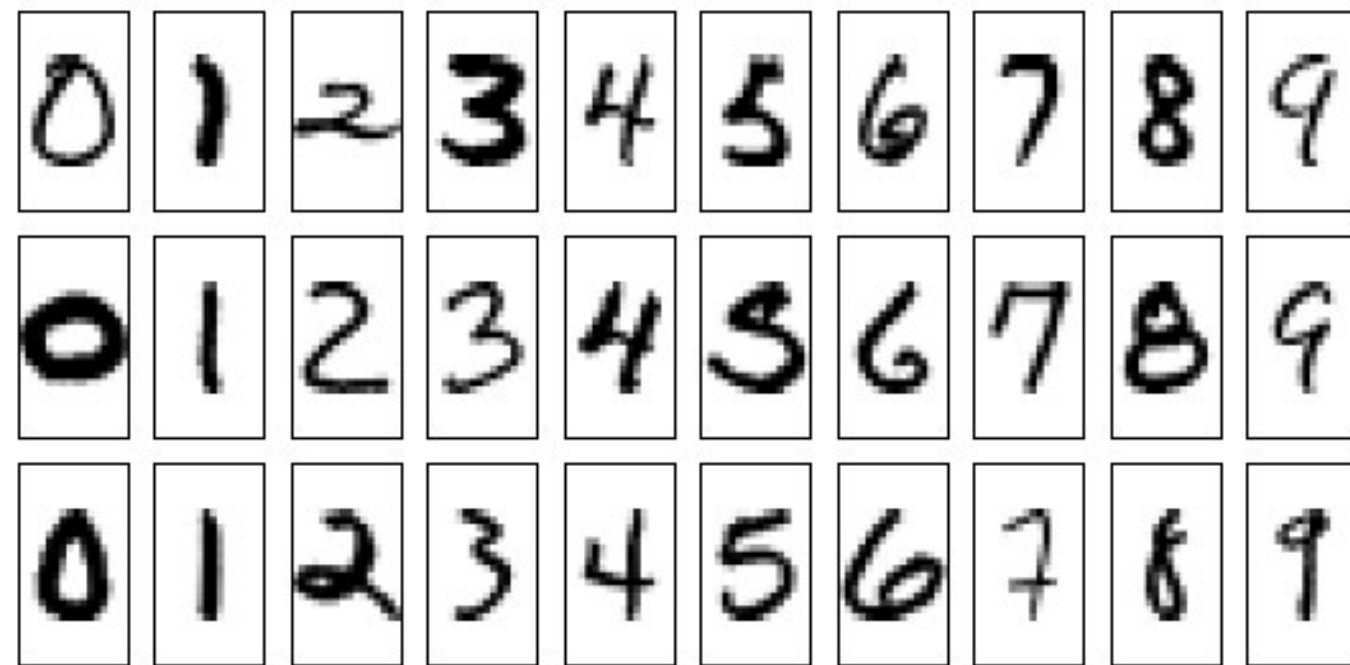
2.1.8 Conclusions

## **2.2. Maximum likelihood (ML) estimation and Bayesian estimation**

# 1.1 INTRODUCTION

Why a probabilistic approach to decision taking?

1. We might have incomplete representations of reality (e.g., we do not have the DNA of the caught fishes)
2. We face problems intrinsically random (e.g., the identification of handwritten characters)



- Measurements: random vectors of size  $d$

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \in \mathbb{R}^d$$

- State of nature: a random variable taking  $c$  possible values

$$\omega_1, \dots, \omega_c$$

- Prior probabilities

$$\Pr(\omega_1), \dots, \Pr(\omega_c) \quad \sum_{i=1}^c \Pr(\omega_i) = 1$$

- Conditioned density functions

$$f_{\mathbf{x}|\omega_1}(\mathbf{x}|\omega_1), \dots, f_{\mathbf{x}|\omega_c}(\mathbf{x}|\omega_c)$$

- Posterior probabilities

$$\Pr(\omega_j|\mathbf{x}) = \frac{f_{\mathbf{x}|\omega_j}(\mathbf{x}|\omega_j)\Pr(\omega_j)}{f_{\mathbf{x}}(\mathbf{x})}$$

$$f_{\mathbf{x}}(\mathbf{x}) = \sum_{i=1}^c f_{\mathbf{x}|\omega_i}(\mathbf{x}|\omega_i)\Pr(\omega_i)$$

$$POSTERIOR = \frac{LIKELIHOOD \times PRIOR}{EVIDENCE}$$

- **POSTERIOR**: Probability of a certain nature state given the observed feature vector.
- **LIKELIHOOD**: Contains the characterization of data for a given class.
- **PRIOR**: Prior knowledge of the class.
- **EVIDENCE**: Scaling factor independent of the class.

Two interpretations of probability appearing in our formulation...

- Frequentist

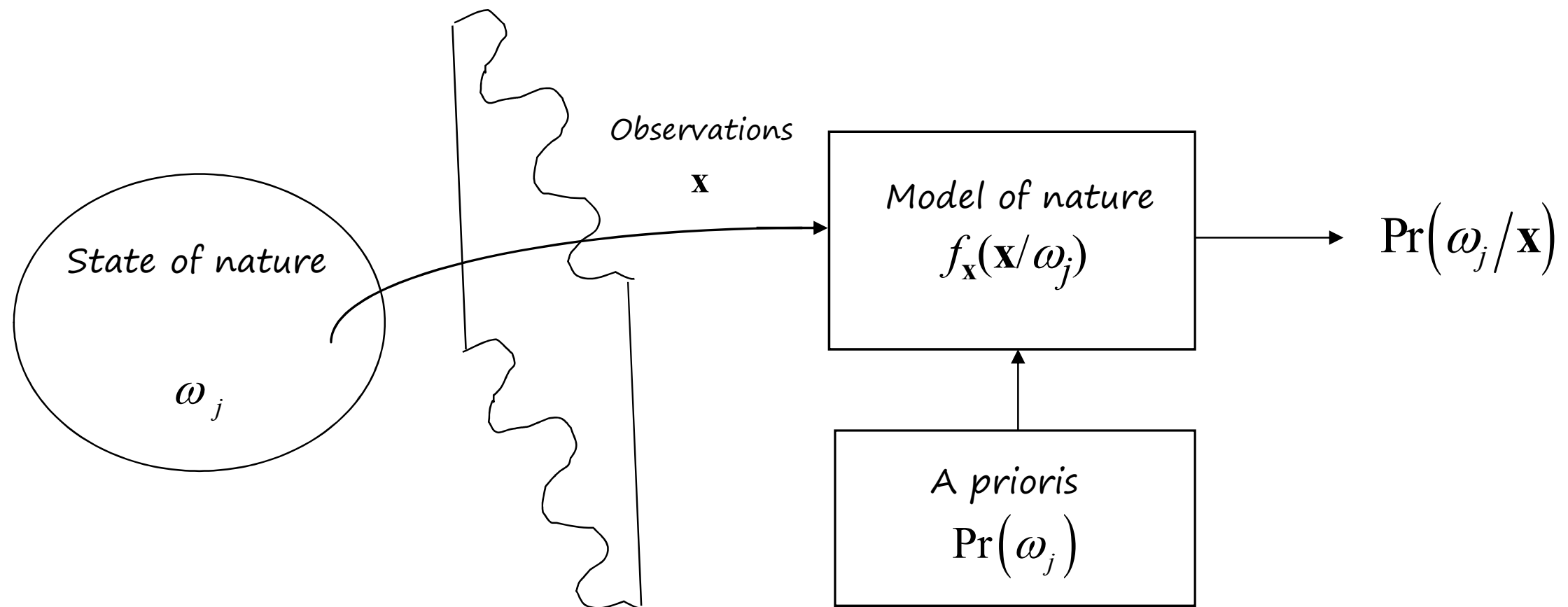


- Reliability (o belief)



What are we computing?

$$\Pr(\omega_j | \mathbf{x}) = \frac{f_{\mathbf{x}}(\mathbf{x} | \omega_j) \Pr(\omega_j)}{f_{\mathbf{x}}(\mathbf{x})}$$







### Thomas Bayes (1702-1761)

In 1763, two years after his death, *Essay Towards Solving a Problem in the Doctrine of Chances* is published. It contains a theory about the causes inferred through the observed effects. More precisely formulated later by Pierre-Simon Laplace.



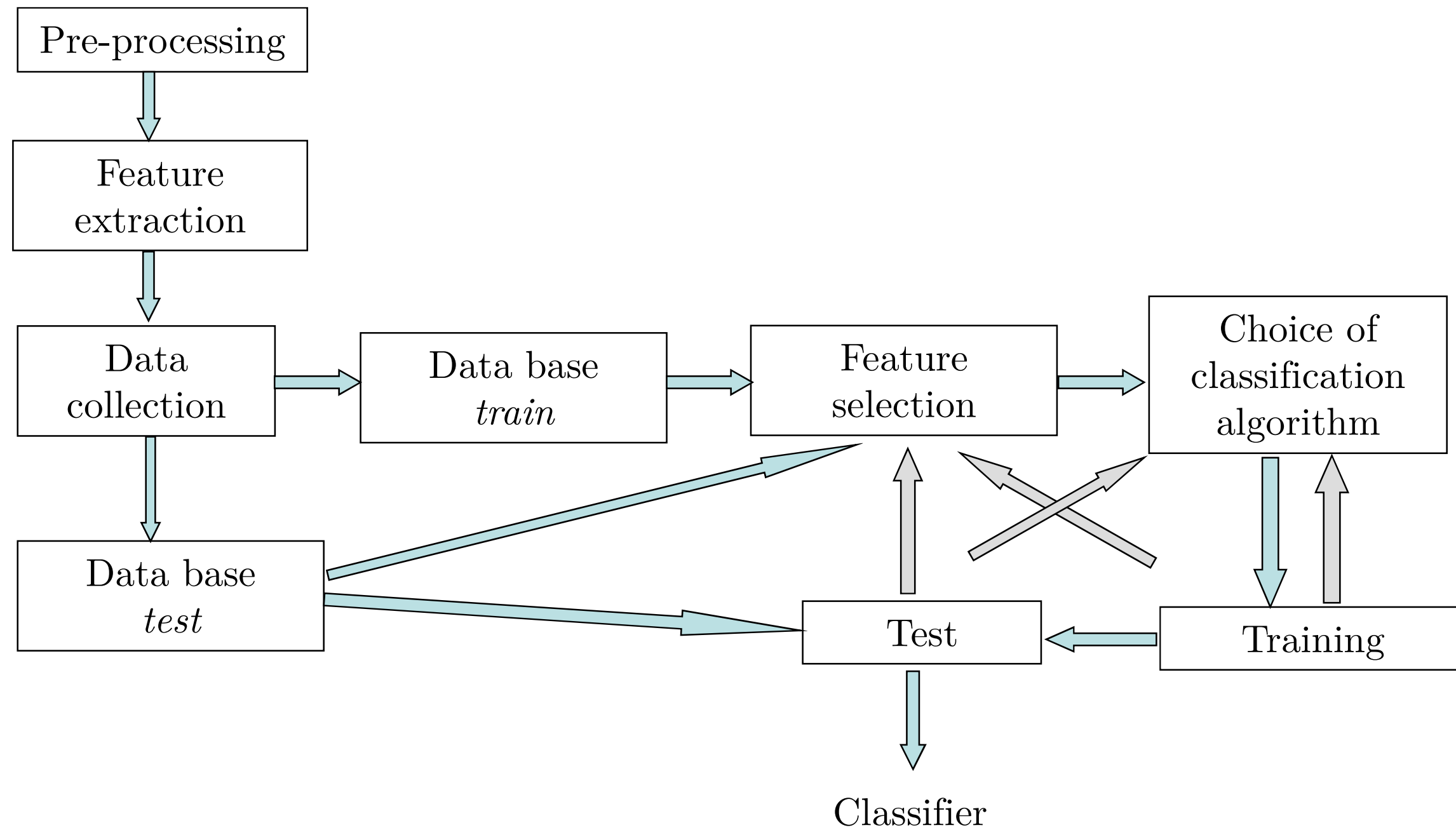
### Immanuel Kant (1724-1804)

Elaborates a synthesis between Descartes' rationalism and Hume's empirism, and publishes it in its treaty *Critique of Pure Reason*, in 1781. There he discuss about the *a priori*, the knowledge not based on experience and applies it to metaphysics.





# Stages in the design of a supervised classifier

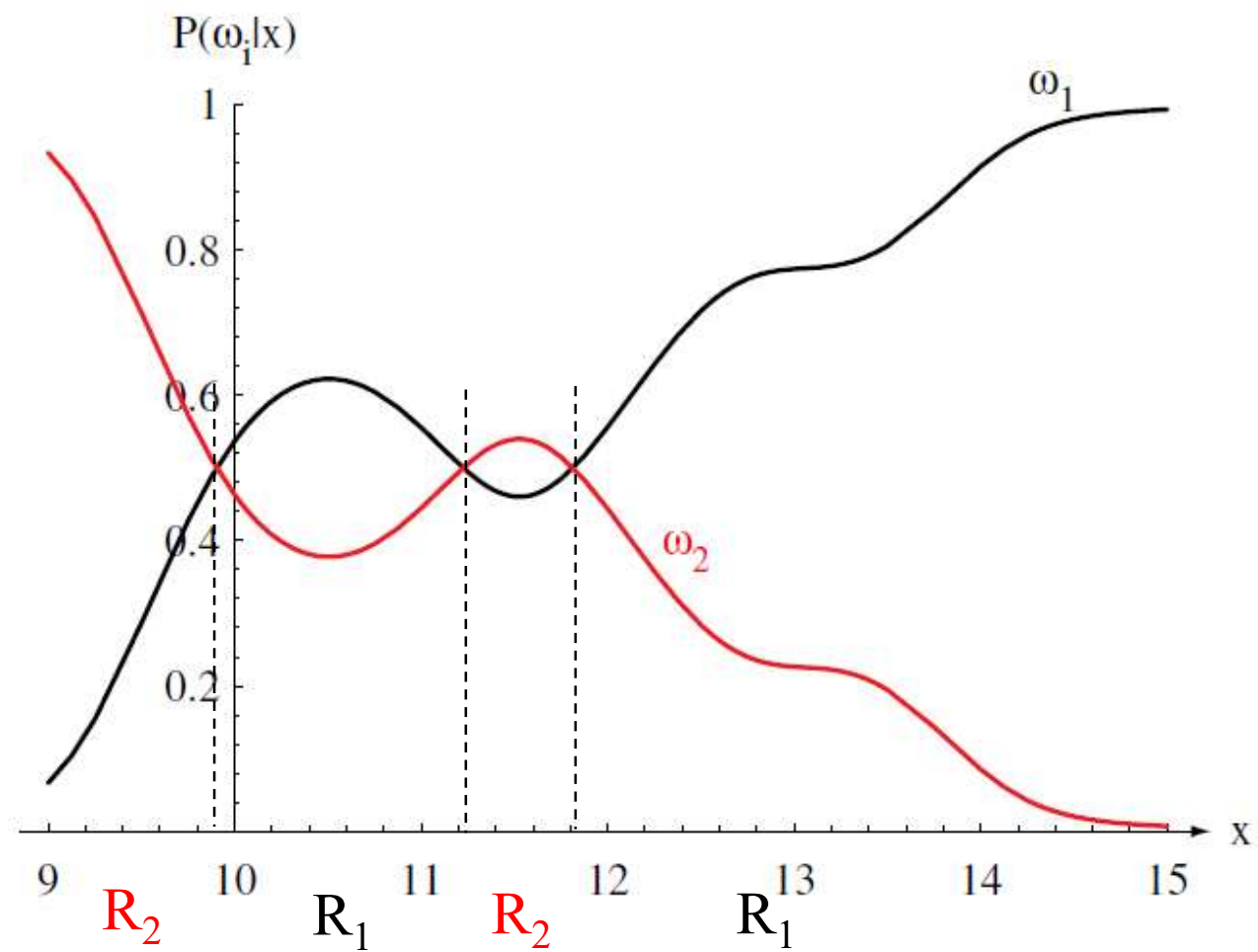


- The classifier assigns a **category (or class)** to the observed **feature vector**.
- **Training** entails adjusting a number of parameters using training feature vectors.
- The evaluation of the classifier has to be done as a function of the **chosen criterion**: minimum classification error, minimum risk, etc.
- The **difficulty in the design** of the classifier depends on the variability of observations among classes:
  - Low inter-class variability
  - High intra-class variability
- **Computational efficiency** has different impact in training and in testing.

## 1.2 MAP DECISION RULE

- Decision rule, given vector  $\mathbf{x}$   $\omega_{MAP} = \arg \max_{\omega_i} \Pr(\omega_i | \mathbf{x})$
- Probability of error, for two classes...  
conditioned to  $\mathbf{x}$ 
$$\Pr(e | \mathbf{x}) = \begin{cases} \Pr(\omega_1 | \mathbf{x}) & \text{if } \omega_2 \text{ is decided} \\ \Pr(\omega_2 | \mathbf{x}) & \text{if } \omega_1 \text{ is decided} \end{cases}$$
$$\Pr(e) = \Pr(\omega_2 | \omega_1) \Pr(\omega_1) + \Pr(\omega_1 | \omega_2) \Pr(\omega_2)$$

these are different, the second is the average of the first, over the possible values of  $\mathbf{x}$ .
- MAP minimizes the misclassification error rate.



$$\Pr(\omega_1 | \mathbf{x}) \underset{\omega_2}{\overset{\omega_1}{\gtrless}} \Pr(\omega_2 | \mathbf{x})$$

$$\omega = \arg \max_{\omega_i} \Pr(\omega_i | \mathbf{x}) = \arg \max_{\omega_i} \Pr(\omega_i) f_{\mathbf{x}}(\mathbf{x} | \omega_i)$$

### Particular cases:

- If one observation  $\mathbf{x}_0$  does not bring information about the state of nature (class)

$$f(\mathbf{x}_0 | \omega_1) = f(\mathbf{x}_0 | \omega_2) \Rightarrow \omega = \arg \max_{\omega_i} \Pr(\omega_i | \mathbf{x}) = \arg \max_{\omega_i} \Pr(\omega_i)$$

- If priors have the same value, the decision is uniquely based on the likelihood

$$\Pr(\omega_1) = \Pr(\omega_2) \Rightarrow \omega = \arg \max_{\omega_i} \Pr(\omega_i | \mathbf{x}) = \arg \max_{\omega_i} f(\mathbf{x} | \omega_i)$$

Is minimum error rate a valid criterion?

- In **biomedical diagnose**, should I equally penalize the errors healthy/ill and ill/healthy?
- **Spam** email classification
- **Optical character recognition**, is equally important the error incurred in a vowel or in a consonant?
- **RADAR**, big difference in priors if targets are present or not in an scenario



## 1.3 MINIMUM RISK CLASSIFIERS

For the two classes case, the error rate probability is

$$\begin{aligned} P(e) &= \Pr(\omega_1 \text{ happens and } \omega_2 \text{ is decided}) + \Pr(\omega_2 \text{ happens and } \omega_1 \text{ is decided}) = \\ &= \Pr(\omega_2 | \omega_1) \Pr(\omega_1) + \Pr(\omega_1 | \omega_2) \Pr(\omega_2) \end{aligned}$$

whereas for  $c$  classes: 
$$P(e) = \sum_{i=1}^c \sum_{\substack{j=1 \\ j \neq i}}^c \Pr(\omega_i | \omega_j) \Pr(\omega_j)$$

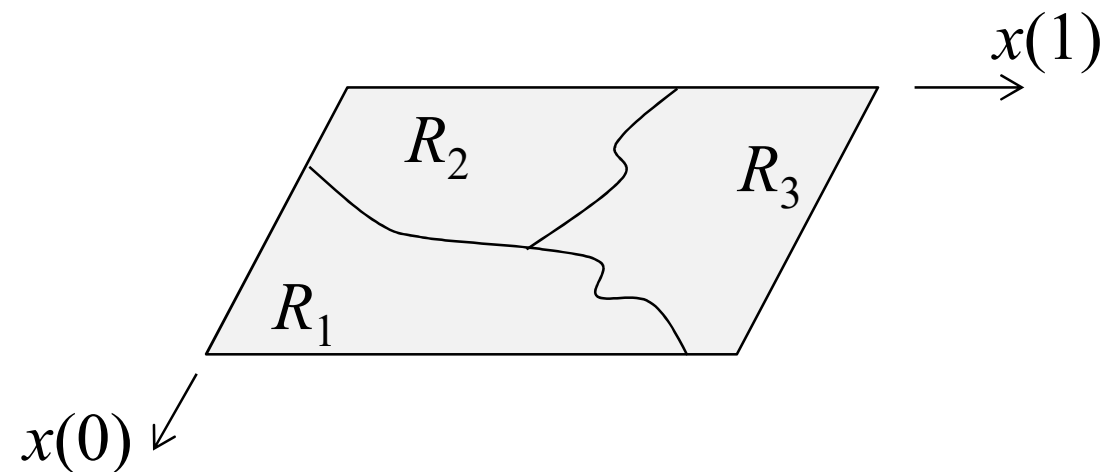
Decisions adopted may entail certain costs, and not all decisions be equally expensive. We can define the Bayes' risk function as:

$$\mathcal{R} = \sum_{i=1}^c \sum_{j=1}^c \pi_{ij} \Pr(\omega_i | \omega_j) \Pr(\omega_j)$$

where:

- $\pi_{ij}$  is the cost associated to decide  $\omega_i$  when  $\omega_j$  happens.
- $\pi_{ii}$  is the cost associated to correctly decide class  $\omega_i$

**Objective:** Design decision regions  $R_1, \dots, R_c$  so that the Bayes' risk is minimized.



$$\mathcal{R} = \sum_{i=1}^c \sum_{j=1}^c \pi_{ij} \Pr(\omega_i | \omega_j) \Pr(\omega_j)$$

$$R^d = R_1 \cup R_2 \dots \cup R_c = \bigcup_{i=1}^c R_i$$

$$R_i \cap R_j = \emptyset$$

**Theorem:** For each  $\mathbf{x}$ , the class  $\omega_i$  minimizing the risk  $\mathcal{R}$  is that associated to the least conditional risk:

$$C(\omega_i | \mathbf{x}) \triangleq \sum_{j=1}^c \pi_{ij} \Pr(\omega_j | \mathbf{x})$$

**Proof.** The space of feature vectors  $\mathbf{x}$  is split in  $M$  disjoint decision regions defined as  $R_i = \{\mathbf{x} \mid \text{decide } \omega_i\}$  for  $i = 0, \dots, c - 1$ .

$$\mathcal{R} = \sum_{i=1}^c \sum_{j=1}^c \pi_{ij} \Pr(\omega_i | \omega_j) \Pr(\omega_j) = \sum_{i=1}^c \sum_{j=1}^c \pi_{ij} \left( \int_{\mathbf{x} \in R_i} f(\mathbf{x} | \omega_j) d\mathbf{x} \right) \Pr(\omega_j) =$$

$$= \sum_{i=1}^c \int_{\mathbf{x} \in R_i} \sum_{j=1}^c \pi_{ij} f(\mathbf{x} | \omega_j) \Pr(\omega_j) d\mathbf{x} =$$

*Bayes  
theorem*

$$= \sum_{i=1}^c \int_{\mathbf{x} \in R_i} \sum_{j=1}^c \pi_{ij} \Pr(\omega_j | \mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

$C_i(\mathbf{x}) \triangleq$  Risk associated class  $\omega_i$

Therefore, to minimize  $\mathcal{R}$  :  $R_i$  is the domain of values of  $\mathbf{x}$  where  $C_i(\mathbf{x})$  is smaller than  $C_k(\mathbf{x})$  for all  $k \neq i$ :

$$R_i = \left\{ \mathbf{x} \in \mathbb{C}^{N \times 1} \mid \sum_{j=1}^c \pi_{ij} \Pr(\omega_j | \mathbf{x}) < \sum_{j=1}^c \pi_{kj} \Pr(\omega_j | \mathbf{x}), \quad \forall k \neq i \right\}$$

◆

Minimizing Bayesian risk implies minimizing...  $\omega_{MRB} = \arg \min_j C(\omega_j | \mathbf{x})$

- Binary case

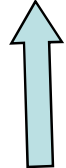
$c = 2$  categories

Conditional risk

$$C(\omega_1 | \mathbf{x}) = \pi_{11} \Pr(\omega_1 | \mathbf{x}) + \pi_{12} \Pr(\omega_2 | \mathbf{x})$$

$$C(\omega_2 | \mathbf{x}) = \pi_{21} \Pr(\omega_1 | \mathbf{x}) + \pi_{22} \Pr(\omega_2 | \mathbf{x})$$

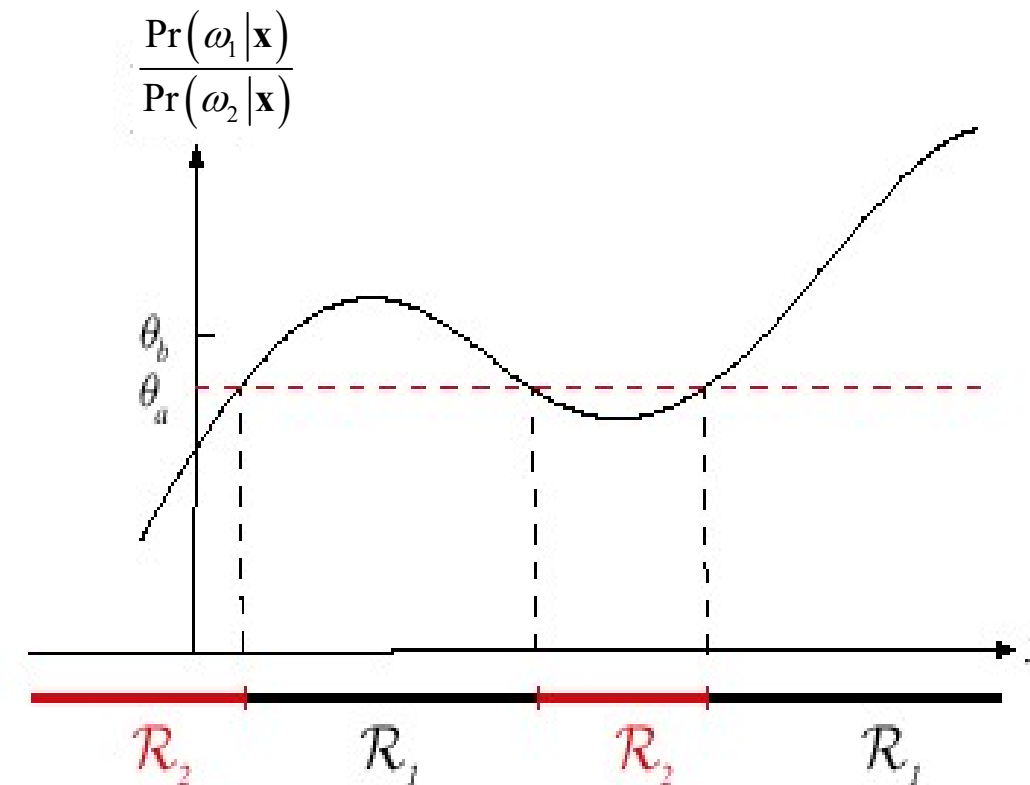
Likelihood ratio...

$$\frac{\Pr(\omega_1 | \mathbf{x})}{\Pr(\omega_2 | \mathbf{x})} \underset{\omega_2}{\overset{\omega_1}{\gtrless}} \left( \frac{\pi_{12} - \pi_{22}}{\pi_{21} - \pi_{11}} \right) = \gamma$$


Minimum risk criterion  
biases the likelihood ratio  
with respect to MAP

Threshold independent of  $\mathbf{X}$

- Increasing the threshold, more decisions on  $R_2$  are taken (and viceversa)



- Minimum risk = Minimum  $\Pr(e)$  if  $\pi_{ij} = \begin{cases} \alpha & i = j \\ \pi & i \neq j \end{cases}$

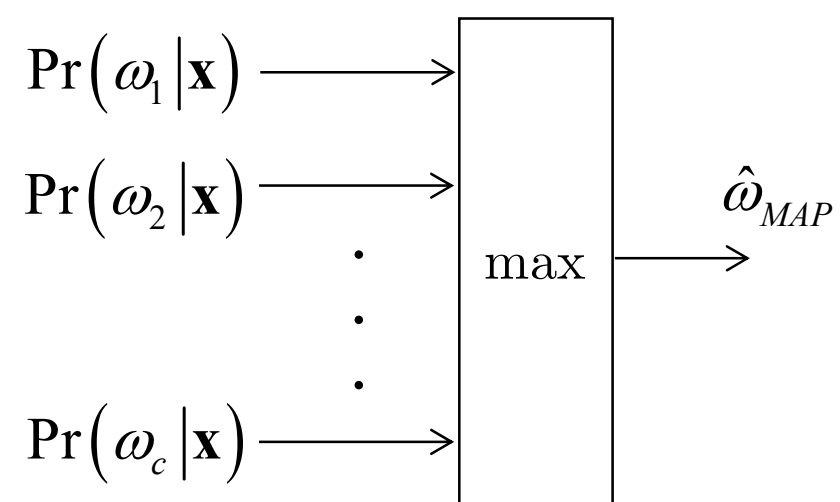
Proof:

$$C(\omega_i | \mathbf{x}) = \sum_{j=1}^c \pi_{ij} \Pr(\omega_j | \mathbf{x}) = \alpha \Pr(\omega_i | \mathbf{x}) + \sum_{\substack{j=1 \\ j \neq i}}^c \pi \Pr(\omega_j | \mathbf{x})$$

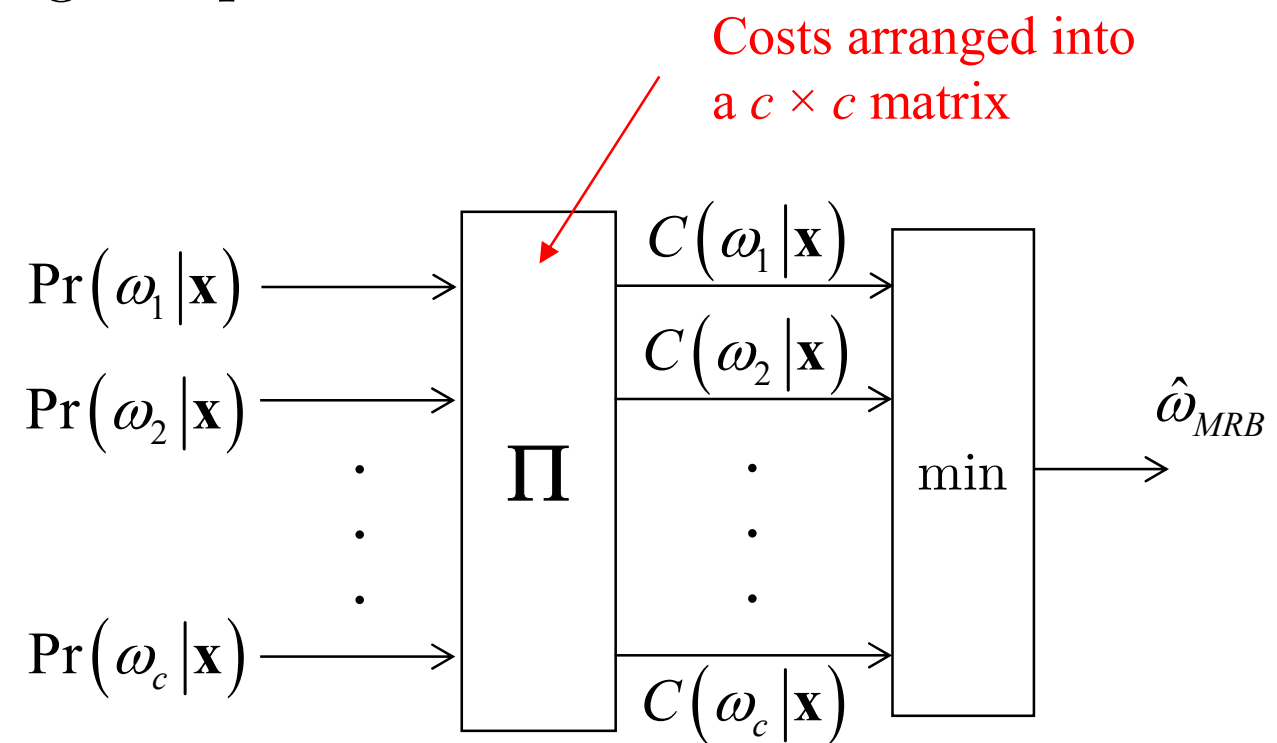
$$= \alpha \Pr(\omega_i | \mathbf{x}) + \pi (1 - \Pr(\omega_i | \mathbf{x})) = \pi - (\pi - \alpha) \Pr(\omega_i | \mathbf{x})$$

MBR criterion  
becomes MAP

- In brief...
  - MAP is a particular case of the minimum bayesian risk.
  - Both are implemented using the posteriors:



MAP



Minimum Bayesian risk





Other criteria for decision can be defined...

- NEYMAN-PEARSON
  - Total risk is minimized subject to a restriction, e.g. upper bounding the classification error for class  $i$ , priors do not intervene:

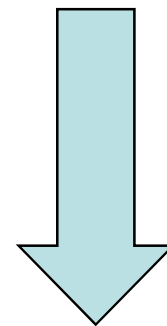
$$\int_{\mathbf{x} \in R_i} C(\omega_i | \mathbf{x}) f(\mathbf{x}) d\mathbf{x} < \text{constant}$$

- MINMAX
  - Used when a priori probabilities are not known, maybe because they are changing over time in an unknown way.
  - Minimizes the **worst total risk**, choosing the decision regions in such a way that the risk function does not depend on the a priori probabilities.
  - Example for  $c = 2$  categories...



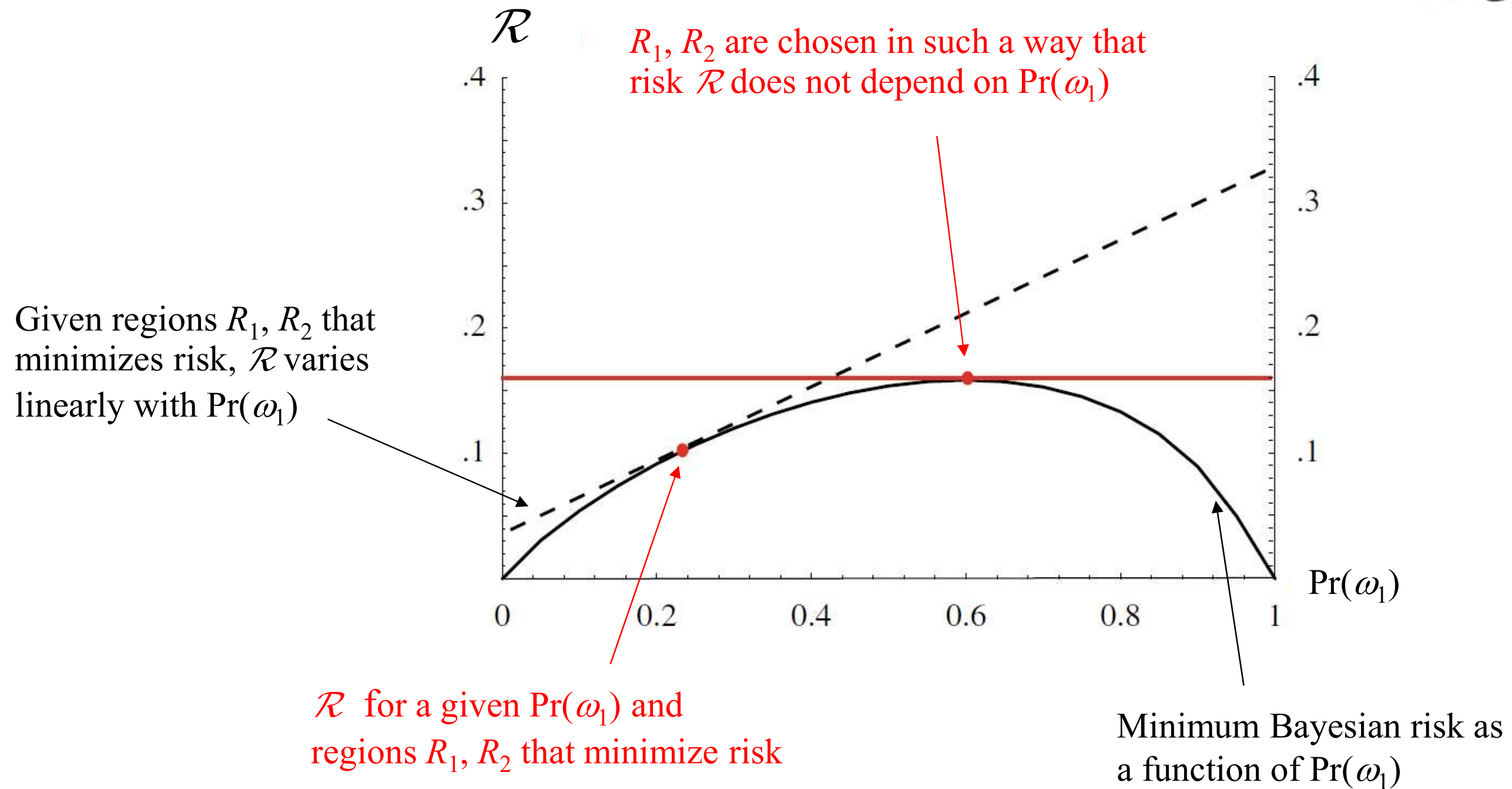
$R_1$  : region where  $\omega_1$  is decided

$$\begin{aligned}
 \mathcal{R} &= \int_{R_1} C(\omega_1 | \mathbf{x}) f(\mathbf{x}) d\mathbf{x} + \int_{R_2} C(\omega_2 | \mathbf{x}) f(\mathbf{x}) d\mathbf{x} \\
 &= \int_{R_1} (\pi_{11} \Pr(\omega_1) f_{\mathbf{x}}(\mathbf{x} | \omega_1) + \pi_{12} \Pr(\omega_2) f_{\mathbf{x}}(\mathbf{x} | \omega_2)) d\mathbf{x} + \\
 &\quad + \int_{R_2} (\pi_{21} \Pr(\omega_1) f_{\mathbf{x}}(\mathbf{x} | \omega_1) + \pi_{22} \Pr(\omega_2) f_{\mathbf{x}}(\mathbf{x} | \omega_2)) d\mathbf{x} = \\
 &= \left\{ \begin{array}{l} \Pr(\omega_1) = 1 - \Pr(\omega_2) \\ \int_{R_1} f_{\mathbf{x}}(\mathbf{x} | \omega_1) d\mathbf{x} = 1 - \int_{R_2} f_{\mathbf{x}}(\mathbf{x} | \omega_1) d\mathbf{x} \end{array} \right\} = \pi_{22} + (\pi_{12} - \pi_{22}) \int_{R_1} f_{\mathbf{x}}(\mathbf{x} | \omega_2) d\mathbf{x} + \\
 &\quad + \Pr(\omega_1) \left( \pi_{11} - \pi_{22} + (\pi_{21} - \pi_{11}) \int_{R_2} f_{\mathbf{x}}(\mathbf{x} | \omega_1) d\mathbf{x} - (\pi_{12} - \pi_{22}) \int_{R_1} f_{\mathbf{x}}(\mathbf{x} | \omega_2) d\mathbf{x} \right)
 \end{aligned}$$

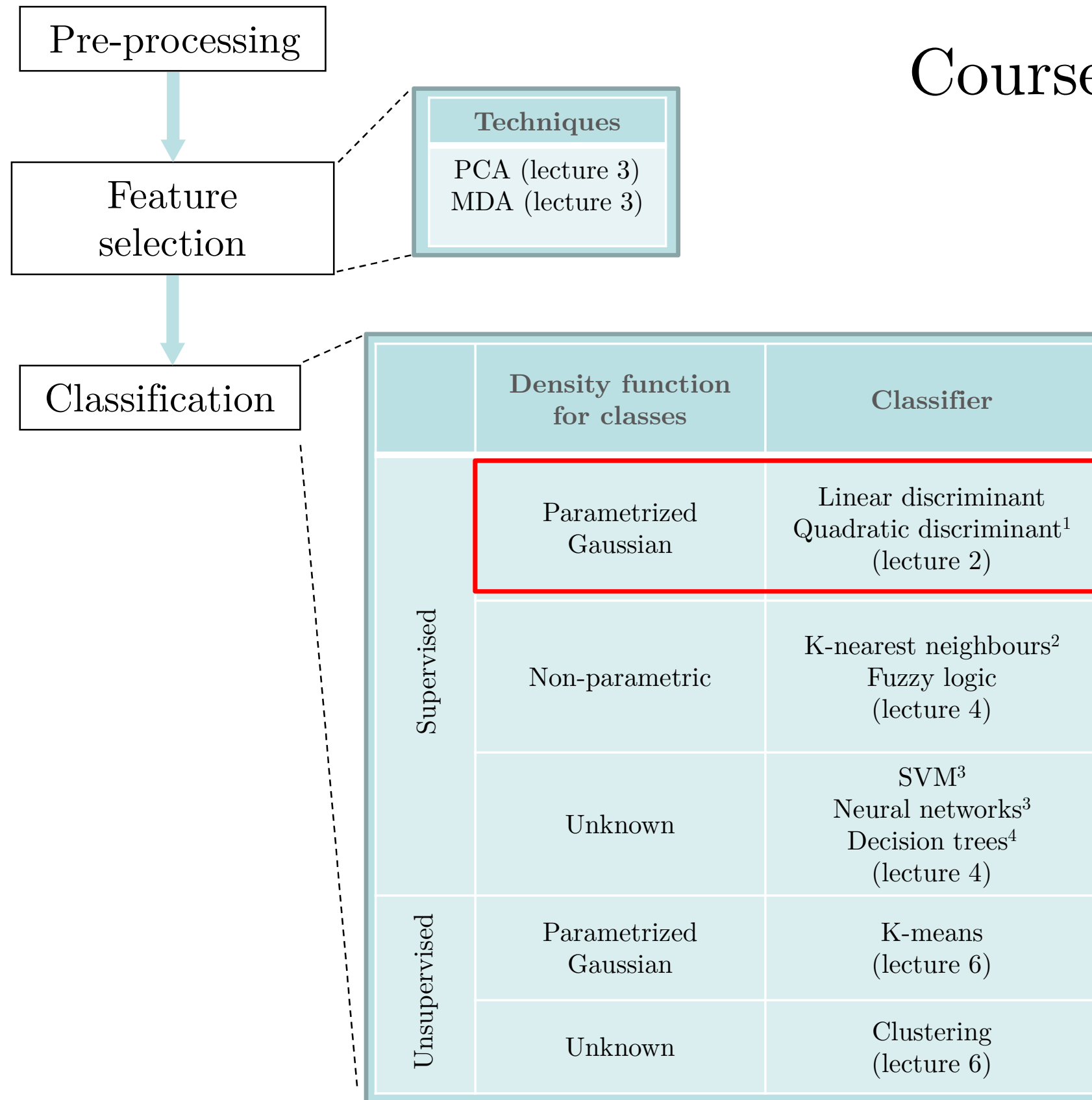


If choosing regions  $R_1$  and  $R_2$  in such a way that this term is cancelled,  $\mathcal{R}$  does not depend on  $\Pr(\omega_1)$

$$\mathcal{R}_{\text{mini-max}} = \pi_{22} + (\pi_{12} - \pi_{22}) \int_{R_1} f_{\mathbf{x}}(\mathbf{x} | \omega_2) d\mathbf{x}$$



# Course overview



1. Useful only if covariance matrices are not rank deficient.
2. Useful with the number of features is very large, even larger than the number of training vectors.
3. Imposes a structure to the classifier irrespective of the training data base.
4. Useful when non-numeric features are present.

## 1.4 DISCRIMINANTS AND DECISION REGIONS

- Definition of a discriminant function  $g_i$  for class  $i$ :
  - The classifier uses it to assign a class  $\omega_i$  to a feature vector  $\mathbf{x}$ .
  - Classification criterion: decide class  $\omega_i$  if  $g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad \forall j \neq i$

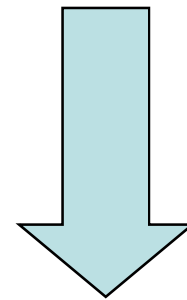
Minimum risk  $g_i(\mathbf{x}) = -C(\omega_i | \mathbf{x})$

MAP  $g_i(\mathbf{x}) = \Pr(\omega_i | \mathbf{x})$

A non-decreasing function can be applied and the criterion does not change

$$g_i(\mathbf{x}) = \Pr(\omega_i | \mathbf{x})$$

$\ln(\cdot)$  is an increasing function



$$h_i(\mathbf{x}) = \ln \Pr(\omega_i | \mathbf{x}) = \ln f_{\mathbf{x}}(\mathbf{x} | \omega_i) + \ln \Pr(\omega_i)$$

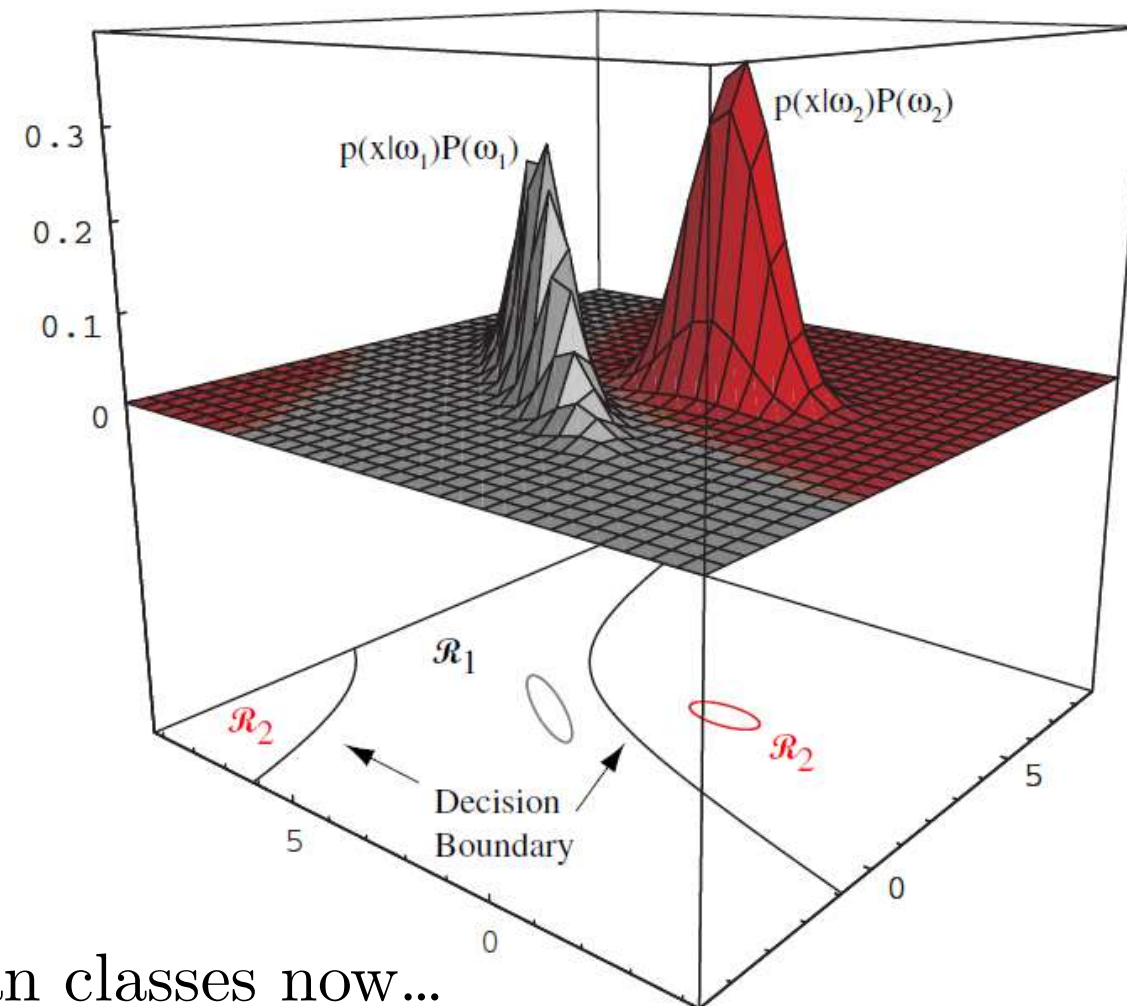


- For  $c = 2$  categories: the decision regions are given by...

$$g(\mathbf{x}) \equiv g_1(\mathbf{x}) - g_2(\mathbf{x}) \underset{\omega_2}{\overset{\omega_1}{\gtrless}} 0$$

- Decision boundary:  $g(\mathbf{x}) = 0$

- Examples
  - Binary communications  
BPSK, FSK -2
  - Detection of illness: Y/N



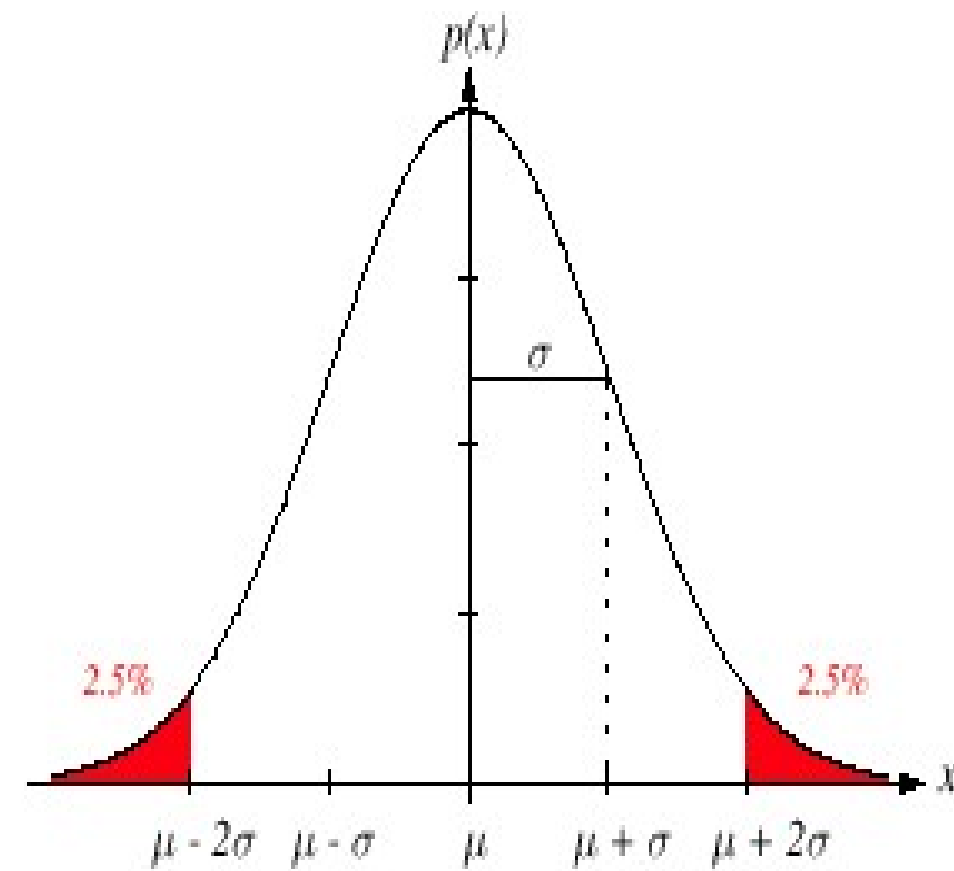
Let us consider the case of Gaussian classes now...

## 1.5 GAUSSIAN DENSITY FUNCTION

- Univariate case...

$$f_x(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right)$$

$$\mu = E\{x\} \quad \sigma^2 = E\{(x-\mu)^2\}$$



- **Multivariate case...**

- Statistical moments

$$\mathbf{x} \in \mathbb{R}^d \quad \boldsymbol{\mu} = E\{\mathbf{x}\} \in \mathbb{R}^d \quad \mathbf{C} = E\left\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\right\} \in \mathbb{R}^{d \times d}$$

- Covariance matrix is positive semi-definite (real non-negative eigenvalues)
- Density function of vector  $\mathbf{x}$ :

$$f_{\mathbf{x}}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\mathbf{C}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

- Linear transformations of Gaussian random variables are also Gaussian:

$$\mathbf{A} \in \mathbb{R}^{d \times k} \quad \mathbf{y} = \mathbf{A}^T \mathbf{x} \in \mathbb{R}^k$$

$$\boldsymbol{\mu}_y = E\{\mathbf{y}\} = E\{\mathbf{A}^T \mathbf{x}\} = \mathbf{A}^T E\{\mathbf{x}\} = \mathbf{A}^T \boldsymbol{\mu}_x$$

$$\begin{aligned} \mathbf{C}_y &= E\left\{(\mathbf{y} - \boldsymbol{\mu}_y)(\mathbf{y} - \boldsymbol{\mu}_y)^T\right\} = \\ &= E\left\{(\mathbf{A}^T \mathbf{x} - \mathbf{A}^T \boldsymbol{\mu}_x)(\mathbf{A}^T \mathbf{x} - \mathbf{A}^T \boldsymbol{\mu}_x)^T\right\} = \\ &= E\left\{\mathbf{A}^T (\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T \mathbf{A}\right\} = \mathbf{A}^T \mathbf{C}_x \mathbf{A} \end{aligned}$$

## Whitening of a feature vector

- Spectral decomposition of  $\mathbf{C}$ 
$$\mathbf{C}\mathbf{u}_i = \lambda_i\mathbf{u}_i$$
$$\mathbf{C} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$
- Orthonormal eigenvectors  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d]$   $\mathbf{U}\mathbf{U}^T = \mathbf{I}$
- Eigenvalues  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$
- A linear transform “whitens” the elements of vector  $\mathbf{X}$

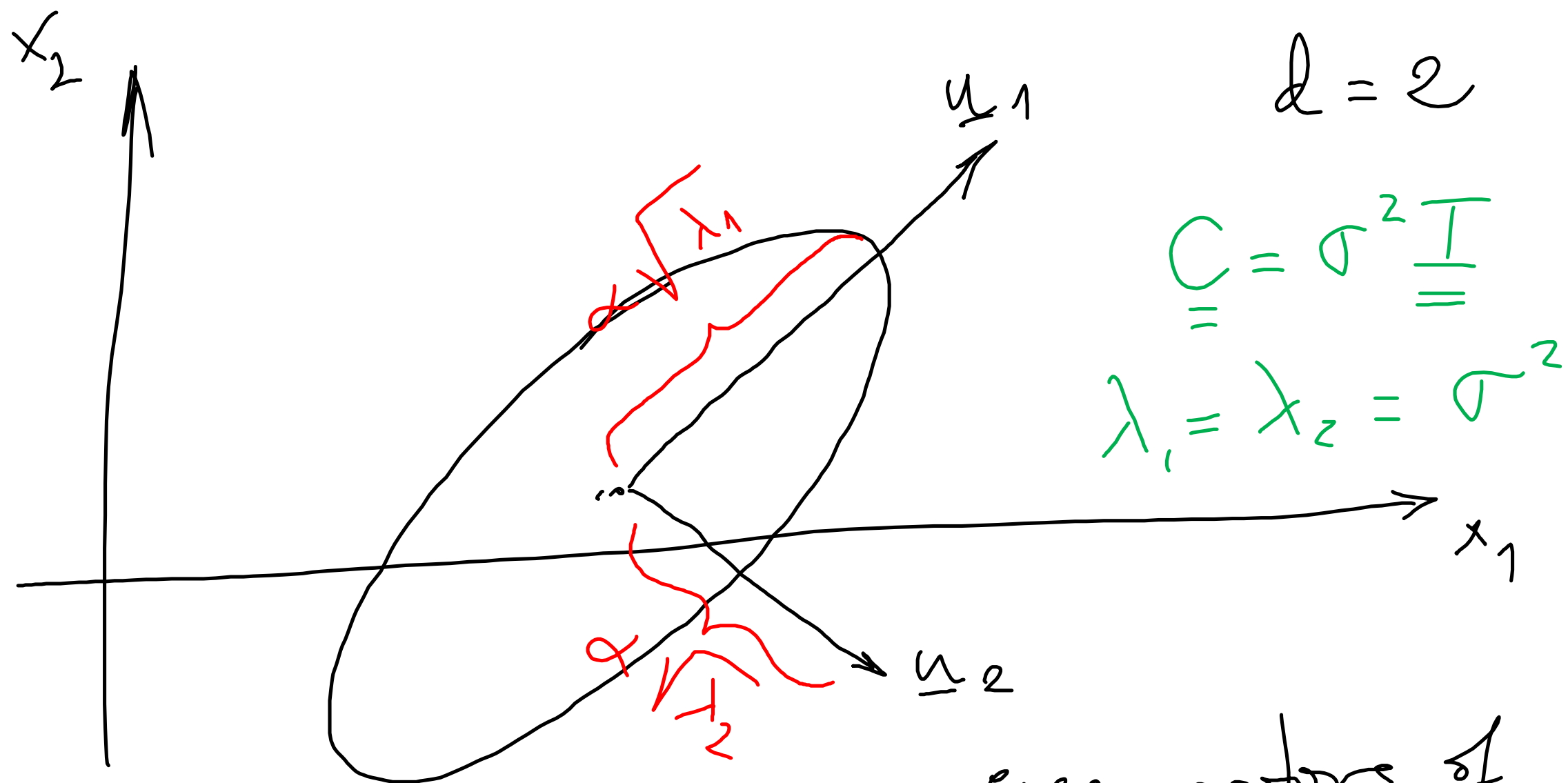
$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}^{-1/2} \quad \mathbf{\Lambda}^{-1/2} = \text{diag}\left(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_d}\right)$$

also

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}^{-1/2}\mathbf{U}^T$$

- The sections of the gaussians are hiper-ellipsoids in a space of dimension  $d$ .
- Whitening turns hiper-ellipsoids into hiper-spheres.
- The samples of a Gaussian cluster are grouped around  $\mu$
- The main axis of the hiper-ellipsoids follow the direction of the eigenvectors of  $\mathbf{C}$ .
- The length of the principal axis of the hiper-ellipsoids are proportional to the square root of the eigenvalues.





$$\underline{\underline{C}} = \sigma^2 \underline{\underline{I}}$$

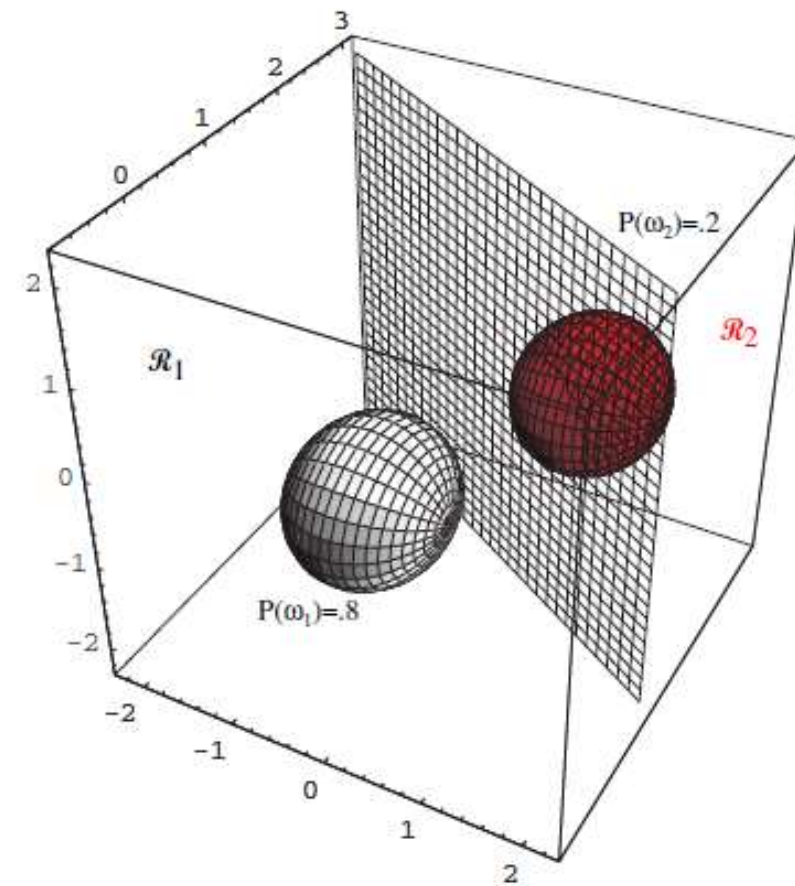
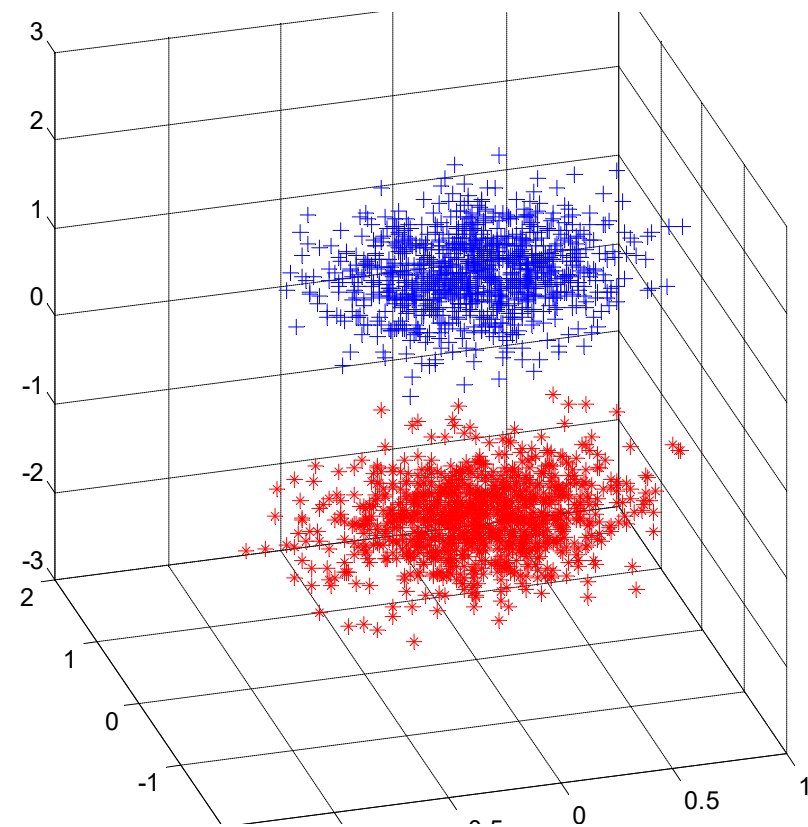
$$\lambda_1 = \lambda_2 = \sigma^2$$

$$\underline{\underline{C}}_i \in \mathbb{R}^{2 \times 2}$$

eigenvectors of  $\underline{\underline{C}}_i$  of  $\omega_i$

$$\underline{\underline{C}} \underline{u}_k = \lambda_k \underline{u}_k \quad k=1,2$$

- Gaussian clusters...



## 1.6 DISCRIMINANTS FOR GAUSSIAN CLASSES

- Density function for class  $i$ :  $f_{\mathbf{x}}(\mathbf{x}|\omega_i) \sim N(\boldsymbol{\mu}_i, \mathbf{C}_i)$

- A priori probability:  $\Pr(\omega_i)$

- Discriminant function for MAP

$$\begin{aligned} h_i(\mathbf{x}) &= \ln f_{\mathbf{x}}(\mathbf{x}|\omega_i) + \ln \Pr(\omega_i) = \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \mathbf{C}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_i| + \ln \Pr(\omega_i) \end{aligned}$$

- Three cases for the covariance matrix...

- Case 1  $\mathbf{C}_i = \sigma^2 \mathbf{I}$

- Case 2  $\mathbf{C}_i = \mathbf{C}$

- Case 3  $\mathbf{C}_i$  arbitrary

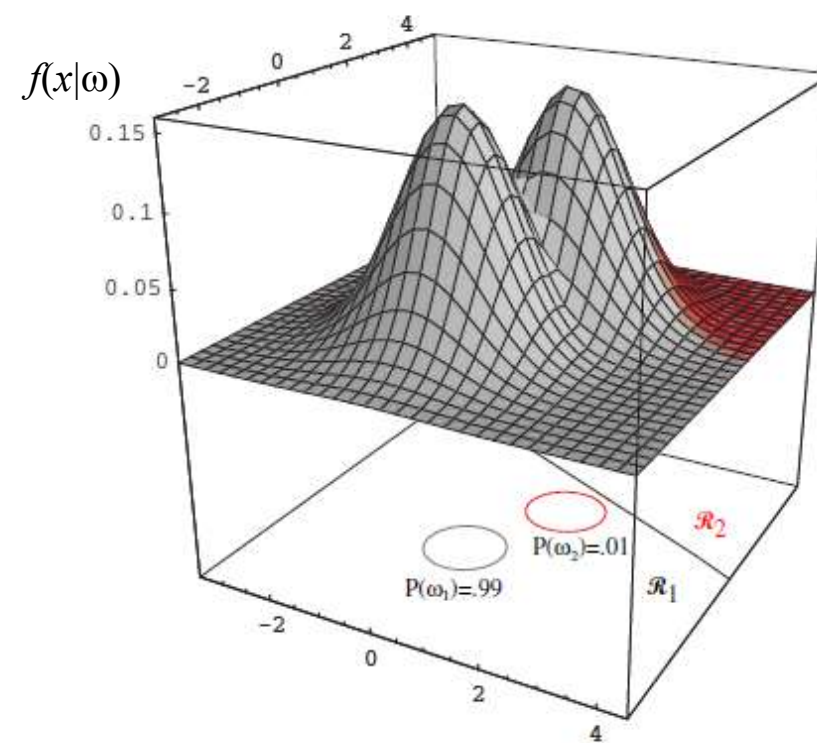
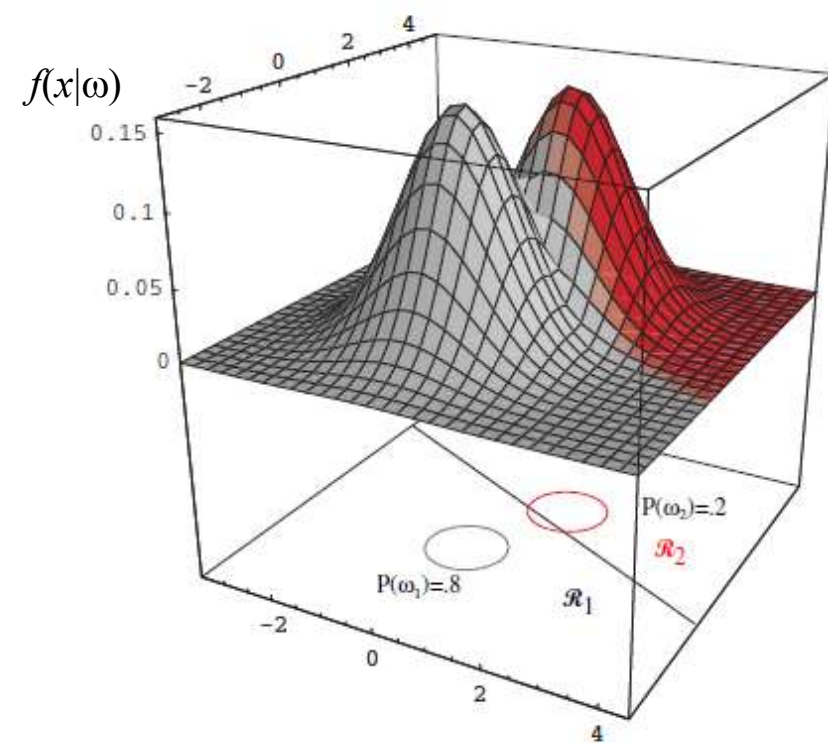
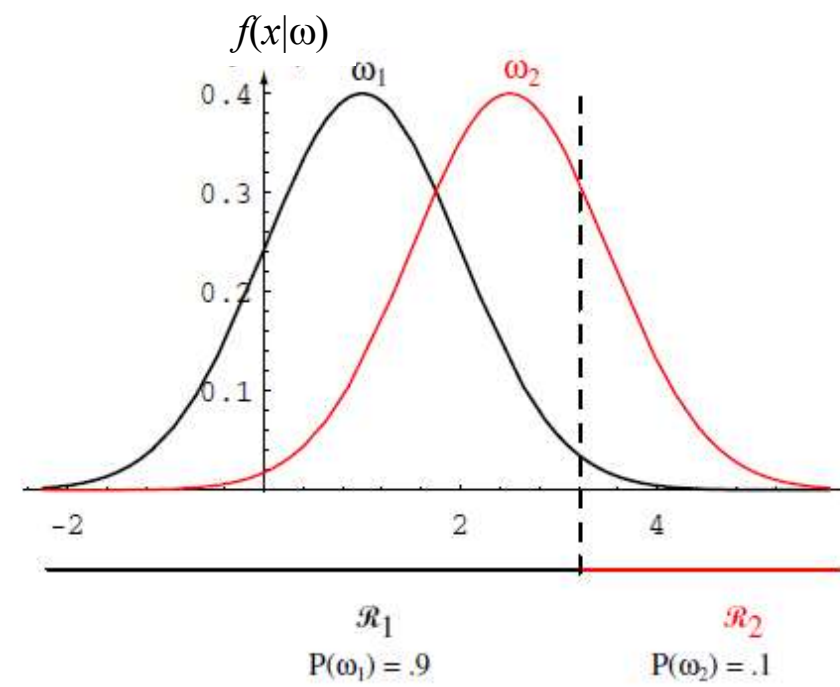
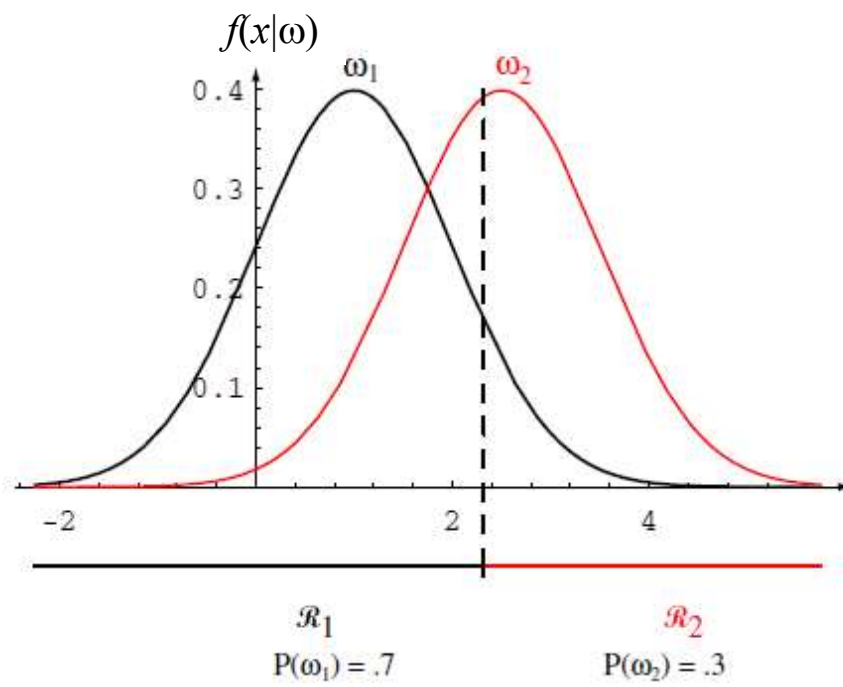
**Case 1**      $\mathbf{C}_i = \sigma^2 \mathbf{I}$

- The discriminant depends on the euclidean distance

$$h_i(\mathbf{x}) = -\frac{1}{2\sigma^2} (\mathbf{x} - \boldsymbol{\mu}_i)^T (\mathbf{x} - \boldsymbol{\mu}_i) + \ln \Pr(\omega_i)$$

- Decision boundaries are **hyperplanes**

$$h(\mathbf{x}) = h_i(\mathbf{x}) - h_j(\mathbf{x}) = 0 \quad \Rightarrow \quad \mathbf{w}^T (\mathbf{x} - \mathbf{x}_0) = 0$$





**Exercise.** Prove that:

$$h(\mathbf{x}) = h_i(\mathbf{x}) - h_j(\mathbf{x}) = 0 \quad \Rightarrow \quad \mathbf{w}^T (\mathbf{x} - \mathbf{x}_0) = 0$$

$$\mathbf{w} = \boldsymbol{\mu}_i - \boldsymbol{\mu}_j$$

$$\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\sigma^2}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \ln \frac{\Pr(\omega_i)}{\Pr(\omega_j)} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$$

It is also called minimum distance classifier.

## Case 2 $\mathbf{C}_i = \mathbf{C}$

- The discriminant depends on the Mahalanobis distance:

$$h_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln \Pr(\omega_i)$$

- Decision boundaries are **hyperplanes**

$$h(\mathbf{x}) = h_i(\mathbf{x}) - h_j(\mathbf{x}) = 0 \quad \Rightarrow \quad \mathbf{w}^T (\mathbf{x} - \mathbf{x}_0) = 0$$



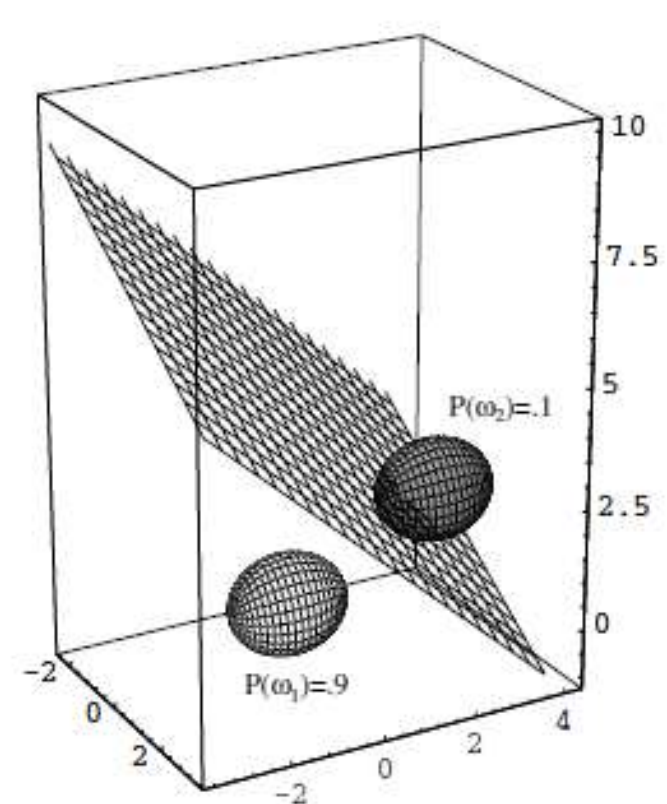
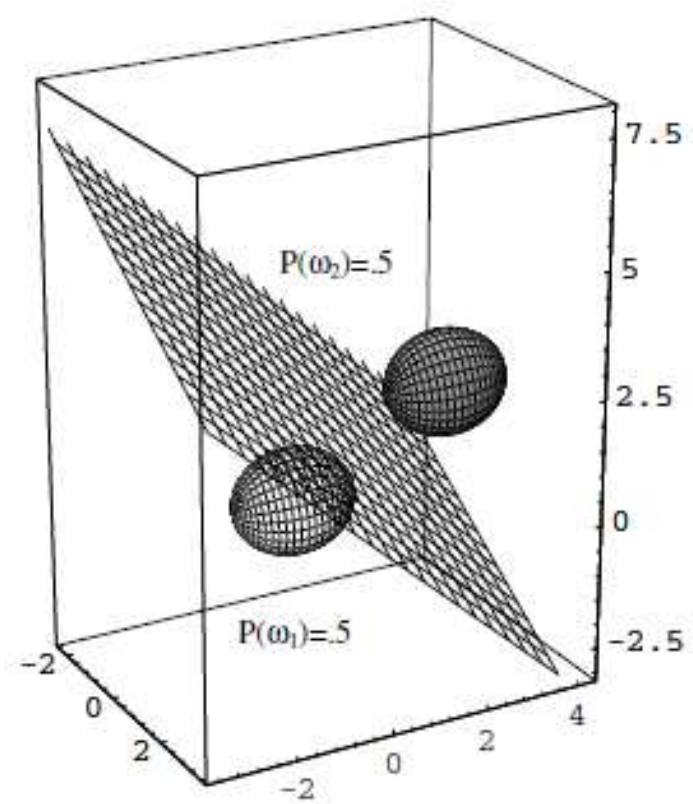
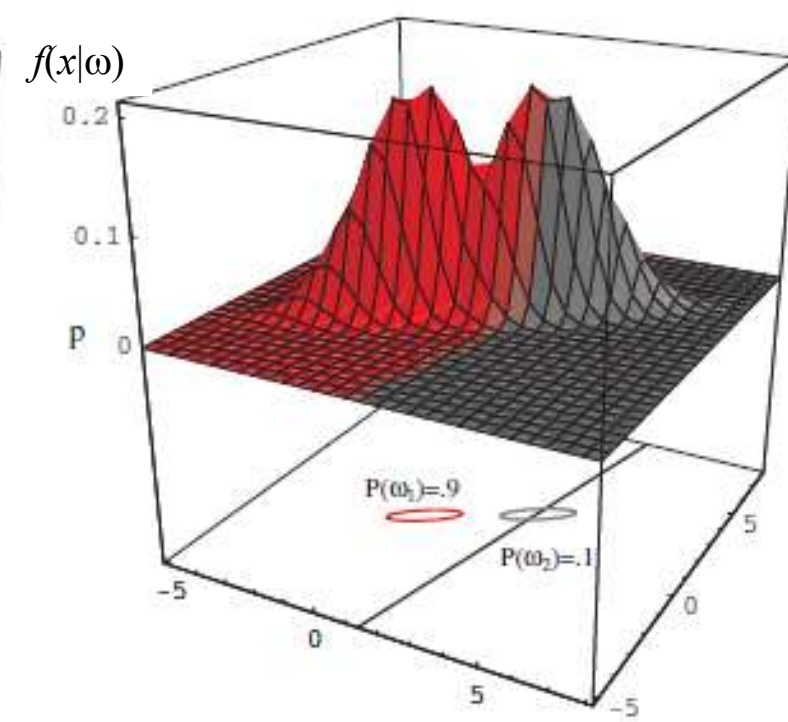
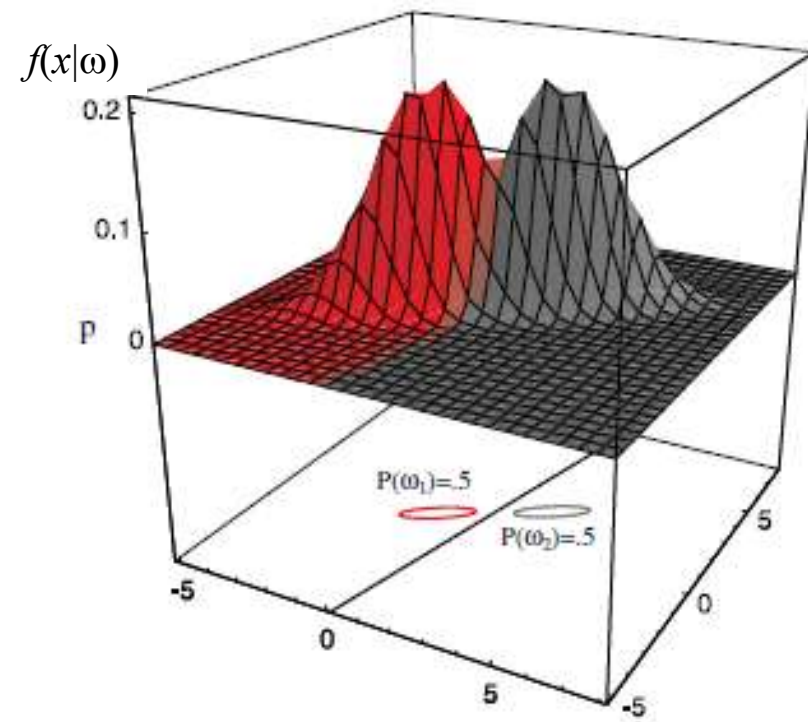


**Exercise.** Prove that:

$$h(\mathbf{x}) = h_i(\mathbf{x}) - h_j(\mathbf{x}) = 0 \quad \Rightarrow \quad \mathbf{w}^T (\mathbf{x} - \mathbf{x}_0) = 0$$

$$\mathbf{w} = \mathbf{C}^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$$

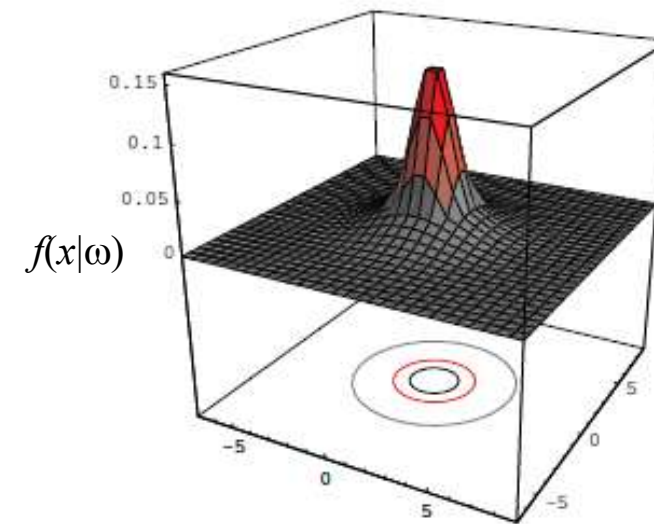
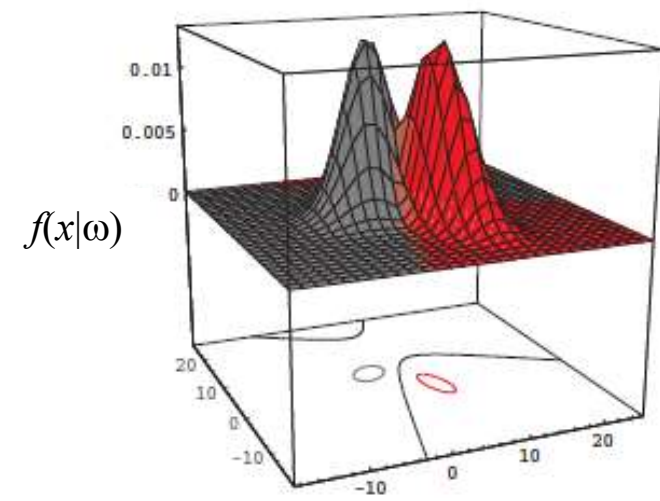
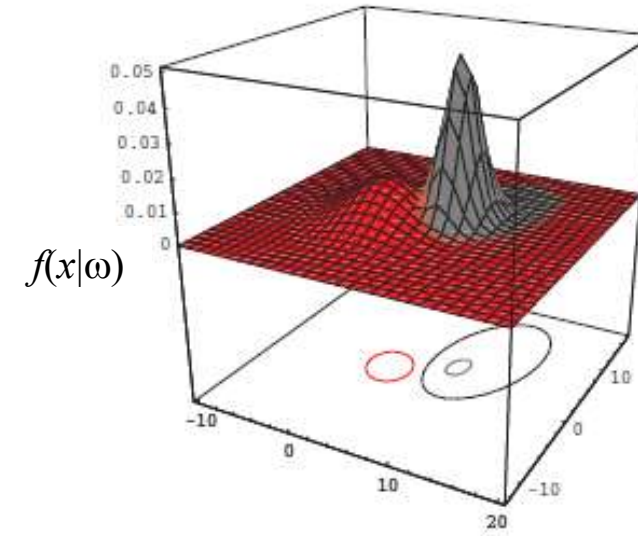
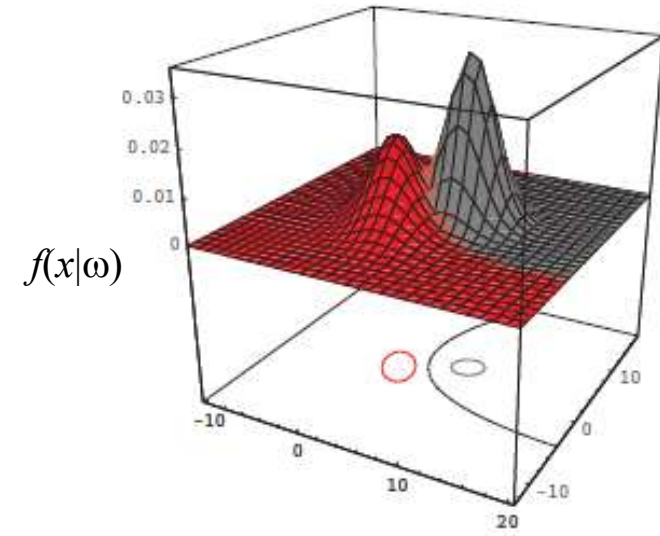
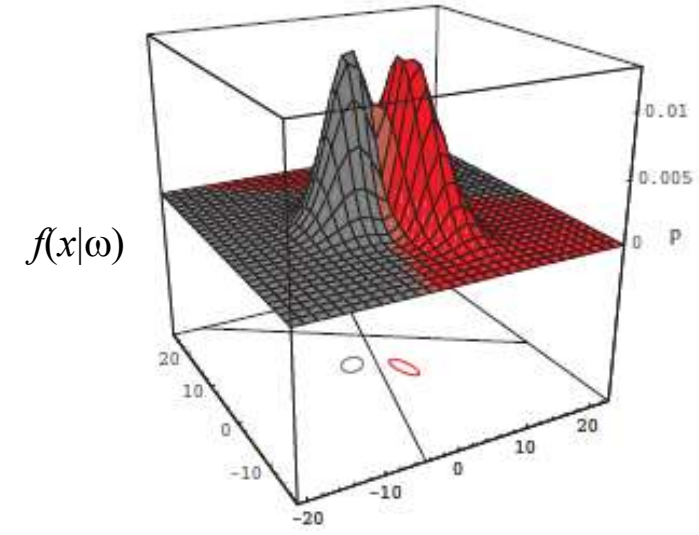
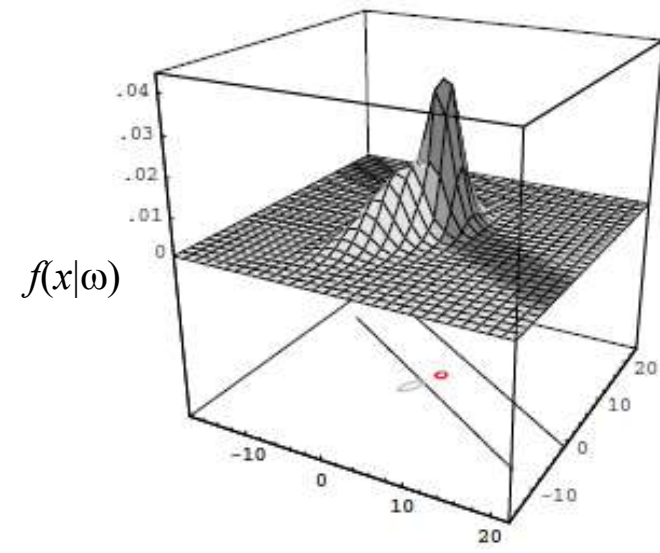
$$\mathbf{x}_0 = \frac{1}{2} (\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\ln(\Pr(\omega_i)) - \ln(\Pr(\omega_j))}{(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^T \mathbf{C}^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$$



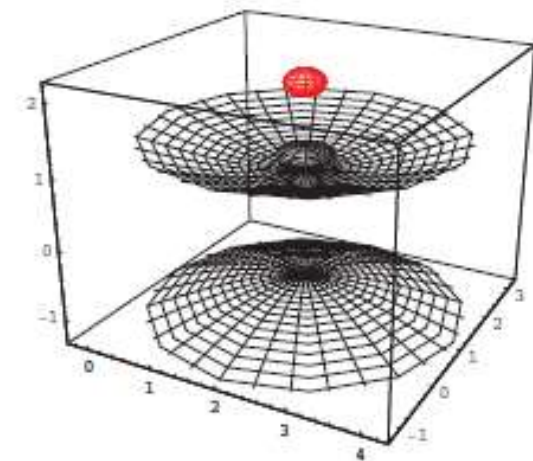
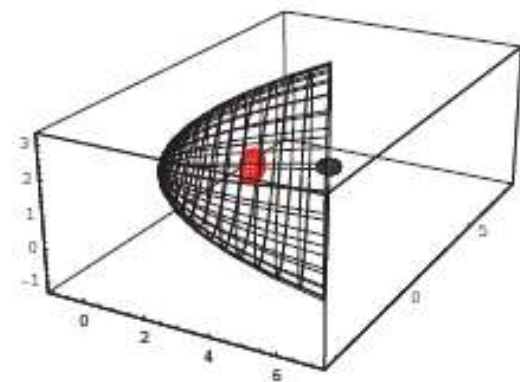
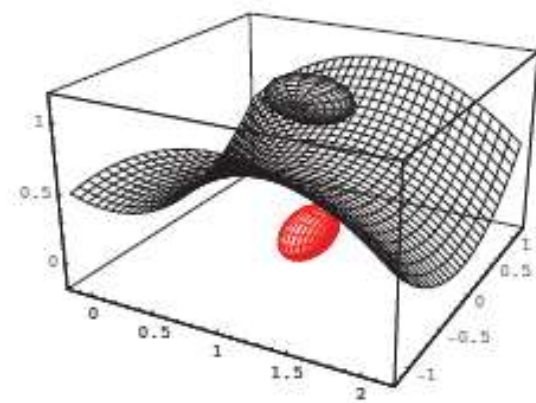
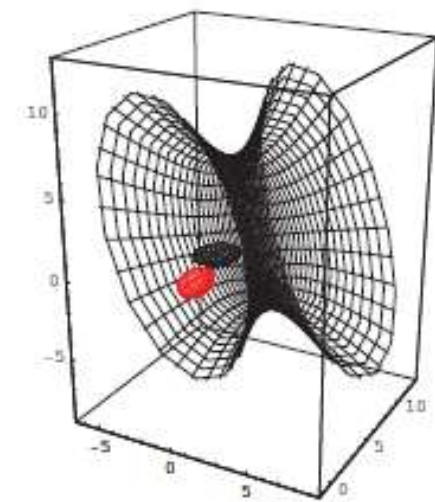
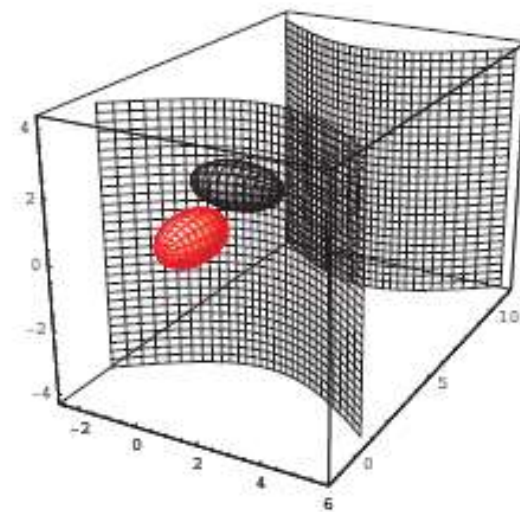
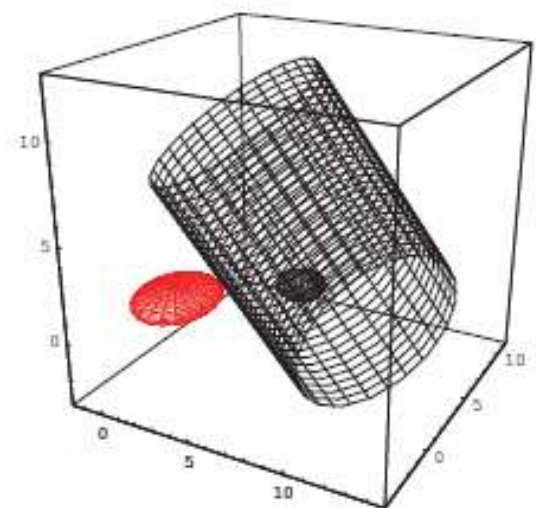
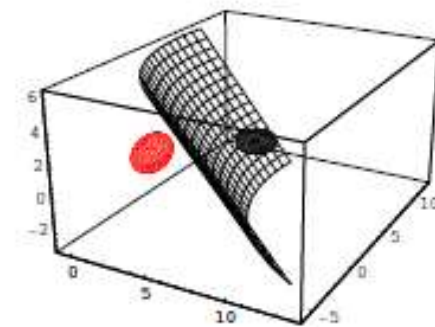
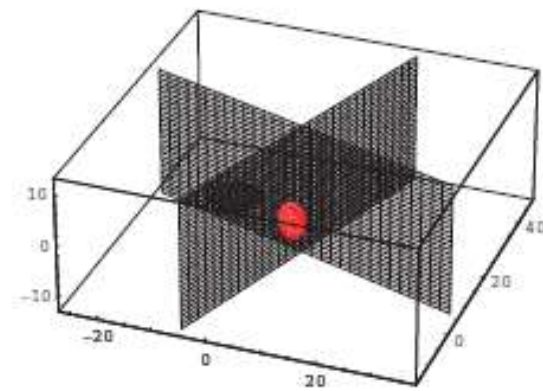
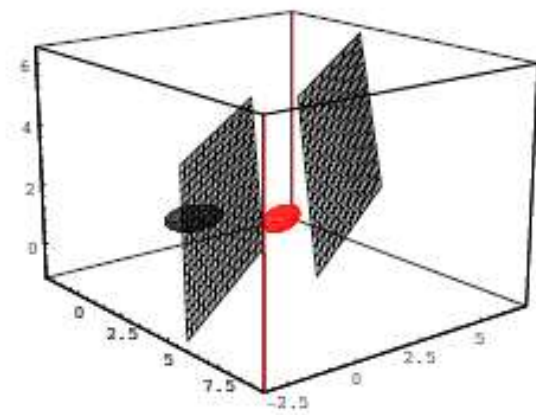
### Case 3 $\mathbf{C}_i$ arbitrary

$$h_i(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^T \mathbf{C}_i^{-1} \mathbf{x} + \boldsymbol{\mu}_i^T \mathbf{C}_i^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i^T \mathbf{C}_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\mathbf{C}_i| + \ln \Pr(\omega_i)$$

- Decision boundaries separating two classes are hiper-quadratic:
  - **Hyper-planes**
  - **Hyper-spheres**
  - **Hyper-ellipsoids**
  - **Hyper-paraboloids**
  - **Hyper-hyperboloids**









**Exercise.** Determine the surfaces separating 2 regions

$$h_i(\mathbf{x}) = h_j(\mathbf{x}) \Rightarrow$$

$$-\frac{1}{2} \mathbf{x}^T \mathbf{C}_i^{-1} \mathbf{x} + \boldsymbol{\mu}_i^T \mathbf{C}_i^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i^T \mathbf{C}_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\mathbf{C}_i| + \ln \Pr(\omega_i) \\ + \frac{1}{2} \mathbf{x}^T \mathbf{C}_j^{-1} \mathbf{x} - \boldsymbol{\mu}_j^T \mathbf{C}_j^{-1} \mathbf{x} + \frac{1}{2} \boldsymbol{\mu}_j^T \mathbf{C}_j^{-1} \boldsymbol{\mu}_j + \frac{1}{2} \ln |\mathbf{C}_j| - \ln \Pr(\omega_j) = 0$$

$\Rightarrow$

$$\mathbf{x}^T \left( \frac{1}{2} \mathbf{C}_j^{-1} - \frac{1}{2} \mathbf{C}_i^{-1} \right) \mathbf{x} + \left( \boldsymbol{\mu}_i^T \mathbf{C}_i^{-1} - \boldsymbol{\mu}_j^T \mathbf{C}_j^{-1} \right) \mathbf{x} \\ - \frac{1}{2} \boldsymbol{\mu}_i^T \mathbf{C}_i^{-1} \boldsymbol{\mu}_i + \frac{1}{2} \boldsymbol{\mu}_j^T \mathbf{C}_j^{-1} \boldsymbol{\mu}_j - \frac{1}{2} \ln \frac{|\mathbf{C}_i|}{|\mathbf{C}_j|} + \ln \frac{\Pr(\omega_i)}{\Pr(\omega_j)} = 0$$

$$h(\mathbf{x}) = h_i(\mathbf{x}) - h_j(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{v}^T \mathbf{x} + e = 0$$

Gaussian classes are not the only ones for which linear discriminants are optimal. Let us take the case of discrete random variables...

- Components of  $\mathbf{X}$ , are discrete

$$\mathbf{X} = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$$

- For  $c=2$  categories and vectors of dimension  $d$

$$p_i = \Pr(x_i = 1 | \omega_1) = 1 - \Pr(x_i = 0 | \omega_1)$$

$$q_i = \Pr(x_i = 1 | \omega_2) = 1 - \Pr(x_i = 0 | \omega_2)$$

- If features are statistically independent

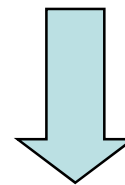
$$f(\mathbf{x} | \omega_1) = \prod_{i=1}^d \Pr(x_i | \omega_1) = \prod_{i=1}^d (p_i)^{x_i} (1 - p_i)^{1-x_i}$$

$$f(\mathbf{x} | \omega_2) = \prod_{i=1}^d \Pr(x_i | \omega_2) = \prod_{i=1}^d (q_i)^{x_i} (1 - q_i)^{1-x_i}$$



**Exercise.** Evaluate the discriminant and prove linearity in  $x_i$

$$h(\mathbf{x}) \equiv \ln \Pr(\omega_1 | \mathbf{x}) - \ln \Pr(\omega_2 | \mathbf{x}) \underset{\omega_2}{\overset{\omega_1}{\gtrless}} 0$$



$$h(\mathbf{x}) = \sum_{i=1}^d \left( x_i \ln \frac{p_i}{q_i} + (1 - x_i) \ln \frac{1-p_i}{1-q_i} \right) + \ln \frac{\Pr(\omega_1)}{\Pr(\omega_2)}$$



## 1.7 PERFORMANCE INDICATORS

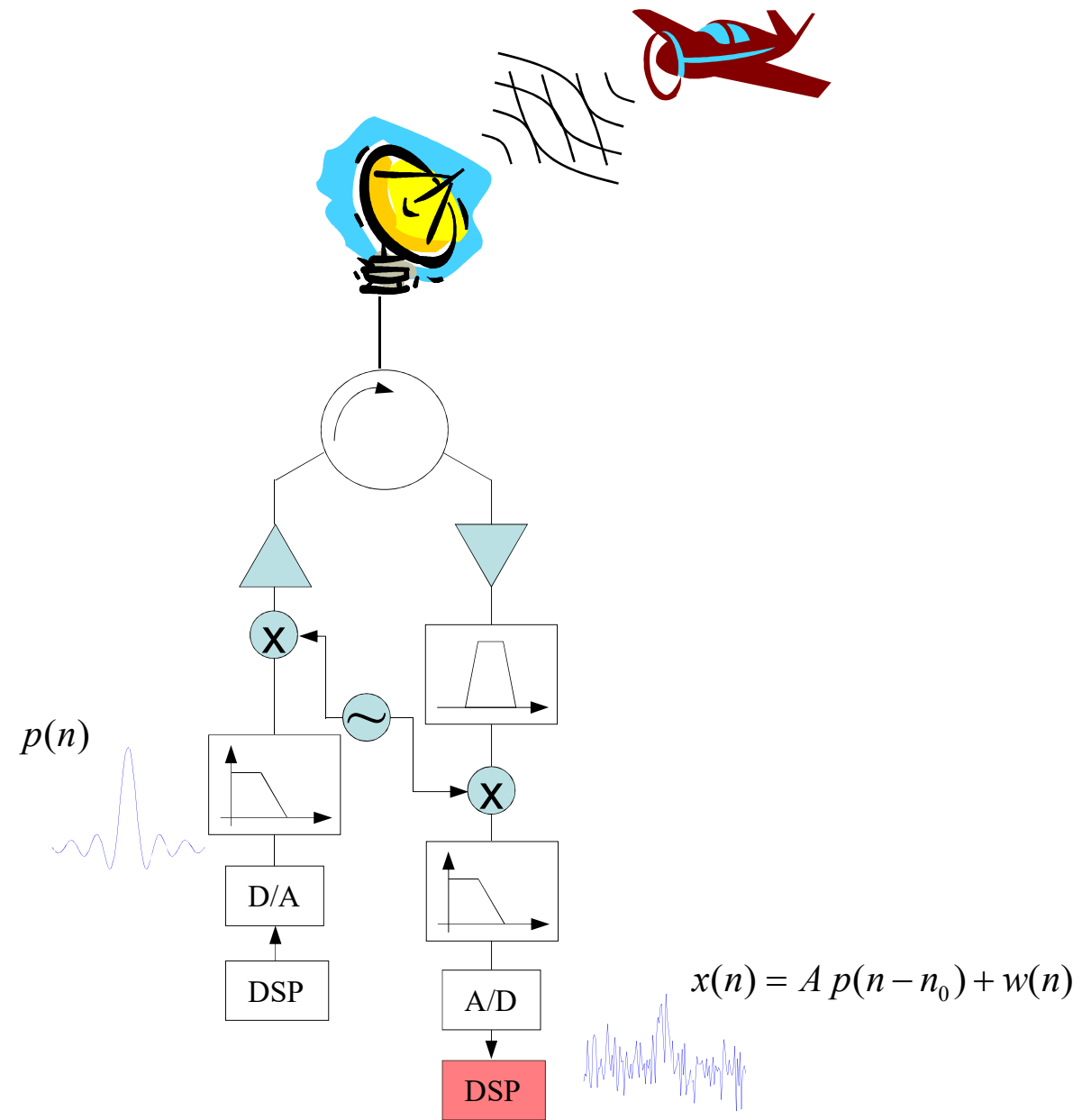
- Assume scalar values for  $\mathbf{x}$  and two classes ( $d = 1, c = 2$ )

$$f_x(x|\omega_1) \sim N(\mu_1, \sigma^2)$$

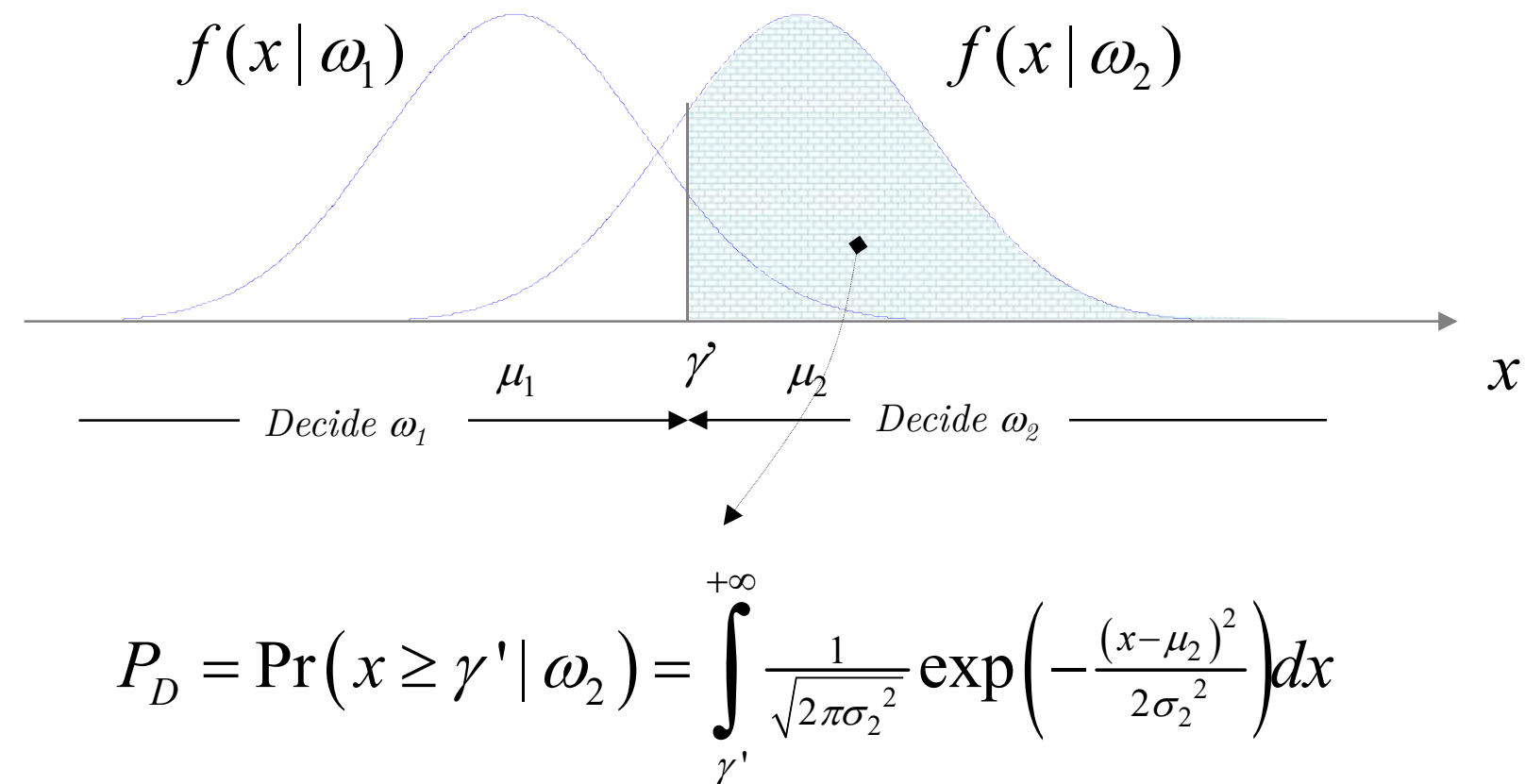
$$f_x(x|\omega_2) \sim N(\mu_2, \sigma^2)$$

- The classifier uses a threshold  $\gamma'$  that defines decision regions
- The performance of the classifier can be measured in terms of 4 probabilities (in radar wording):

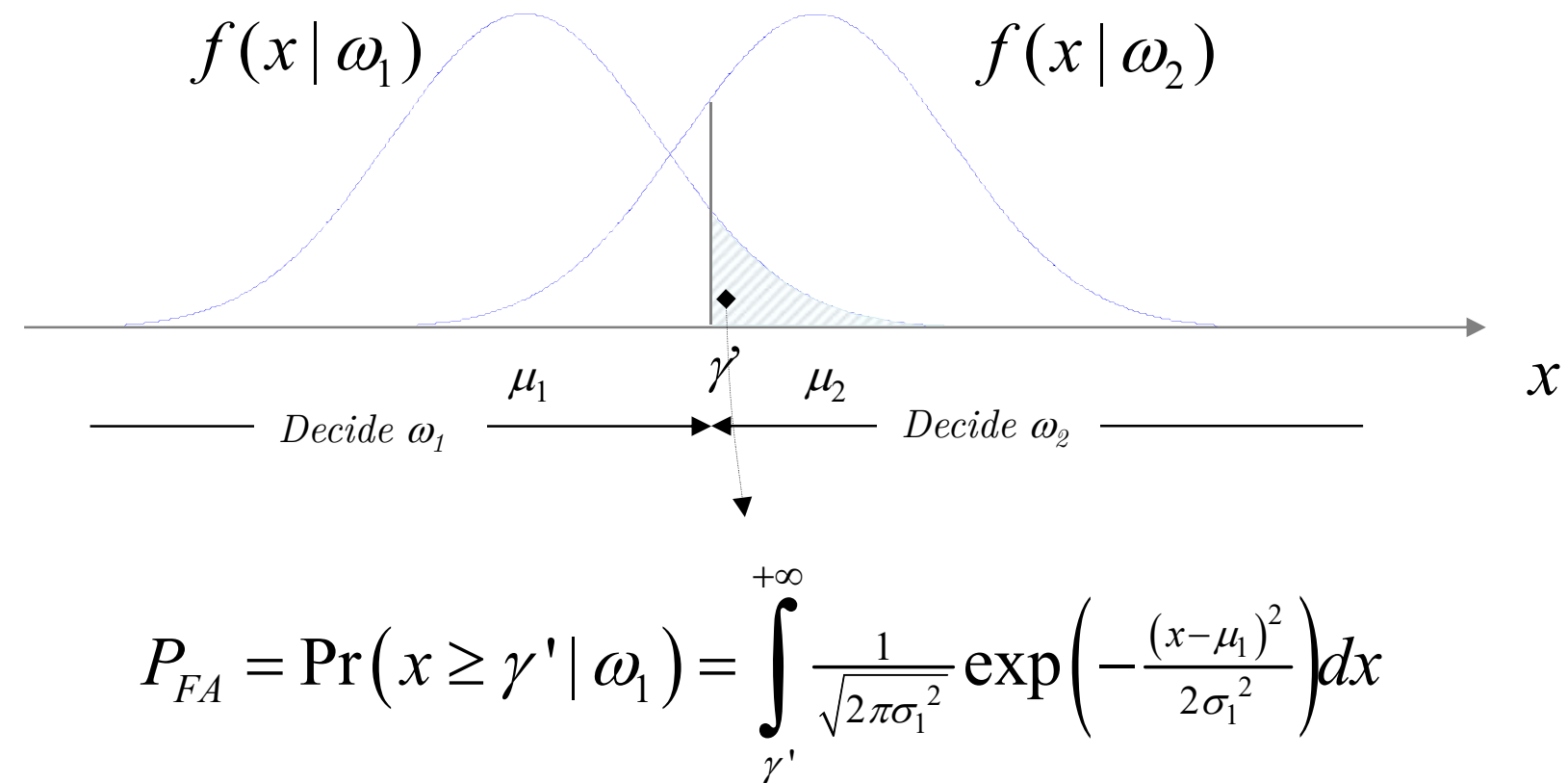
- |                  |                                  |  |
|------------------|----------------------------------|--|
| – Detection      | $\Pr(x \geq \gamma'   \omega_2)$ | Sensitivity of the detector              |
| – False alarm    | $\Pr(x \geq \gamma'   \omega_1)$ |  |
| – Miss           | $\Pr(x \leq \gamma'   \omega_2)$ |  |
| – Correct reject | $\Pr(x \leq \gamma'   \omega_1)$ | Specificity (or power) of the classifier |



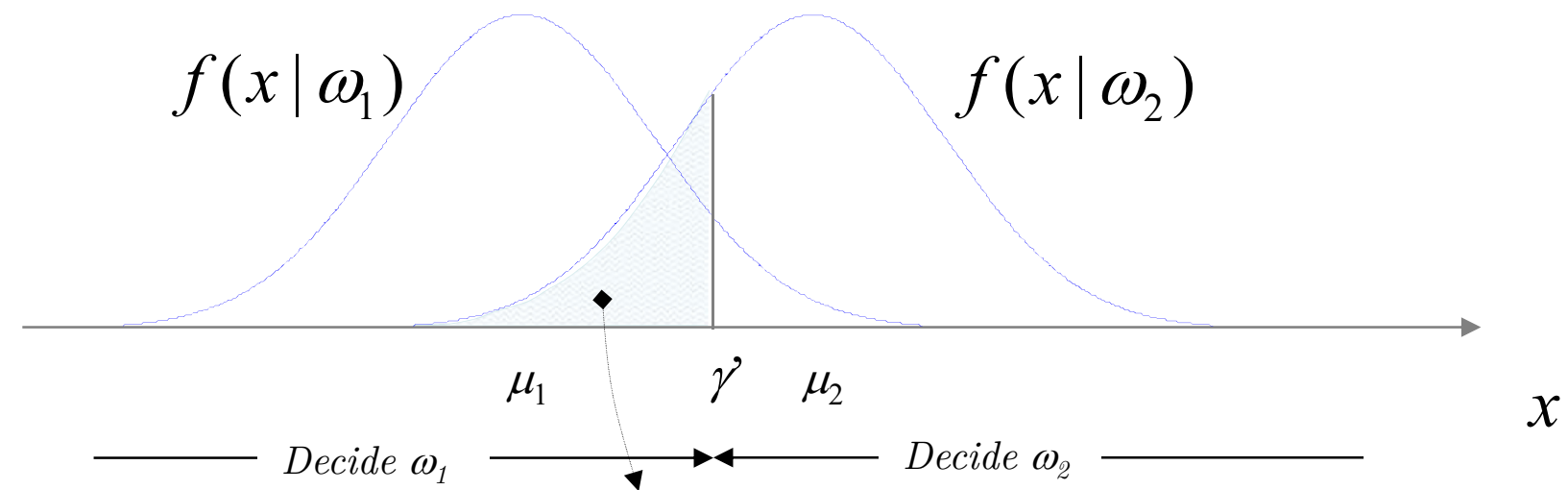
Probability of detection (sensitivity), for Gaussian model:



Probability of false alarm, for Gaussian model :

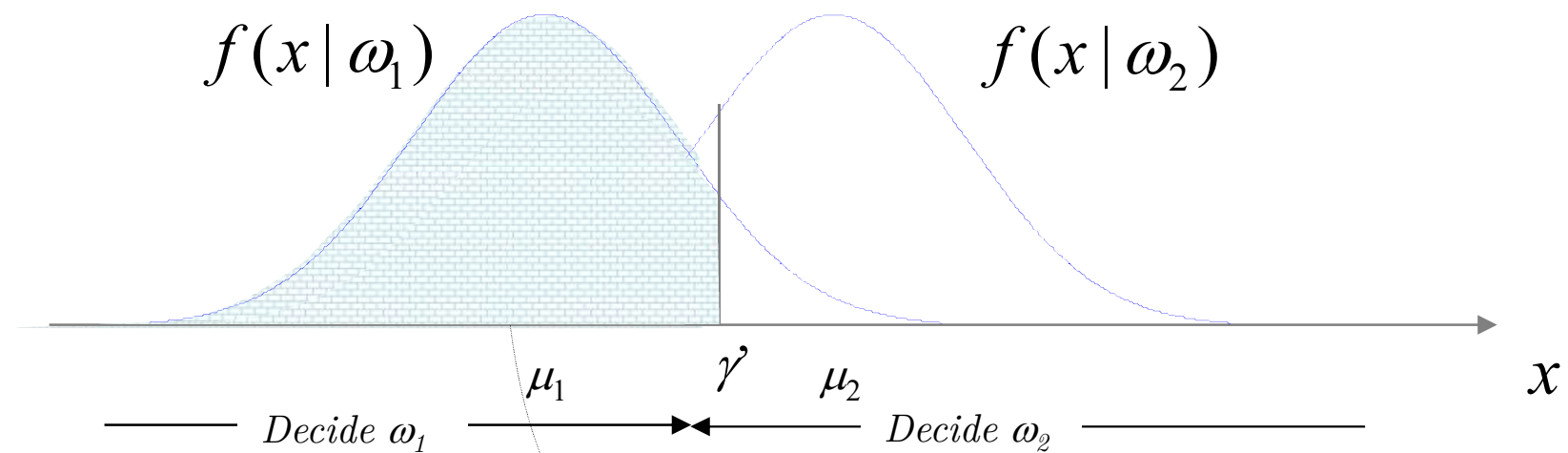


Probability of miss, for Gaussian model:



$$P_p = \Pr(x \leq \gamma' | \omega_2) = \int_{-\infty}^{\gamma'} \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(x-\mu_2)^2}{2\sigma_2^2}\right) dx$$

Probability of correct reject (specificity), for Gaussian model:



$$P_R = \Pr(x \leq \gamma' | \omega_1) = \int_{-\infty}^{\gamma'} \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(x-\mu_1)^2}{2\sigma_1^2}\right) dx$$

# The receiver operating characteristic (ROC)

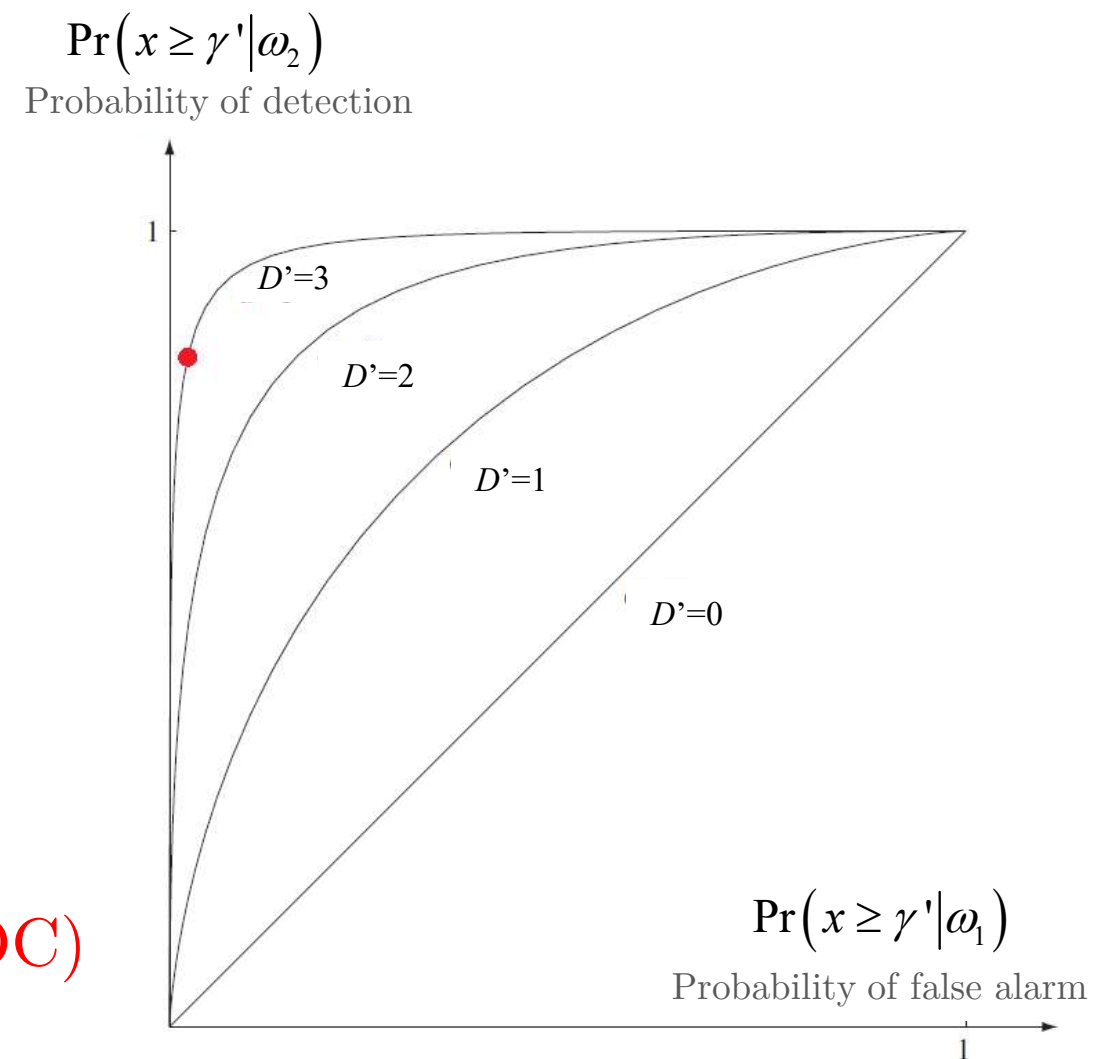
- ROC represents probability of detection vs. probability of false alarm.

- It depends on the discriminability among classes.

- A measure of discriminability, if variances for each class are equal:

$$D' = \frac{|\mu_2 - \mu_1|}{\sigma}$$

- The **area under the ROC (AUROC)** is often used as a measure of discriminability (or classifier performance).



- For more than two classes ( $c > 2$ ), for a given probability of detection there are several possible values for the false alarm probability.
- A simple extension of discriminability for  $d > 1$

$$D(\omega_i, \omega_j) = \sum_{k=1}^d \left| \frac{(\boldsymbol{\mu}_i)_k}{(\boldsymbol{\sigma}_i)_k} - \frac{(\boldsymbol{\mu}_j)_k}{(\boldsymbol{\sigma}_j)_k} \right|$$

- Discriminability based on the Mahalanobis distance:

$$D_M(\mathbf{x} \in \omega_i, \omega_j) = \sum_{\mathbf{x} \in \omega_i} (\mathbf{x} - \boldsymbol{\mu}_j) \mathbf{C}_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j)$$



## Error metrics in more common wording...

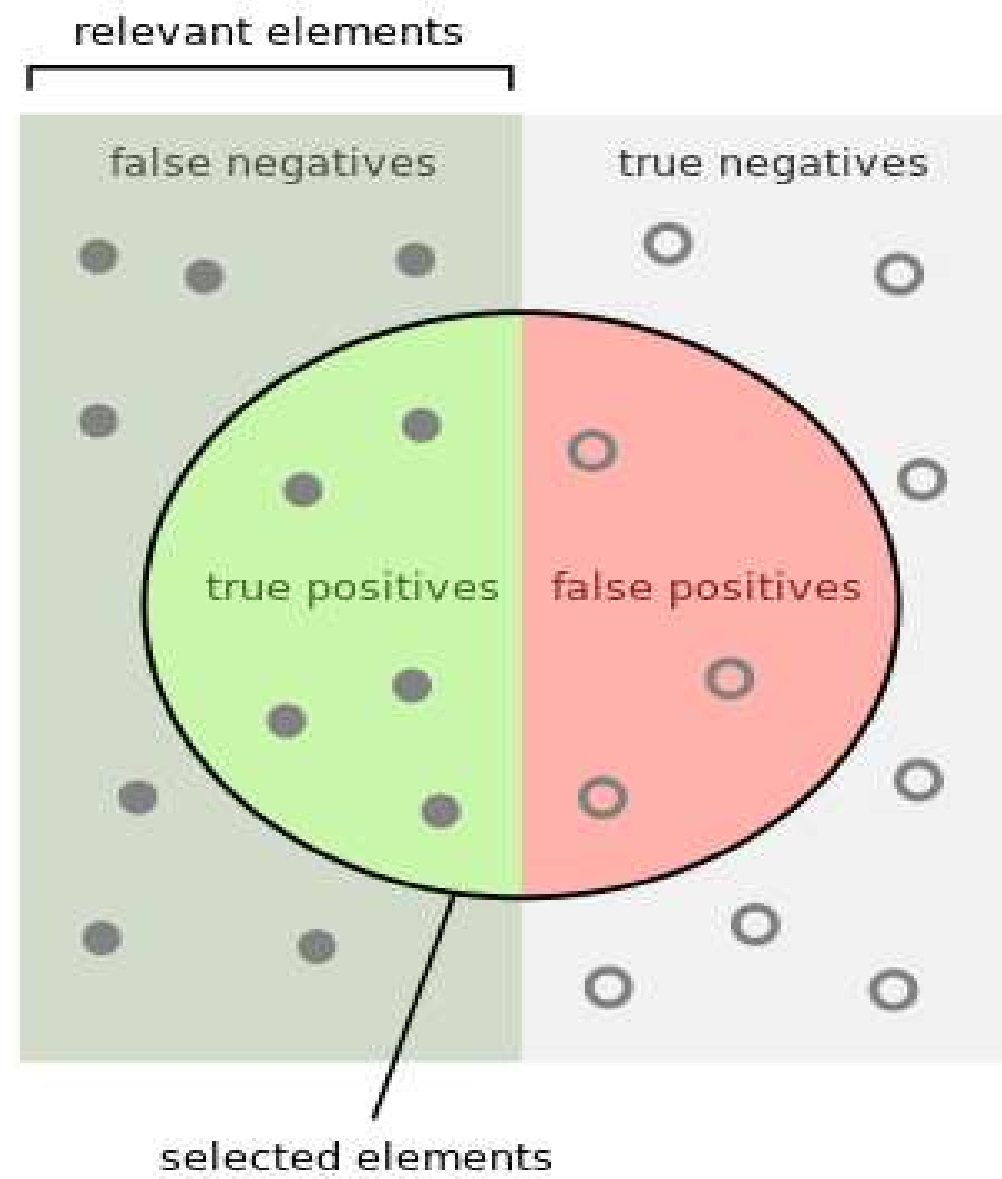
<u>True positive</u> : Sick people correctly diagnosed as sick	$\Pr(\mathbf{x} \in R_2   \omega_2)$
<u>False positive</u> : Healthy people incorrectly identified as sick	$\Pr(\mathbf{x} \in R_2   \omega_1)$
<u>True negative</u> : Healthy people correctly identified as healthy	$\Pr(\mathbf{x} \in R_1   \omega_1)$
<u>False negative</u> : Sick people incorrectly identified as healthy	$\Pr(\mathbf{x} \in R_1   \omega_2)$

- Sensitivity (the true positive probability, or “recall”,  $R$ ):  
Probability that a sick patient be diagnosed as sick (true positive rate)
- Specificity: (the true negative probability,  $S$ ): Probability that a non-infected patient be diagnosed as healthy (true negative rate)
- ➡ • Precision (also “positive predictive value”, or  $P$ ): Out of all patients decided sick, what fraction is actually sick?  $\Pr(\omega_2 | \mathbf{x} \in R_2)$
- ➡ • F-score: It measures the test accuracy considering both the precision and the recall

$$F_{score} = \frac{2P \cdot R}{P + R}$$

- ➡ • Accuracy: fraction of total correct decisions for all classes.

$$A = \Pr(\omega_1)S + \Pr(\omega_2)R$$



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Different combinations of precision (P) and recall (R) have the following meanings:

- **high recall + high precision** : class  $\omega_2$  is perfectly handled by the model
- **low recall + high precision** : the model cannot detect class  $\omega_2$  well but is highly trustable when it does
- **high recall + low precision** : class  $\omega_2$  is well detected but the model also include points of classes  $\omega_1$  in it
- **low recall + low precision** : class  $\omega_2$  is poorly handled by the model

# The confusion matrix

		Predicted Value		
		$\omega_2$	$\omega_1$	
Actual Value	$\omega_2$	16 True Positives	30 False Negatives	Recall
	$\omega_1$	10 False Positives	144 True Negatives	Specificity
		Prevalence	Precision	Negative Predictive Value
				Accuracy

Sensitivity

Precision



It is crucial to distinguish between  $\Pr(\mathbf{x} \in R_2 | \omega_2)$  and  $\Pr(\omega_2 | \mathbf{x} \in R_2)$   
and between  $\Pr(\mathbf{x} \in R_1 | \omega_1)$  and  $\Pr(\omega_1 | \mathbf{x} \in R_1)$

Specificity

**How are they related?**

When getting the result of a medical test on a disease, what can be said about the patient's condition? Take as random variables:

Patient condition:  $\omega \in \{\text{healthy}, \text{sick}\}$

Result of the test:  $y \in \{-, +\}$

and the prevalence of the disease:  $\Pr(\text{sick})$



Use Bayes' theorem to determine the chances you are sick given that the result of the test is positive:

$$\Pr(sick|+) = \frac{\Pr(+|sick)\Pr(sick)}{\Pr(+|healthy)\Pr(healthy) + \Pr(+|sick)\Pr(sick)}$$

Prob. of getting a positive test if you are sick

Prevalence



Clinical trials at the pharmaceutical lab...

$\omega$  : healthy



$y$  : - - - - - + - - - - -

$\Pr(- | \text{healthy})$   
Specificity

$\omega$  : sick



$y$  : + - + + + + + + -

$\Pr(+ | \text{sick})$   
Sensitivity

Me at the doctor's...



$y$  : +

$\Pr(\text{sick} | +)$  ?

Get values of specificity and sensitivity for AIDS test from this [link](#),  
and values of prevalence from [AIDS prevalence data in Spain](#).



Using Bayes' theorem

$$\Pr(sick|+) = \frac{\overset{0.997}{\Pr(+|sick)} \overset{0.001}{\Pr(sick)}}{\Pr(+|healthy) \Pr(healthy) + \Pr(+|sick) \Pr(sick)} = 0.062$$

Chances are very low!

$$\uparrow (1-0.985) \times 0.999 + 0.997 \times 0.001$$

What are the chances if you belong to a risk population from [AIDS prevalence data in Spain](#)?

Now, what are the chances of being sick if a second independent test is positive too,  $+_2$ ? Let us apply Bayes' theorem again...





... but note that the prior is now  $\Pr(sick|+_1) = 0.062$ , as obtained from the first test:

$$\begin{aligned}\Pr(sick|+_1, +_2) &= \frac{\Pr(+_2|sick, +_1)\Pr(sick|+_1)}{\Pr(+_2|healthy, +_1)\Pr(healthy|+_1) + \Pr(+_2|sick, +_1)\Pr(sick|+_1)} = \\ &= \frac{\Pr(+_2|sick)\Pr(sick|+_1)}{\Pr(+_2|healthy)\Pr(healthy|+_1) + \Pr(+_2|sick)\Pr(sick|+_1)} = 0.82\end{aligned}$$

Now the probability is much higher!

## 1.8 CONCLUSIONS

1. Classification is done by maximising a discriminant function:

$$\hat{\omega}_i = \max_i \{g_i(\mathbf{x})\} \quad i = 1, \dots, c$$

In the MAP case  $g_i(\mathbf{x}) = \Pr(\omega_i | \mathbf{x})$

2. Linear boundaries are only optimal in some specific cases.
3. Clusters of vectors in the  $d$ -dimensional space identify classes. The shape of clusters depend on the eigenvalues of the covariance matrix.
4. The ROC is useful to measure the behaviour of the classifier. It can be used to compute the optimal threshold in a experimental way, from a labeled data base. The area under the ROC and Mahalanobis distance provide a measure of the discriminability of our problem.

# CONTENTS

## **2.1 Bayesian decision**

## **2.2. Maximum likelihood (ML) estimation and Bayesian estimation**

### 2.2.1 Introduction

### 2.2.2 ML estimation

### 2.2.3 Bayesian estimation

### 2.2.4 Conclusions

## 2.1 INTRODUCTION

Bayesian decision requires precise knowledge of  $f_{\mathbf{x}}(\mathbf{x} | \omega_i)$  and  $\Pr(\omega_i)$ . The computation of these magnitudes require:

- Having a previously-labeled reliable data base (train data base).
- Having an estimator of the pdf and the a priori probabilities.

The estimation of  $f_{\mathbf{x}}(\mathbf{x} | \omega_i)$  requires many data unless we can use of function that depends on a few parameters  $\theta_i$ .

**Gaussian case:**  $\theta_i$  contains the mean and the covariance matrix

$$f_{\mathbf{x}}(\mathbf{x} | \omega_i, \theta_i) = \frac{1}{(2\pi)^{d/2} |\mathbf{C}_i|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \mathbf{C}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right\}$$

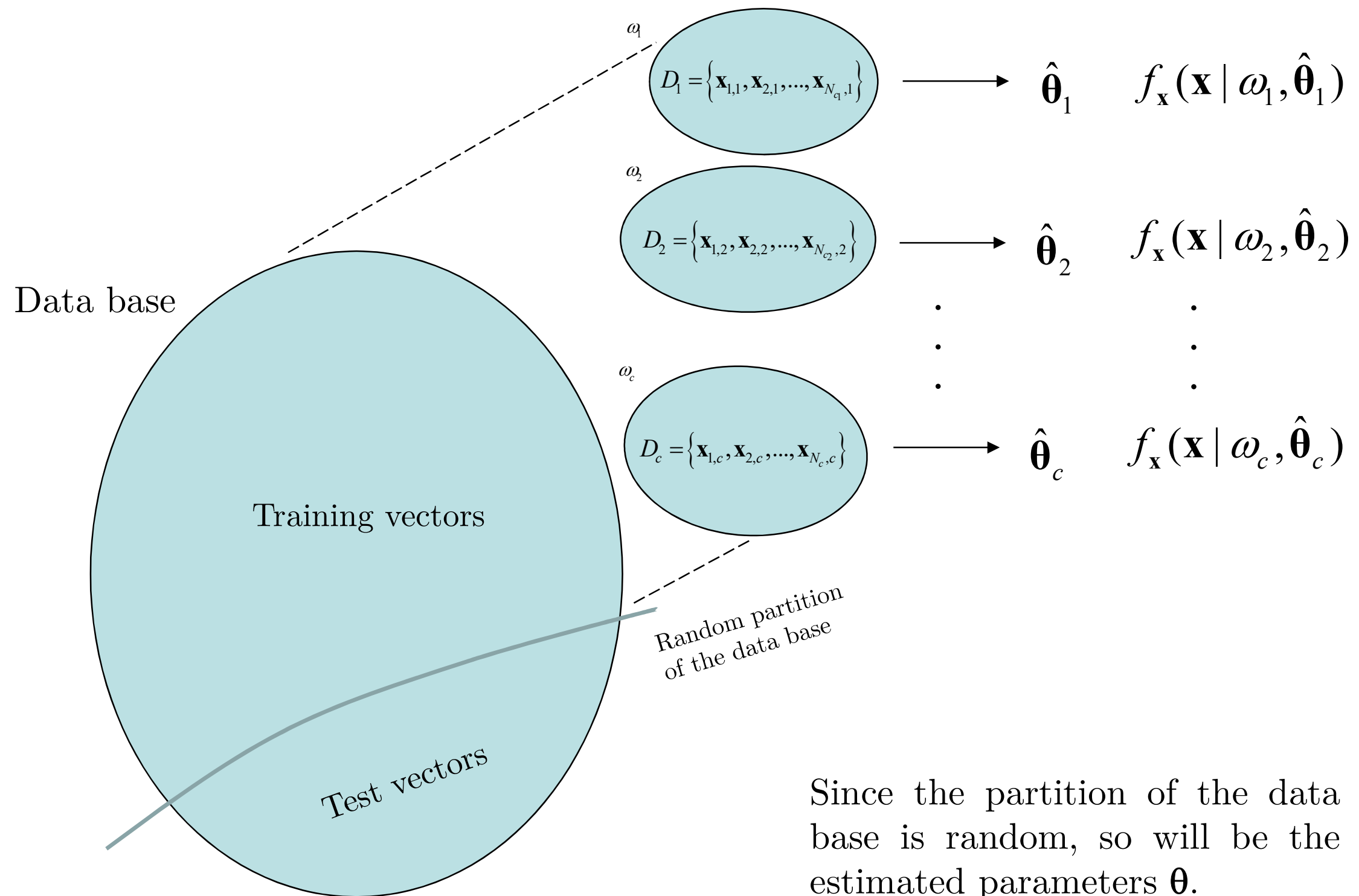
We have two possibilities:

1. **Maximum likelihood estimation (ML).** The parameters are considered deterministic (although unknown).
2. **Bayesian estimation.** The parameters are random variables of which an a-priori knowledge is available (related to the concept of “belief”) in the form of a pdf.

In all cases, we will assume that a labeled data base is available. Using a partition (the **training data base**) we have to determine  $f_{\mathbf{x}}(\mathbf{x}|\omega_i)$ .

The rest of vectors will be used to evaluate the performance of the classifier (the **test data base**).

See over...



## 2.2 MAXIMUM LIKELIHOOD ESTIMATION (ML)

Assume that, on each class  $i$ , the observed vectors  $\mathbf{x}_{k,i} \in D_i$  are statistically independent. The likelihood function is given by:

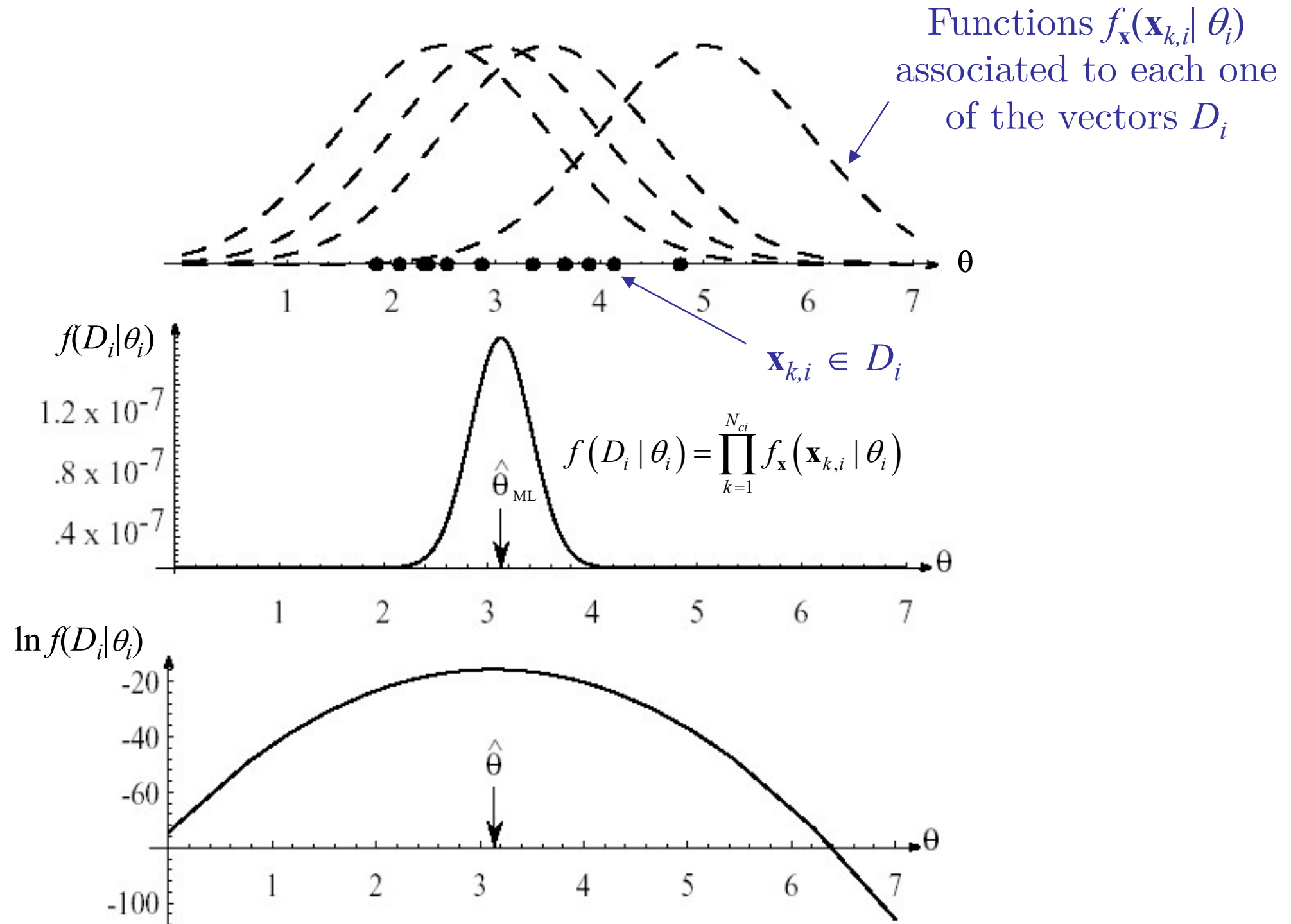
$$f(D_i | \boldsymbol{\theta}_i) = \prod_{k=1}^{N_{c_i}} f_{\mathbf{x}}(\mathbf{x}_{k,i} | \boldsymbol{\theta}_i)$$

The ML estimator maximizes this function (or a non-decreasing function of it):

$$\hat{\boldsymbol{\theta}}_{i,ML} = \arg \max_{\boldsymbol{\theta}_i} f(D_i | \boldsymbol{\theta}_i) = \arg \max_{\boldsymbol{\theta}_i} \ln f(D_i | \boldsymbol{\theta}_i)$$

A necessary condition to obtain the estimator is therefore given by:

$$\nabla_{\boldsymbol{\theta}_i} \ln f(D_i | \boldsymbol{\theta}_i) = \mathbf{0}$$





## Characterization of an estimator

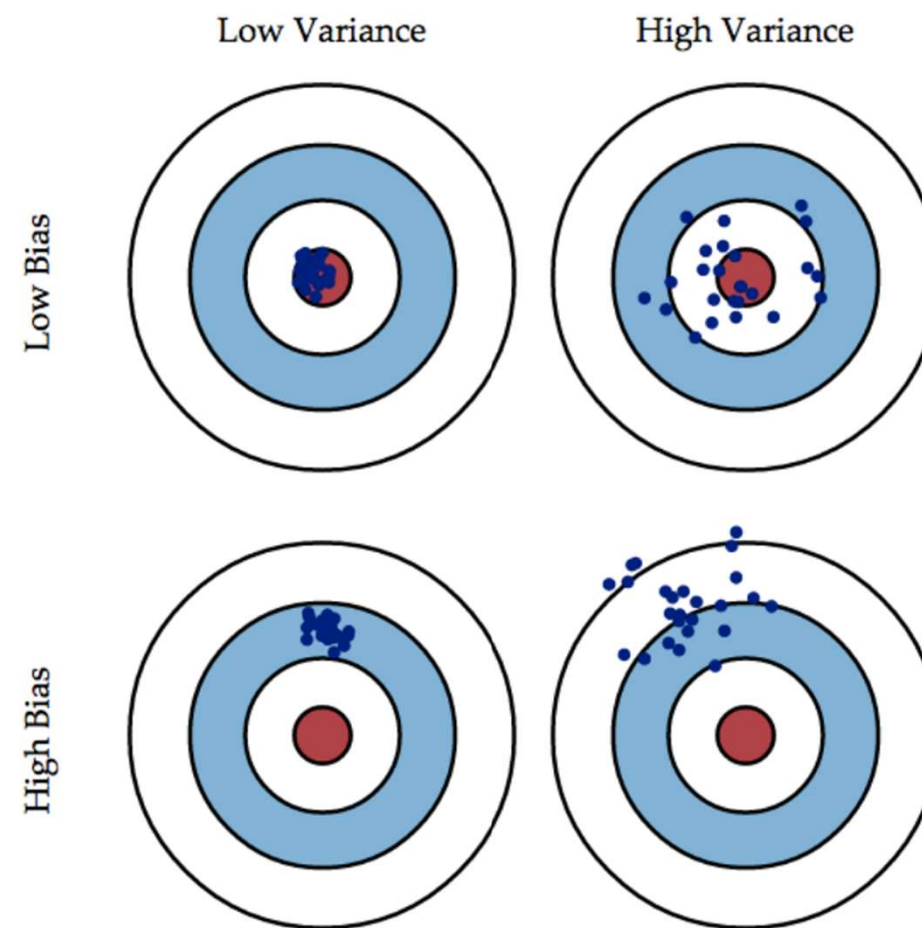
An estimator is a function that applies on an ensemble of feature vectors  $\mathbf{x}_{k,i}$  from the training data base. If the selection of vectors is random, so will be the outcomes of the estimator: for each possible partition  $l$  of the data base, we obtain a different estimate  $\hat{\boldsymbol{\theta}}_{l,i}$

1. **Bias:** is the difference between the true value of the parameter and the average of all possible estimates obtained from all possible random partitions of the data base. It measures the systematic error of the estimator.

$$B\{\hat{\boldsymbol{\theta}}_i\} = \boldsymbol{\theta} - \frac{1}{L} \sum_{l=1}^L \hat{\boldsymbol{\theta}}_{l,i}$$

**2. Variance:** is the deviation of the estimated values from the average value. It measures how the outcome of the estimator depends on a specific partition of the data base. For a scalar parameter:

$$\text{var}\{\hat{\theta}_i\} = \frac{1}{L} \sum_{l=1}^L \left( \hat{\theta}_{l,i} - \frac{1}{L} \sum_{s=1}^L \hat{\theta}_{s,i} \right)^2$$



## Properties of the ML estimator:

1. It is asymptotically unbiased (in many cases it is unbiased for low value of  $N$ )
2. It is asymptotically efficient (for large  $N$ , its variance attains the minimum variance given by the Crámer-Rao bound)

## However...

1. It does not necessarily provides the least misclassification error if used in

$$f_{\mathbf{x}}(\mathbf{x} | \omega_i, \hat{\boldsymbol{\theta}}_{i,ML})$$

2. If the assumed pdf is far from the real one, the estimations may be of very low quality.



### Example 1:

ML estimator of the mean  $\boldsymbol{\mu}_i$  if the covariance matrix  $\mathbf{C}_i$  is known, in the multivariate Gaussian case. Prove that:

$$\hat{\boldsymbol{\mu}}_{i,ML} = \frac{1}{N_{c_i}} \sum_{k=1}^{N_{c_i}} \mathbf{x}_k$$

### Example 2:

ML estimator of both the mean  $\boldsymbol{\mu}_i$  and the covariance matrix  $\mathbf{C}_i$  in the multivariate Gaussian case. Prove that:

$$\hat{\boldsymbol{\mu}}_{i,ML} = \frac{1}{N_{c_i}} \sum_{k=1}^{N_{c_i}} \mathbf{x}_k \quad \hat{\mathbf{C}}_{i,ML} = \frac{1}{N_{c_i}} \sum_{k=1}^{N_{c_i}} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_{i,ML})(\mathbf{x}_k - \hat{\boldsymbol{\mu}}_{i,ML})^T$$



### Example 3:

ML estimator of the probability  $p_k$  for ‘1’ in each component of a binary-valued vector  $\mathbf{x} \in \{0,1\}^d$ :

$$f_{\mathbf{x}}(D \mid \omega, \mathbf{p}) = \prod_{j=1}^{N_i} \prod_{k=1}^d p_k^{x_{k,j}} (1 - p_k)^{1-x_{k,j}}$$

$$\mathbf{p} = [p_1, \dots, p_d]$$

## 2.3 BAYESIAN ESTIMATION

Sometimes we can take advantage of a priori knowledge about the possible values of  $\theta_i$ . This knowledge will be included in  $f(\theta_i)$ , which has all properties of a pdf and expresses our “belief” about the possible values of  $\theta_i$ . Two approaches are possible:

1. Improve the ML estimation of  $\theta_i$  (using MAP principles)

$$\hat{\theta}_{i,MAP} = \arg \max_{\theta_i} f(D_i | \theta_i) f(\theta_i) = \arg \max_{\theta_i} [\ln f(D_i | \theta_i) + \ln f(\theta_i)]$$

2. Directly estimate the a posteriori probabilities  $\Pr(\omega_i|\mathbf{x})$

Computing  $f_{\mathbf{x}}(\mathbf{x}|\omega_i)$  and  $\Pr(\omega_i)$ . This is the recommended procedure in a classification application.

## ML AND BAYESIAN ESTIMATION

### Comparison:

The function  $f(D_i | \boldsymbol{\theta}_i)$  will peak around  $\boldsymbol{\theta}_i = \hat{\boldsymbol{\theta}}_i$ , the more as larger is  $N_i$ .

If  $f(\boldsymbol{\theta}_i)$  is non-zero and it does not change much around  $\boldsymbol{\theta}_i = \hat{\boldsymbol{\theta}}_i$ , then

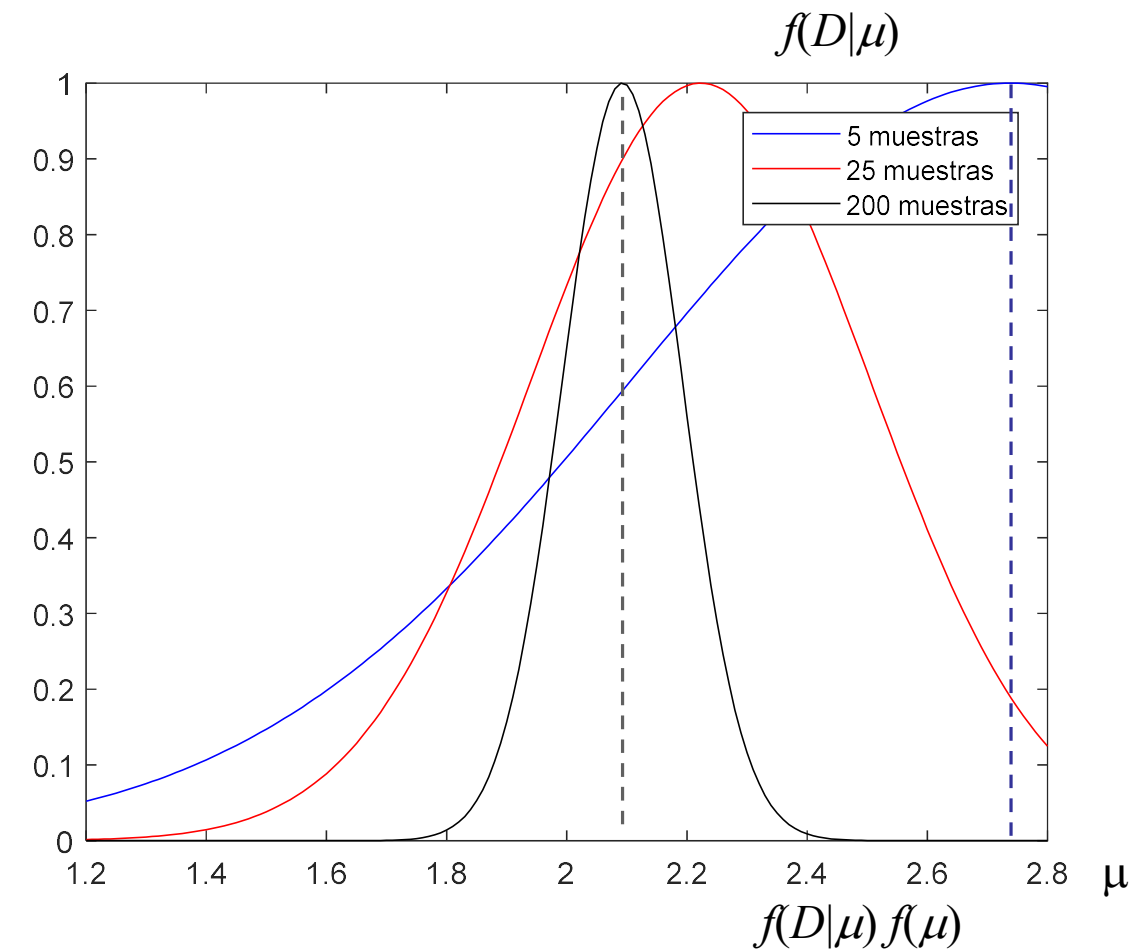
$$f(\boldsymbol{\theta}_i | D_i) = \frac{f(D_i | \boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i)}{f(D_i)}$$

also peaks in  $\boldsymbol{\theta}_i = \hat{\boldsymbol{\theta}}_i$  and the estimates obtained through ML and Bayesian principles coincide.

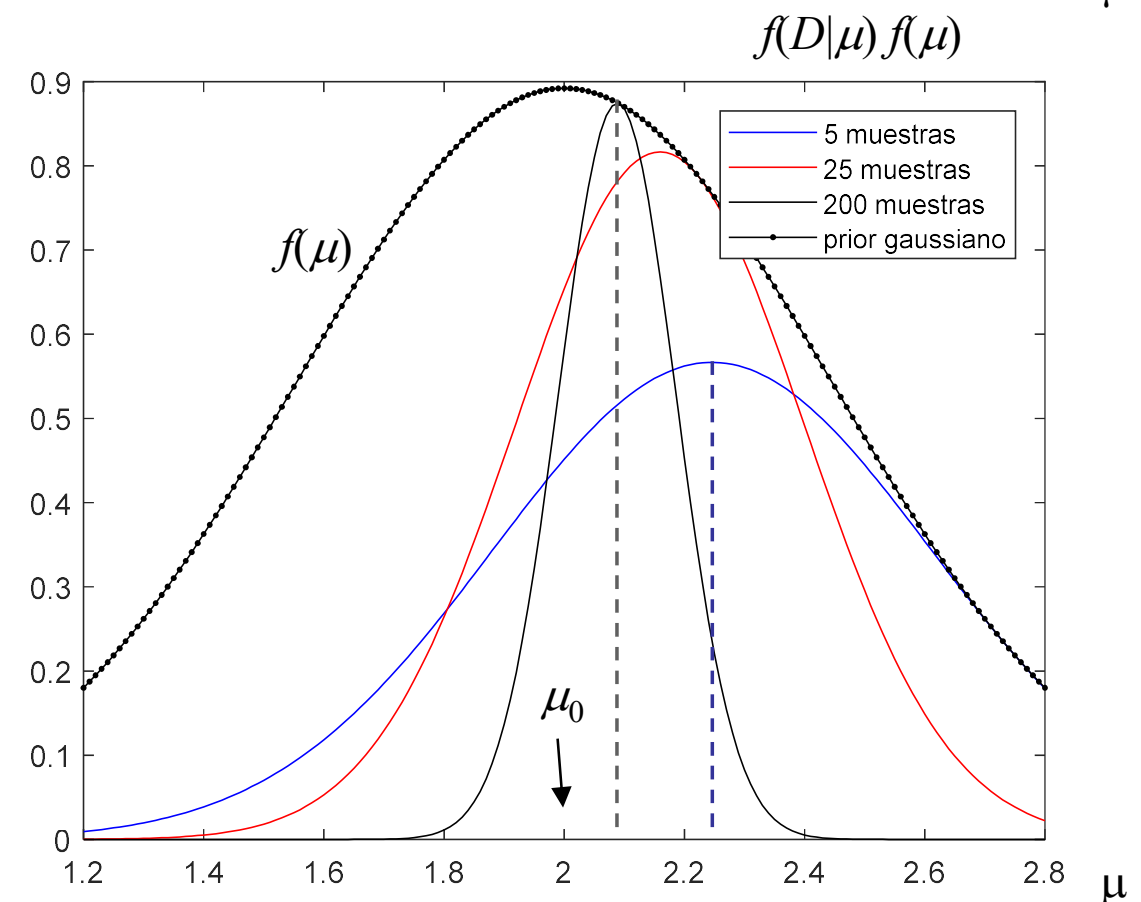
In practice, if the number of vectors in  $D_i$  is small, Bayesian estimate is preferred. When many vectors are available, both coincide...

### Example 5:

**ML** estimation of the mean ( $\mu_0=2$ ) obtained from a number of Gaussian samples (likelihood functions have been normalised to the máximo for clarity)

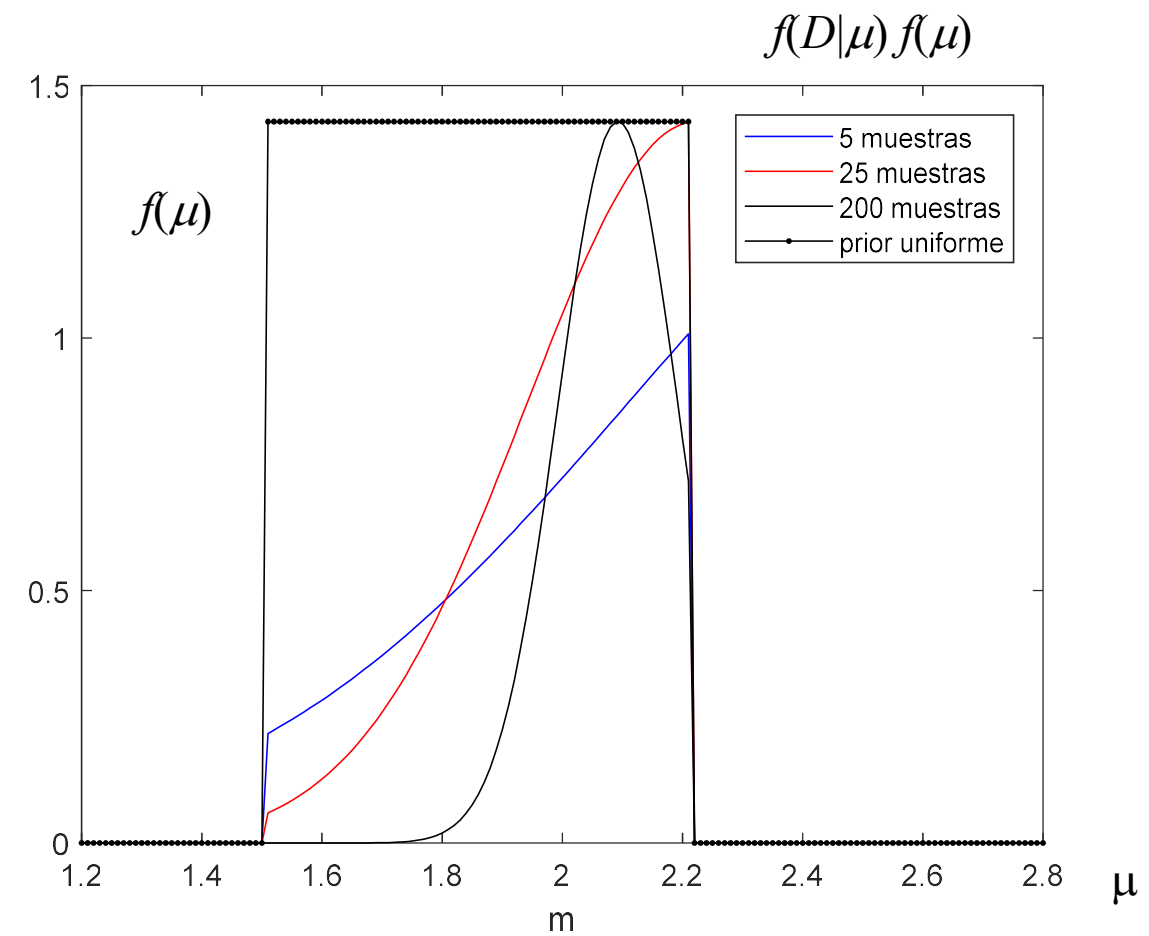


**Bayesian** estimate of the mean ( $\mu_0=2$ ) obtained from a number of Gaussian samples.  
The a priori pdf of  $\mu$  is Gaussian.





**Bayesian** estimate of the mean ( $\mu_0=2$ ) obtained from a number of Gaussian samples.  
The a priori pdf of  $\mu$  is uniform.





2. Directly estimate the a posteriori probabilities  $\Pr(\omega_i|\mathbf{x})$

### Assumptions

- The shape of  $f_{\mathbf{x}}(\mathbf{x}|\theta_i)$  is known, but not the parameter  $\theta_i$
- Our a priori knowledge of  $\theta_i$  is in  $f(\theta_i)$
- The rest of our knowledge on  $\theta_i$  is given by data in  $D_i$



## Procedure:

1. Average the likelihood function with respect to the a posteriori probability of our parameter:

$$f_{\mathbf{x}}(\mathbf{x} | \omega_i) \cong f_{\mathbf{x}}(\mathbf{x} | D_i) = \int f(\mathbf{x} | \boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i | D_i) d\boldsymbol{\theta}_i$$

2. Compute the a posteriori probability of our parameter as:

$$f(\boldsymbol{\theta}_i | D_i) = \frac{f(D_i | \boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i)}{\int f(D_i | \boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i} \propto f(D_i | \boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i)$$

3. Assuming independence of data in  $D_i$

$$f(D_i | \boldsymbol{\theta}_i) = \prod_{k=1}^{N_i} f(\mathbf{x}_{k,i} | \boldsymbol{\theta}_i)$$



#### Example 4:

Bayesian estimate of  $f_{\mathbf{x}}(\mathbf{x}|D)$  if

$$f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\mu}) = N(\boldsymbol{\mu}, \mathbf{C}) \qquad f(\boldsymbol{\mu}) = N(\boldsymbol{\mu}_0, \mathbf{C}_0)$$

where  $\boldsymbol{\mu}_0$ ,  $\mathbf{C}_0$  and  $\mathbf{C}$  are known, and feature vectors are available  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \omega$

From **2** and **3** we can write:

$$\begin{aligned} f(\boldsymbol{\mu} | D) &= \alpha \prod_{k=1}^N f_{\mathbf{x}}(\mathbf{x}_k | \boldsymbol{\mu}) f(\boldsymbol{\mu}) = \\ &= \alpha' \exp \left[ -\frac{1}{2} \left( \boldsymbol{\mu}^T (N\mathbf{C}^{-1} + \mathbf{C}_0^{-1}) \boldsymbol{\mu} + 2\boldsymbol{\mu}^T \left( \mathbf{C}^{-1} \sum_{k=1}^N \mathbf{x}_k + \mathbf{C}_0^{-1} \boldsymbol{\mu}_0 \right) \right) \right] \end{aligned}$$



This equation can also be written as:

$$f(\boldsymbol{\mu}|D) = \alpha'' \exp \left[ -\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_N)^T \mathbf{C}_N^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_N) \right]$$

And equating both expressions:

$$\boldsymbol{\mu}^T (N\mathbf{C}^{-1} + \mathbf{C}_0^{-1}) \boldsymbol{\mu} - 2\boldsymbol{\mu}^T \left( \mathbf{C}^{-1} \sum_{k=1}^N x_k + \mathbf{C}_0^{-1} \boldsymbol{\mu}_0 \right) = \boldsymbol{\mu}^T \mathbf{C}_N^{-1} \boldsymbol{\mu} - 2\boldsymbol{\mu}_N^T \mathbf{C}_N^{-1} \boldsymbol{\mu} + K$$

where the terms that do not depend on  $\boldsymbol{\mu}$  and other constants are lumped into  $K$ . Comparing the quadratic term in  $\boldsymbol{\mu}$ :

$$\mathbf{C}_N^{-1} = N\mathbf{C}^{-1} + \mathbf{C}_0^{-1} \quad (1)$$

And comparing the linear terms in  $\boldsymbol{\mu}$ :

$$\mathbf{C}_N^{-1} \boldsymbol{\mu}_N = \mathbf{C}^{-1} \sum_{k=1}^N x_k + \mathbf{C}_0^{-1} \boldsymbol{\mu}_0 \quad (2)$$



From (1) y using the equality:  $\left(\mathbf{A}^{-1} + \mathbf{B}^{-1}\right)^{-1} = \mathbf{A}(\mathbf{A} + \mathbf{B})^{-1} \mathbf{B}$

$$\mathbf{C}_N = \mathbf{C}_0 \left( \mathbf{C} + N\mathbf{C}_0 \right)^{-1} \mathbf{C} \quad (3)$$

Having in mind that if  $\mathbf{A}$  y  $\mathbf{B}$  are invertible

$$\mathbf{A}(\mathbf{A} + \mathbf{B})^{-1} \mathbf{B} = \mathbf{B}(\mathbf{A} + \mathbf{B})^{-1} \mathbf{A}$$

we can use (3) in (2) to obtain

$$\boldsymbol{\mu}_N = \mathbf{C}_0 \left( \mathbf{C}_0 + \frac{1}{N} \mathbf{C} \right)^{-1} \mathbf{m}_N + \frac{1}{N} \mathbf{C} \left( \mathbf{C}_0 + \frac{1}{N} \mathbf{C} \right)^{-1} \boldsymbol{\mu}_0$$

$$\mathbf{m}_N = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k$$



Note that the estimated mean is a linear combination of a priori knowledge of the mean  $\boldsymbol{\mu}_0$  and the information provided by the data  $\mathbf{m}_N$ . Integrating equation (1):

$$f_{\mathbf{x}}(\mathbf{x} | \omega) \cong f_{\mathbf{x}}(\mathbf{x} | D) = \int f(\mathbf{x} | \boldsymbol{\mu}) f(\boldsymbol{\mu} | D) d\boldsymbol{\mu} \sim N(\boldsymbol{\mu}_N, \mathbf{C} + \mathbf{C}_N)$$

When  $N \rightarrow \infty$  the estimation of  $\boldsymbol{\mu}$  from  $f(\boldsymbol{\mu}|D)$  tends to be ML

$$\boldsymbol{\mu}_N = \mathbf{m}_N \qquad \mathbf{C}_N = \frac{1}{N} \mathbf{C}$$

## 2.3 CONCLUSIONS

- If we can assume a parametric function for  $f_{\mathbf{x}}(\mathbf{x}|\omega_i)$  then the training phase of the classifier reduces to the estimation of the parameters.
- We can use two approaches for the estimation: ML (computationally simpler) or Bayesian (if we have a priori knowledge of parameters)