

Detection & Pattern Recognition

José Heinz

Electrical Engineering & Information Technology

heinzmollersantos@gmail.com

Resume

Here is the exam preparation for the course Detection & Pattern Recognition.

1. Mathematical notations and symbols

- true class: ω
- estimated class: $\hat{\omega}$
- probability of an event A : $P(A)$
- a priori of class ω_i : $P(\omega = \omega_i) = P(\omega_i)$
- PDF of a RV \underline{x} : $p(\underline{x})$
- likelihood of \underline{x} : $p(\underline{x} | \omega_i)$
- a posterior prob. of class ω_i : $P(\omega_i | \underline{x})$
- expectation of a RV \underline{x} : $\mu = E(\underline{x}) = \int \underline{x} p(\underline{x}) d\underline{x}$
- covariance matrix of \underline{x} : $C = E[(\underline{x} - \mu)(\underline{x} - \mu)^T]$
- discriminant function for class ω_i : $f_i(\underline{x})$

There are for binary detection four cases. Let H_0 and H_1 be the null hypothesis and the alternative hypothesis for binary detection. Let $H \in \{H_0, H_1\}$ be the unknown true state. Let $\hat{H}(\underline{x}) \in \{H_0, H_1\}$ be the estimated state of H .

$\hat{H} H$	$H = H_0$	$H = H_1$
$\hat{H} = H_0$	correct rejection true negative (TN)	miss false negative (FN) type II error
$\hat{H} = H_1$	false alarm false positive (FP) type I error	detection, hit true positive (TP)

Figure 1: Four possible combinations for \hat{H}

One can define the performance measures of this example for $\hat{H}(\underline{x})$ as for the

TN-Rate: $P_{TN} = P(\hat{H} = H_0 | H = H_0)$ with the cond. probability.

Analog for the FP-Rate $P_{FP} = P(\hat{H} = H_1 | H = H_0)$.

For the FN-Rate $P(\hat{H} = H_0 | H = H_1)$

for the TP-Rate $P(\hat{H} = H_1 | H = H_1)$

Magically this happens:

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\ P(A \cap B) &= 0 \iff A \cap B = \emptyset \\ \Rightarrow P_{TN} + P_{FP} &= 1 \\ \Rightarrow P_{FN} + P_{TP} &= 1 \end{aligned}$$

Now we can define the confusion matrix as:

$$n_{ij} = \#(\hat{\omega} = \omega_i, \omega = \omega_j) \in \mathbb{N} \\ [n_{ij}]_{ij} \in (c \times c),$$

where c is the number of classes. The properties of this crazy ass matrix are

- column sum: $n_{:j} = \sum_{i=1}^c n_{ij} = \#(\omega = \omega_j)$
- row sum: $n_{i:} = \sum_{j=1}^c n_{ij} = \#(\hat{\omega} = \omega_i)$
- Matrix sum: $N = \sum_{i=1}^c \sum_{j=1}^c n_{ij} = \text{total}$

Normalization of n_{ij} to make a probability:

- marginal probability: $P_j = P(\omega = \omega_j) = \frac{n_{:j}}{N}$ and $P(\hat{\omega} = \omega_i) = \frac{n_{i:}}{N}$ with $\sum_{j=1}^c P_j = 1$
- joint probability: $P(\hat{\omega} = \omega_i, \omega = \omega_j) = \frac{n_{ij}}{N}$ with $\sum_{i=1}^c \sum_{j=1}^c P_{ij} = 1$
- conditional probability (column sum normalization): $\hat{P}_{ij} = P(\hat{\omega} = \omega_i | \omega = \omega_j) = \frac{n_{ij}}{n_{:j}}$ which can be interpreted as the TN, FP, FN, TP- Rate.
- there is also row sum normalization.

Example:

1.4 Confusion matrix

E1.10: Confusion matrix of fish sorting (2)

$$[n_{ij}] = \begin{bmatrix} n_{11} & n_{12} \\ n_{21} & n_{22} \end{bmatrix} = \begin{bmatrix} 50 & 5 \\ 10 & 35 \end{bmatrix}, \quad \sum_i \sum_j n_{ij} = 100$$

Marginal probabilities

column: 60% salmon and 40% sea bass
row: 55% classified salmon and 45% classified sea bass

Joint probabilities

$$[P_{ij}] = [P(\hat{\omega} = \omega_i, \omega = \omega_j)] = \begin{bmatrix} 50\% & 5\% \\ 10\% & 35\% \end{bmatrix},$$

Conditional probabilities

$$\begin{aligned} [\bar{P}_{ij}] &= [P(\hat{\omega} = \omega_i | \omega = \omega_j)] = \begin{bmatrix} 50/60 & 5/40 \\ 10/60 & 35/40 \end{bmatrix} = \begin{bmatrix} P_{TN} & P_{FN} \\ P_{FP} & P_{TP} \end{bmatrix}, \\ [P(\omega = \omega_j | \hat{\omega} = \omega_i)] &= \begin{bmatrix} 50/55 & 5/55 \\ 10/45 & 35/45 \end{bmatrix}. \end{aligned}$$

Figure 2: Confusion matrix example

2. Bayesian decision theory

Given two events, for the general case; the Bayesian theorem/rule/formula is:

$$P(A | B) = \frac{P(A, B)}{P(B)} = \frac{P(A | B)P(B)}{P(A)}.$$

This nice formula allows the conversion between $P(A | B)$ and $P(B | A)$.

Now we may consider a RV \underline{x} and classes $\omega_i, i = 1, 2, \dots, c$.

$$P(\omega_j | \underline{x}) = \frac{p(\underline{x}, \omega_j)}{p(\underline{x})} = \frac{p(\underline{x} | \omega_j)P(\omega_j)}{\sum_{k=1}^c p(\underline{x} | \omega_k)P(\omega_k)}$$

i.e.,

$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}}$$

2.1 Minimum bayesian risk decision

The joint prob. matrix $[P_{ij}]$ can be derived from the confusion matrix $[n_{ij}]_{ij}$ with the following formula as stated: $P_{ij} = P(\hat{\omega} = \omega_i, \omega = \omega_j) = \frac{n_{ij}}{N}$. These are $c^2 - 1$ measures. With a few formulations we can derive the recognition rate and error rate:

$$\begin{aligned} 1 &= \sum_{i=1}^c \sum_{j=1}^c P_{ij} = \sum_{i=1}^c P_{ii} + \sum_{i \neq j}^c P_{ij} \\ &= \underbrace{P(\hat{\omega} = \omega)}_{\text{Recognition Rate}} + \underbrace{P(\hat{\omega} \neq \omega)}_{\text{Error Rate (ER)}}. \end{aligned}$$

We can now define the loss of a decision when the class is not correctly chosen:

$$l_{ij} = l(\hat{\omega} = \omega_i, \omega = \omega_j) \geq 0, \quad 1 \leq i, j \leq c,$$

where typically $l_{ii} = 0$ and $l_{ij} > 0 \forall i \neq j$.

With the loss we can penalize the decisions like this and call it Bayesian Risk BR:

$$BR = \sum_i \sum_j l_{ij} P_{ij} = E_{\hat{\omega}, \omega}[l(\hat{\omega}, \omega)] \geq 0$$

and if the loss function is the 0/1-loss: $l_{ij} = 1 \forall i \neq j$, meaning that all errors are equally bad, then the BR simplifies to the Error Rate ER:

$$BR = \sum_i \sum_{j \neq i} P_{ij} = ER.$$

as expectation of the loss over $\hat{\omega}$ and ω .

Now we get to define if the loss is dependent of the input vector \underline{x} . Knowing also that the estimated class $\hat{\omega}$ is also dependent of the RV \underline{x} . Then it follows that the loss is also dependent of the RV such that $l = l(\hat{\omega}(\underline{x}), \omega)$.

Then it can be derived that the expectation function over the loss is defined over (\underline{x}, ω) , meaning; computing the average loss for all (\underline{x}, ω) . In the following definition we define d as the number of dimensions or characteristics of $\underline{x} \in \mathbb{R}^d$:

$$\begin{aligned} BR &= E[l(\hat{\omega}(\underline{x}), \omega)] = \int_{\mathbb{R}^d} \sum_{j=1}^c l(\hat{\omega}(\underline{x}), \omega = \omega_j) p(\underline{x}, \omega_j) d\underline{x} \\ &= \underbrace{\int_{\mathbb{R}^d} \sum_{j=1}^c l(\hat{\omega}(\underline{x}), \omega = \omega_j) P(\omega_j, \underline{x}) p(\underline{x}) d\underline{x}}_{\text{conditional Risk: } R(\hat{\omega} | \underline{x})} \\ &\quad \int_{\mathbb{R}^d} R(\hat{\omega} | \underline{x}) p(\underline{x}) d\underline{x}. \end{aligned}$$

Meaning that in order to minimize the Bayesian Risk, we have to minimize the Conditional Risk, and therefore minimizing the loss of the system:

$$\min BR \Rightarrow \min_{\hat{\omega}} R(\hat{\omega} | \underline{x})$$

$$\min_{\hat{\omega}} \sum_{j=1}^c l(\hat{\omega}(\underline{x}), \omega = \omega_j) P(\omega_j, \underline{x}).$$

However, we are not done yet, bro... We have the following dependencies:

- $l(\hat{\omega}, \omega)$ which is chosen by YOU
- $P(\omega_j | \underline{x})$ calculated with BAYESIAN THEOREM
- $\hat{\omega}(\underline{x})$ TO BE DESIGNED!

2.2 MBR: First design of $\hat{\omega}(\underline{x})$

Let $\hat{\omega}(\underline{x}) = \hat{\omega}_{MBR}(\underline{x})$ be the MBR decision/detector or classifier be defined as

$$\hat{\omega}_{MBR}(\underline{x}) = \operatorname{argmin}_{\hat{\omega}} R(\hat{\omega} | \underline{x}) = \operatorname{argmin}_{\omega_i} R(\hat{\omega}_i = \omega_i | \underline{x}).$$

For this calculation one can do the following algorithm:

i.e. for any given \underline{x} :

- calculate c values $R(\hat{\omega}_i | \underline{x}) = \sum_{j=1}^c l_{ij} P(\omega_j | \underline{x})$ for $1 \leq i \leq c$
- look for minimum

or mathematically easier:

$$\begin{bmatrix} l_{11} & \dots & l_{1c} \\ \vdots & \dots & \vdots \\ l_{c1} & \dots & l_{cc} \end{bmatrix} \begin{bmatrix} P(\omega_1 | \underline{x}) \\ \vdots \\ P(\omega_c | \underline{x}) \end{bmatrix} = \begin{bmatrix} R(\hat{\omega} = \omega_1 | \underline{x}) \\ \vdots \\ R(\hat{\omega} = \omega_c | \underline{x}) \end{bmatrix}$$

For the Bayesian Rule vector, ignore the denominator $p(\underline{x})$: $P(\omega_j | \underline{x}) = \frac{p(\underline{x}, \omega_j)P(\omega_j)}{p(\underline{x})} \propto p(\underline{x}, \omega_j)P(\omega_j)$ because of independency of ω_j

2.3 Two-class MBR: Likelihood Ratio Test

For two classes $c = 2$, the algorithm can be simplified as the likelihood ratio test (LRT):

$$LR = \frac{p(\underline{x} | \omega_2)}{p(\underline{x}, \omega_1)} \stackrel{\omega_2}{\underset{\omega_1}{\lessdot}} \frac{l_{21} - l_{11}}{l_{12} - l_{22}} \cdot \frac{P(\omega_1)}{P(\omega_2)} = \frac{\gamma}{\text{Ratio}}$$

and applying the 0/1-loss:

$$LR = \frac{p(\underline{x} | \omega_2)}{p(\underline{x}, \omega_1)} \stackrel{\omega_2}{\underset{\omega_1}{\lessdot}} \frac{P(\omega_1)}{P(\omega_2)} = \frac{\gamma}{\text{Ratio}}$$

and also the conditional Risk changes when applied the 0/1-loss to

$$\begin{aligned} R(\hat{\omega} = \omega_i | \underline{x}) &= \sum_{j=1}^c l_{ij} P(\omega_j | \underline{x}) \\ \stackrel{0/1\text{-loss}}{=} & \sum_{j=1, j \neq i}^c P(\omega_j | \underline{x}) + P(\omega_i | \underline{x}) - P(\omega_i | \underline{x}) \\ &= 1 - P(\omega_i | \underline{x}) \end{aligned}$$

Hence $\min_{\omega_i} ER$ is equivalent to $\max_{\omega_i} P(\omega_i | \underline{x})$ which is the Max. Aposterior prob. (MAP).

2.4 Maximum Likelihood; Special case for MAP

Let the probability of c classes be equal, meaning $P(\omega_i) = \frac{1}{c}, \forall i$. The for the computation of the posterior:

$$\begin{aligned} P(\omega_i | \underline{x}) &= p(\underline{x} | \omega_i) \cdot \frac{P(\omega_i)}{p(\underline{x})} \\ \Rightarrow \max_{\omega_i} P(\omega_i | \underline{x}) &\Leftrightarrow \underbrace{\max_{\omega_i} p(\underline{x} | \omega_i)}_{\text{Maximum Likelihood (ML)}} \end{aligned}$$

which is possible because $p(\underline{x})$ is independent of ω_i

2.5 Discriminant Functions using MBR, MAP, ML Decision

The main goal of this section is to create a unified framework with c discriminant functions $f_i(\underline{x})$ for $\omega_i 1 \leq i \leq c$. The main objective is also to find the best estimated class $\hat{\omega}(\underline{x}) = \operatorname{argmax}_i f_i(\underline{x})$ with $\hat{\omega}(\underline{x}) \in \{\omega_1, \dots, \omega_c\}$.

Defining first a few elements:

- Decision Region: $R_i = \{\underline{x} | f_i(\underline{x}) > f_j(\underline{x}) \forall j \neq i\}$
- Boundary: $\beta_{ij} = \{\underline{x} | f_i(\underline{x}) = f_j(\underline{x}) > f_k(\underline{x}) \forall k \neq i, j\}$

b) $c = 10$ classes

- MBR:** $f_i(\underline{x}) = R(\hat{\omega} = \omega_j \mid \underline{x}) = \sum_{j=1}^c l_{ij} P(\omega_j \mid \underline{x})$ (no log function)
- MAP:** $f_i(\underline{x}) = P(\omega_i \mid \underline{x}) \sim p(\underline{x} \mid \omega_i) \cdot P(\omega_i)$ or with the log fct: $f_i(\underline{x}) = \ln(p(\underline{x} \mid \omega_i)) + \ln(P(\omega_i))$
- ML:** $f_i(\underline{x}) = p(\underline{x} \mid \omega_i)$ or $\ln(p(\underline{x} \mid \omega_i))$

2.6 MAP decision for Gaußian likelihoods

Assumptions:

- Classes $c \geq 2$
- Features $\underline{x} \in \mathbb{R}^d, d \geq 1$
- MAP: 0/1-loss
- Known priors $P(\omega_i) = P_i$
- Known Gaußian Likelihood $(\underline{x}, \omega_i) \sim N(\underline{\mu}_i \mid \mathbf{C}_i)$

The MAP discriminant function for class ω_i is then described by using the log function:

$$f_i(\underline{x}) = \ln p(\underline{x} \mid \omega_i) + \ln P_i = -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_i| - \frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \mathbf{C}_i^{-1} (\underline{x} - \underline{\mu}_i) + \ln P_i$$

Case A: $\mathbf{C}_i = \sigma^2 \mathbf{I}$

For decision region:

$$\begin{aligned} f_i(\underline{x}) &= -\frac{1}{2\sigma^2} \|\underline{x} - \underline{\mu}_i\|^2 + \ln P_i \\ &\text{with } \|\underline{x}\| = \text{const} \\ \Rightarrow f_i(\underline{x}) &= \underbrace{\left(\frac{\underline{\mu}_i}{2}\right)^T \underline{x}}_{w_{i0} \text{ (offset)}} + \underbrace{\ln P_i - \frac{1}{2\sigma^2} \|\underline{\mu}_i\|^2}_{w_i^T \underline{x}} \end{aligned}$$

as an affine function for every class.

For decision boundary β_{ij} between ω_i and ω_j :

$$\begin{aligned} f_i(\underline{x}) - f_j(\underline{x}) &= 0 \\ (\underline{w}_i - \underline{w}_j)^T \underline{x} + (w_{i0} - w_{j0}) &= 0. \end{aligned}$$

Properties:

- a linear flat decision boundary
- $\underline{\mu}_i - \underline{\mu}_j = \sigma^2(\underline{w}_i - \underline{w}_j) \perp \beta_{ij}$
- if $P_i = P_j$: $\max_i f_i(\underline{x}) \triangleq \min_i \|\underline{x} - \underline{\mu}_i\|$ as nearest mean

E2.5: MAP decision for Gaussian likelihood – Case A

- $d = 2$ features $\underline{x} = [x_1, x_2]^T$ with Gaussian likelihood $\underline{x} \mid \omega_i \sim N(\underline{\mu}_i, \mathbf{C}_i)$
- equal covariance matrices $\mathbf{C}_i = \sigma^2 \mathbf{I}$,
- equal priors $P(\omega_i) = \text{const}$

In this case, the MAP decision rule simplifies to the **nearest mean classifier** (see section 4.2). Each point in the following **Voronoi diagram** represents a class center $\underline{\mu}_i$ and the lines represent the decision boundaries.

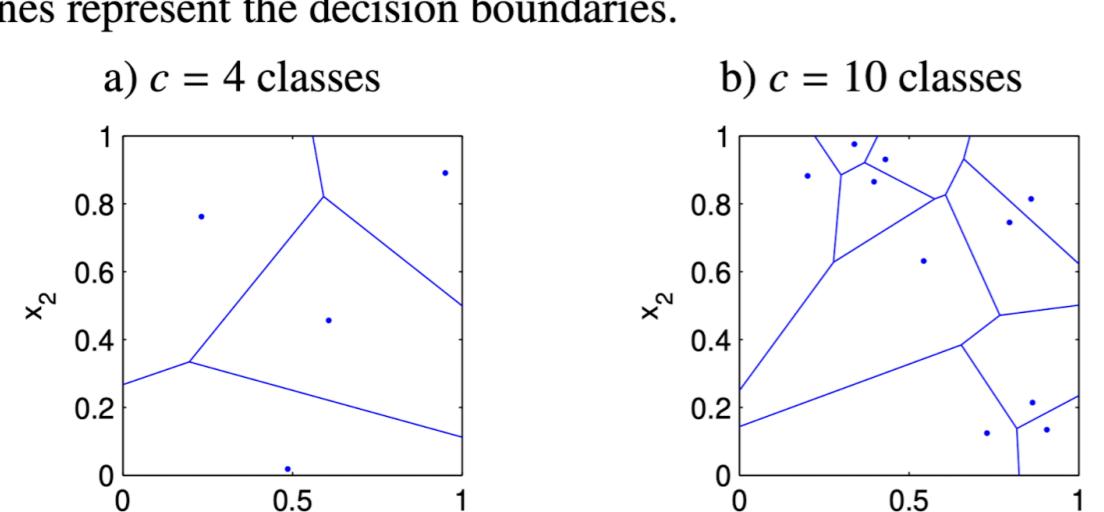


Figure 4: Example of MAP decision

Case B: $\mathbf{C}_i = \mathbf{C} \neq \sigma^2 \mathbf{I}$

$$\begin{aligned} f_i(\underline{x}) &= -\frac{1}{2} \underbrace{(\underline{x} - \underline{\mu}_i)^T \mathbf{C}^{-1} (\underline{x} - \underline{\mu}_i)}_{\text{Mahalanobis Distance}} + \ln P_i \\ &= \underline{w}_i^T \underline{x} + w_{i0} \end{aligned}$$

with $\underline{w}_i = \mathbf{C}^{-1} \underline{\mu}_i$

$$w_{i0} = \ln P_i - \frac{1}{2} \underline{\mu}_i^T \mathbf{C}^{-1} \underline{\mu}_i$$

The **Euclidean distance** between \underline{x} and $\underline{\mu}$ is defined by

$$D_{\text{Eucl}}(\underline{x}, \underline{\mu}) = \sqrt{(\underline{x} - \underline{\mu})^T (\underline{x} - \underline{\mu})}$$

All \underline{x} with the same Euclidean distance to $\underline{\mu}$ lie on a ball centered at $\underline{\mu}$. Given two centers $\underline{\mu}_i$ and $\underline{\mu}_j$, all \underline{x} having the same Euclidean distance to $\underline{\mu}_i$ and $\underline{\mu}_j$ lie on the red solid line.

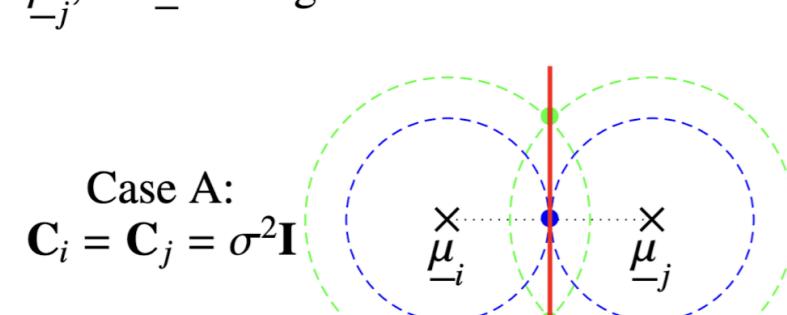


Figure 5: Mahalanobis vs Euclidean distance

Case C: $\mathbf{C}_i \neq \mathbf{C}_j$

E2.6: MAP decision for Gaussian likelihood – Case C

Assume $P_1 = P_2 = 1/2$ and

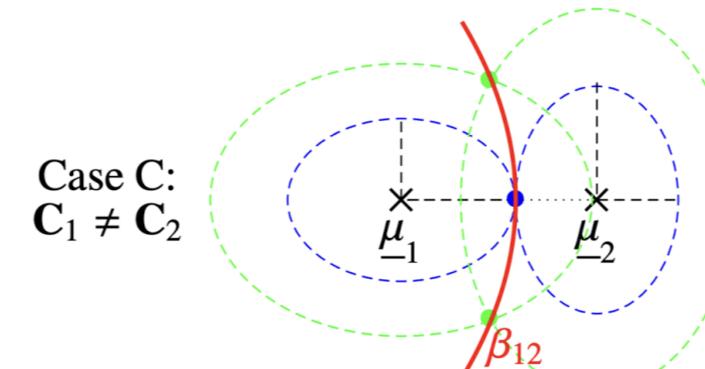
$$\underline{\mu}_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \underline{\mu}_2 = \begin{bmatrix} 10 \\ 0 \end{bmatrix}, \quad \mathbf{C}_1 = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{C}_2 = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}.$$

The MAP decision boundary β_{12} is then given by

$$\ln |\mathbf{C}_1| + (\underline{x} - \underline{\mu}_1)^T \mathbf{C}_1^{-1} (\underline{x} - \underline{\mu}_1) = \ln |\mathbf{C}_2| + (\underline{x} - \underline{\mu}_2)^T \mathbf{C}_2^{-1} (\underline{x} - \underline{\mu}_2)$$

or after some simplifications

$$x_1^2 - x_2^2 - 40x_1 + 200 = 0.$$



Summary of MAP decision boundaries for Gaussian likelihood

	distance of \underline{x} on β_{ij} to $\underline{\mu}_i$ and $\underline{\mu}_j$	decision boundary β_{ij}
Case A: $\mathbf{C}_i = \sigma^2 \mathbf{I}$	same Euclidean distance	affine in \underline{x} , orthogonal to $\underline{\mu}_i - \underline{\mu}_j$
Case B: $\mathbf{C}_i = \mathbf{C}_j \neq \sigma^2 \mathbf{I}$	same Mahalanobis distance	affine in \underline{x} , not orthogonal to $\underline{\mu}_i - \underline{\mu}_j$
Case C: $\mathbf{C}_i \neq \mathbf{C}_j$	different Mahalanobis distances	quadratic in \underline{x}

Figure 6: Summary

For minimax theory and Neyman-Pearson lemma, the lecture notes.

3. Detection

Given:

- data \underline{x}
- likelihood $p(\underline{x} \mid \omega)$
- prior (optional)
- loss (optional)

Desired is a detector with good detector performance for **binary** detection. We use now the notation H_0 and H_1 instead of ω . Defining also the estimated hypothesis $\hat{H}(\underline{x}) \in \{H_0, H_1\}$.

We define a received signal as a vector with a time dependency n . The vector $\underline{x} = [x(0), x(1), \dots, x(d-1)]^T \in \mathbb{R}^d$.

We can define two hypotheses. $H_0 : x(n) = z(n)$ and $H_1 : x(n) = A + z(n)$ and $z(n)$ as gaußian noise $N(0, \sigma^2)$. and the losses l_{ij} for $\hat{H} \in \{H_0, H_1\}$.

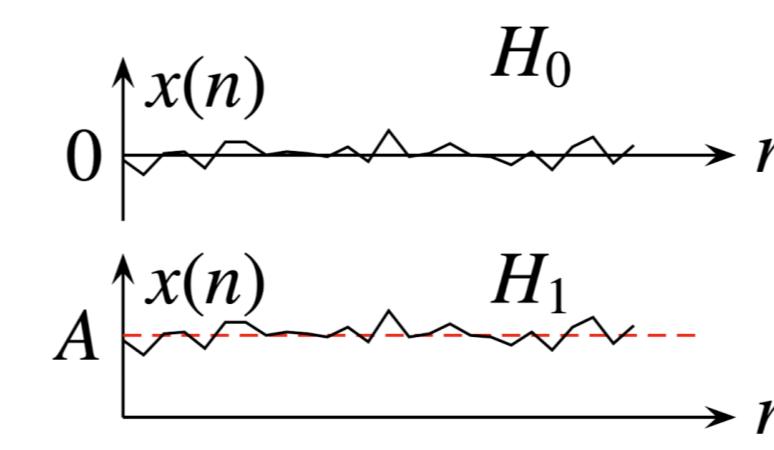


Figure 7: Representation of the signals to be detected

We got also the known priors $P_0 = P(H_0)$ and $P_1 = P(H_1) = 1 - P_0$. Then the likelihoods are

$$\begin{aligned} p(\underline{x} \mid H_0) &= \prod_{n=0}^{d-1} p(x(n) \mid H_0) = \frac{1}{(2\pi\sigma^2)^{d/2}} e^{(-\frac{1}{2\sigma^2} \sum_{n=0}^{d-1} x(n)^2)} \\ p(\underline{x} \mid H_1) &= \prod_{n=0}^{d-1} p(x(n) \mid H_1) = \frac{1}{(2\pi\sigma^2)^{d/2}} e^{(-\frac{1}{2\sigma^2} \sum_{n=0}^{d-1} (x(n) - A)^2)} \end{aligned}$$

One can use a general LRT detector:

$$\begin{aligned} \ln \frac{p(\underline{x} \mid H_1)}{p(\underline{x} \mid H_0)} &= -\frac{1}{2} \sum_{n=0}^{d-1} (x(n) - A)^2 - \sum_{n=0}^{d-1} x(n)^2 \\ &= \underbrace{\frac{dA^2}{\rho > 0} \left[\frac{1}{d} \sum_{n=0}^{d-1} x(n) - \frac{1}{2} \right]}_{t(\underline{x})} \\ &= \rho [t(\underline{x}) - \frac{1}{2} \underbrace{\ln \frac{H_1}{H_0}}_{\text{LR}}] \end{aligned}$$

$$\Rightarrow t(\underline{x}) \leqslant \frac{H_1}{H_0} \frac{1}{\rho} \ln \gamma + \frac{1}{2} = \gamma$$

where γ is the new threshold for $t(\underline{x})$. For this notation is:

- Test: $t(\underline{x}) = \frac{1}{Ad} \sum_{n=0}^{d-1} x(n) = \frac{1}{Ad}$
- SNR (signal to noise ratio) $\frac{A^2}{\sigma^2}$ of $x(n) \mid H_1 = x(n) = Az(n) \sim N(A, \sigma^2)$
- SNR $\rho = \frac{dA^2}{\sigma^2}$ of $\underline{x} \mid H_1 = A + \bar{z} \sim N(A, \frac{\sigma^2}{d})$

Distributions of the function $t(\underline{x})$:

$$\begin{aligned} H_0 : t(\underline{x}) &= \frac{1}{A} \cdot \bar{z} \sim N(0, \frac{\sigma^2}{A^2 d}) = N(0, \frac{1}{\rho}) \\ H_1 : t(\underline{x}) &= \frac{1}{A} (A + \bar{z}) \sim N(1, \frac{1}{\rho}) \end{aligned}$$

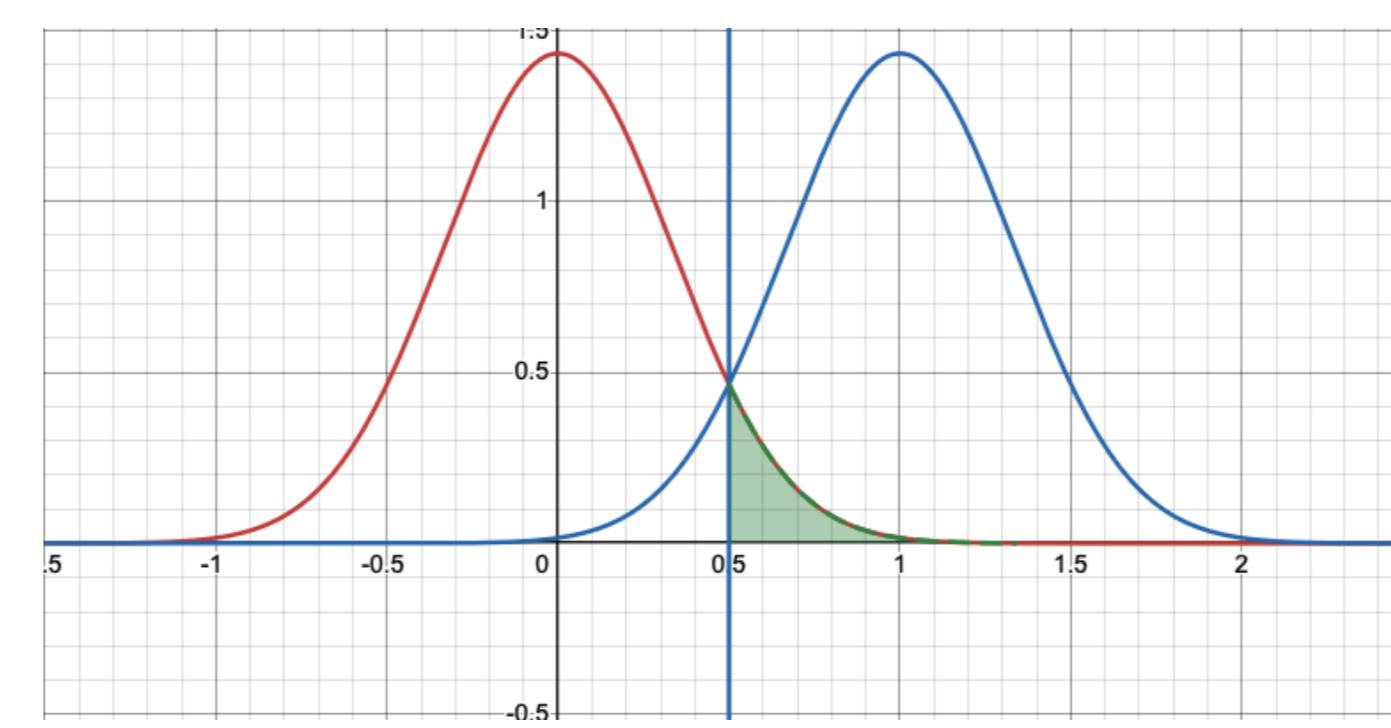


Figure 8: False Alarm rate in red

The next step is to minimize the P_{FA} , which is the false alarm rate. The line on figure 8 is the γ which is described in the last equations. For this we remember from Mathematics the Standard Gaußian distribution and Q function:

3 Detection

Standard Gaussian distribution and Q-function

A random variable x is **standard Gaussian** or **standard normal** distributed if it has the PDF

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$

It has zero mean and unit variance. The short notation is $x \sim N(0, 1)$. The **cumulative distribution function (CDF)** of x is $p(x) = \int_{-\infty}^x p(z) dz$ with $p(x) = d\phi(x)/dx$. The **Q-function** is defined as

$$Q(x) = \int_x^\infty p(z) dz = \int_x^\infty \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz.$$

$Q(\gamma)$ describes the probability $P(x > \gamma)$ that x obtains a value larger than γ . It corresponds to the area below $p(x)$ from γ to ∞ and decreases from 1 to 0 as γ increases.

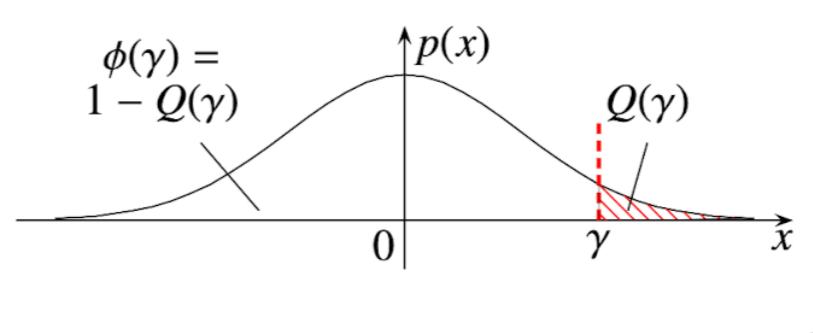


Figure 9: Q-Function

We can then define the false alarm rate as:

$$\begin{aligned} P_{FA} &= P(\hat{H} = H_1 \mid H = H_0) = P(t > \gamma \mid H_0), \text{ with } (t \mid H_0) \sim N(0, \frac{1}{\rho}) \\ &= P(\underbrace{\sqrt{\rho t}}_{N(0,1)} > \sqrt{\rho \gamma} \mid H_0) = Q(\sqrt{\rho \gamma}) \end{aligned}$$

and we also define the Detection Rate:

$$\begin{aligned} P_D(\gamma) &= P(\hat{H} = H_1 \mid H = H_1) = P(t > \gamma \mid H_1), \text{ with } (t \mid H_1) \sim N(1, \frac{1}{\rho}) \\ &= P(\underbrace{\sqrt{\rho}(t-1)}_{N(0,1)} > \sqrt{\rho}(\gamma-1) \mid H_1) = Q(\sqrt{\rho}(\gamma-1)) \end{aligned}$$

3 Detection

E3.1: Detection of DC in AWGN (cont.)

In E3.1 we assumed to know the priors P_0, P_1 and Gaussian likelihoods $P(\underline{x} \mid H_0), P(\underline{x} \mid H_1)$ including the DC value A and noise variance σ^2 . We studied different detectors (MBR, MAP, ML, NP). Are they still possible if A and/or σ^2 are unknown? Remember that the detector is

$$t(\underline{x}) = \frac{\bar{x} H_1}{A H_0} \gamma \quad \text{or} \quad \bar{x} = \frac{1}{d} \sum_{n=0}^{d-1} x(n) \frac{H_1}{H_0} \gamma,$$

where the threshold γ for \bar{x} depends on the type of detector as well as on A and σ^2 .

Detector	γ	$A\gamma</$
----------	----------	-------------

Given
training samples $\{(\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)\}$
Parameter
a suitable distance metric $D(\underline{x}, \underline{\mu})$
Training
• estimate all class centers $\hat{\underline{\mu}}_j = \frac{1}{N_j} \sum_{y_n=\omega_j} \underline{x}_n \quad (1 \leq j \leq c)$
• estimate all class covariance matrices when using the Mahalanobis distance $\hat{\mathbf{C}}_j = \frac{1}{N_j} \sum_{y_n=\omega_j} (\underline{x}_n - \hat{\underline{\mu}}_j)(\underline{x}_n - \hat{\underline{\mu}}_j)^T \quad (1 \leq j \leq c)$
Classification
for each new feature vector \underline{x} ,
• calculate the distance of \underline{x} to all class centers $d_j = D_{\text{Euc}}(\underline{x}, \hat{\underline{\mu}}_j) \text{ or } D_{\text{Maha}}(\underline{x}, \hat{\underline{\mu}}_j, \hat{\mathbf{C}}_j) \quad (1 \leq j \leq c)$
• choose the class with the smallest distance to \underline{x} $\hat{y}(\underline{x}) = \arg \min_{\omega_j} d_j$

Comments

- multiclass classifier
- computational complexity
 - training: $O(Nd)$ for $\hat{\underline{\mu}}_j$ and $O(Nd^2)$ for $\hat{\mathbf{C}}_j$
 - classification: c distance calculations

$p_m(\underline{x}; \underline{\theta}) \sim N(\underline{\mu}, \mathbf{C}_m)$
 $\underline{\theta} \{ \underline{\mu}, \mathbf{C}_m \}, 1 \leq m \leq M$

Introducing a new random mode variable $Z \in \{1, 2, \dots, M\}$ as the label for each mode in a PDF. We can use the law of total probability, adding the new random mode variable:

$$p(\underline{x}, \underline{\theta}) = \sum_{m=1}^M P(z=m) p(\underline{x} | z=m) = \sum_{m=1}^M \alpha_m p_m(\underline{x}; \underline{\theta}_m)$$

also

$$\alpha_m = P(z=m) \geq 0, \quad \sum_{m=1}^M \alpha_m = 1$$

as prior of z , **mode weight** with $m-1$ parameters

$\underline{\theta} = \{\underline{\theta}_1, \dots, \underline{\theta}_m; \alpha_1, \dots, \alpha_{m-1}\}$ contains $M \frac{d(d+3)}{2} + m-1$ parameters if covariance matrix non diagonal,
 $2Md + M - 1$ if the covariance matrix diagonal.

Approach A:

Assuming $d=1$ and $M=2$, i.e.,

$$p(\underline{x}; \underline{\theta}) = \alpha_1 p_1(x) + \alpha_2 p_2(x), \quad \alpha_1 = \alpha, \quad \alpha_2 = 1 - \alpha,$$

$$p_i(x) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2\sigma_i^2}(x - \mu_i)^2\right), \quad i = 1, 2$$

with $\underline{\theta} = [\mu_1, \sigma_1^2, \mu_2, \sigma_2^2, \alpha]^T$. Given N i.i.d. measurements $x_n \sim p(\underline{x}; \underline{\theta})$, $1 \leq n \leq N$, the log-likelihood is

$$L(\underline{\theta}) = \ln \prod_{n=1}^N p(x_n; \underline{\theta}) = \sum_{n=1}^N \ln p(x_n; \underline{\theta}).$$

The gradient vector $\underline{\nabla} L$ of $L(\underline{\theta})$ contains the following elements:

$$\begin{aligned} \frac{\partial L}{\partial \mu_1} &= \sum_{n=1}^N \frac{\alpha_1 p(x_n; \underline{\theta})}{p(x_n; \underline{\theta})} \cdot \frac{x_n - \mu_1}{\sigma_1^2}, \quad i = 1, 2 \\ \frac{\partial L}{\partial \sigma_1^2} &= \sum_{n=1}^N \frac{\alpha_1 p(x_n; \underline{\theta})}{p(x_n; \underline{\theta})} \cdot \frac{(x_n - \mu_1)^2}{\sigma_1^2} - 1, \quad i = 1, 2 \\ \frac{\partial L}{\partial \alpha} &= \sum_{n=1}^N \frac{p_1(x_n) - p_2(x_n)}{p(x_n; \underline{\theta})}. \end{aligned}$$

Obviously, the nonlinear equation system $\underline{\nabla} L = \underline{0}$ has no closed-form solution.

Figure 15: this approach fails often

Approach B: A probability model of N i.i.d. $\underline{x}_n \in \mathbb{R}^d$ and $z_n, 1 \leq n \leq N$. With z_n as additional mode/sub-class measurement. Let

$$\underline{x} = \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_N \end{bmatrix} \quad Nd \times 1, \quad \underline{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix} \quad N \times 1$$

and we build a joint probe mode of \underline{x} and \underline{z} :

$$\begin{aligned} p(\underline{x}, \underline{z}; \underline{\theta}) &\stackrel{i.i.d.}{=} p(\underline{x}_n, z_n; \underline{\theta}) \\ &\Rightarrow \max_{\underline{\theta}} p(\underline{x}_n, z_n; \underline{\theta}), \text{ as GMM} \end{aligned}$$

But for this, we have to know what \underline{x} and \underline{z} exactly are. Therefore this only works if the mode labels are known!

EM algorithm (1)

The **Expectation-Maximization (EM)** algorithm from Dempster (1977) is a popular iterative method for finding ML parameter estimates when we have missing (bad) data.

Let $p(\underline{x}, \underline{z}; \underline{\theta})$ be the likelihood where $\underline{\theta}$ is the unknown parameter to be estimated, \underline{x} is the observed (good) data, and \underline{z} is the missing (bad) data. If we assume to know both \underline{x} and \underline{z} , we could easily find the ML estimate of $\underline{\theta}$ by $\hat{\underline{\theta}} = \arg \max_{\underline{\theta}} p(\underline{x}, \underline{z}; \underline{\theta})$. But \underline{z} is unknown. Can we still estimate $\underline{\theta}$ from \underline{x} ? Yes, it is possible by using the EM algorithm. It is an iterative method which starts with an initial guess $\hat{\underline{\theta}}(0)$ about $\underline{\theta}$ and produces a sequence of improved parameter estimates $\hat{\underline{\theta}}(1), \hat{\underline{\theta}}(2), \dots$

Let $L(\underline{x}, \underline{z}; \underline{\theta}) = \ln p(\underline{x}, \underline{z}; \underline{\theta})$ be the log-likelihood. At each iteration $k = 1, 2, \dots$, the EM algorithm consists of two steps. The first **expectation step (E-step)** computes the expectation of the log-likelihood over the missing data \underline{z} conditioned on the current parameter estimate $\hat{\underline{\theta}}(k)$ in order to "remove" \underline{z} .

EM algorithm (2)

For this purpose, we first calculate the marginal PDF of \underline{z} conditioned on \underline{x} and $\hat{\underline{\theta}}(k)$:

$$p(\underline{z} | \underline{x}, \hat{\underline{\theta}}(k)) = \frac{p(\underline{x}, \underline{z}; \hat{\underline{\theta}}(k))}{p(\underline{x}; \hat{\underline{\theta}}(k))} = \frac{p(\underline{x}, \underline{z}; \hat{\underline{\theta}}(k))}{\int p(\underline{x}, \underline{z}; \hat{\underline{\theta}}(k)) d\underline{z}}$$

Then an expectation of $L(\underline{x}, \underline{z}; \underline{\theta})$ over $\underline{z} | \underline{x}, \hat{\underline{\theta}}(k)$ returns an averaged log-likelihood

$$Q(\underline{\theta}) = E_{\underline{z} | \underline{x}, \hat{\underline{\theta}}(k)} L(\underline{x}, \underline{z}; \underline{\theta}) = \int (\ln p(\underline{x}, \underline{z}; \underline{\theta})) p(\underline{z} | \underline{x}, \hat{\underline{\theta}}(k)) d\underline{z}$$

which only depends on $\underline{\theta}$.

The second step, the **maximization step (M-step)**, maximizes $Q(\underline{\theta})$ over $\underline{\theta}$ and returns an improved parameter estimate

$$\hat{\underline{\theta}}(k+1) = \arg \max_{\underline{\theta}} Q(\underline{\theta})$$

The iteration stops if a given number of iterations is reached or if $\|\hat{\underline{\theta}}(k+1) - \hat{\underline{\theta}}(k)\|$ is below a certain threshold.

Table 4.5: GMM classifier (1)

Given
training samples $\{(\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)\}$
Parameter
no
Training
We assume a Gaussian likelihood $N(\underline{\mu}_j, \mathbf{C}_j)$ for each class ω_j .
• estimate the class means by $\hat{\underline{\mu}}_j = \frac{1}{N_j} \sum_{y_n=\omega_j} \underline{x}_n$
• estimate the class covariance matrices by $\hat{\mathbf{C}}_j = \frac{1}{N_j} \sum_{y_n=\omega_j} (\underline{x}_n - \hat{\underline{\mu}}_j)(\underline{x}_n - \hat{\underline{\mu}}_j)^T$
• estimate the priors by $\hat{P}(\omega_j) = \frac{N_j}{N} \quad (1 \leq j \leq c)$
Classification
use $N(\hat{\underline{\mu}}_j, \hat{\mathbf{C}}_j)$ and $\hat{P}(\omega_j)$ in the Bayesian decision rules (MBR, MAP, ML etc.) from ch. 2

Figure 14: Gaußian classifier alg.

Problem is that it has to be guaranteed that the covariance matrix is invertible, so we can specify the "naive Gaussian classifier" as just the diagonal elements, making the matrix invertible at all feasible points.

4.4 Gaußian Mixture Model

Is a Mixture of Gaußian PDFs. Here we introduce the **Unimodal probability Function**. This means using GMM for one class with m Gaußian Modes (sub-classes).

Given

training samples $\{(\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)\}$

Parameter

model order M_j for all classes ω_j ($1 \leq j \leq c$)

Training

Estimate for each class ω_j the GMM parameters from the N_j training samples of class ω_j by using the EM algorithm from Table 4.4.

Classification

use the estimated GMM models in the Bayesian decision rules (MBR, MAP, ML etc.) from ch. 2

4.5 Parzen Window Method (Kernel Function)

to estimate $p(\underline{x})$ with $p(\underline{x})$

Given
training samples $\{(\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)\}$
Parameter
• a suitable kernel function $\phi(\underline{x})$ with $\int \phi(\underline{x}) d\underline{x} = 1$
• bandwidth $h > 0$ with $\phi_h(\underline{x}) = \frac{1}{h^d} \phi(\frac{\underline{x}}{h})$
Training
no or
estimate the kernel function $\phi_h(\underline{x})$ or the bandwidth h
ML classification
for each new feature vector \underline{x} ,
• estimate the likelihood $p(\underline{x} \omega_j)$ for each class ω_j by using the Parzen window method
$\hat{p}(\underline{x} \omega_j) = \frac{1}{N_j} \sum_{y_n=\omega_j} \phi_h(\underline{x} - \underline{x}_n) \quad (1 \leq j \leq c)$
• choose the class with the largest likelihood
$\hat{y}(\underline{x}) = \arg \max_{\omega_j} \hat{p}(\underline{x} \omega_j)$

Comments

- multiclass classifier
- computational complexity
 - classification: N kernel calculations
 - high classification complexity
 - difficult choice of h
 - needs a huge amount of training data for a large dimension d

Figure 17: Parzen Window alg.

4.6 Discriminant Functions

The main idea is to create c discriminant functions $f_j(\underline{x})$ for c classes ω_j . For each \underline{x} : $\hat{y}(\underline{x}) = \arg \max_{\omega_j} f_j(\underline{x})$.

Training **Bayes Plug-In Discriminant**
 need models for $p(\underline{x} | \omega_j; \underline{\theta})$ $f_j(\underline{x}; \underline{\theta})$
 estimate from $\underline{\theta}_j, P(\omega_j)$ $\underline{\theta}_j$
 derive $f_j(\underline{x})$ from $p(\underline{x}_j, \omega_j; \underline{\theta})$

Tabela 1: Table of Prof. Yang to describe this section comparing the older sections.

Advantage of discriminant functions is that we do not need the $\underline{x} | \omega_j$. We can define our space $\underline{x} \in \mathbb{R}^d$ and our decision Region $R_i = \{\underline{x} | f_i(\underline{x}) > f_j(\underline{x}), i \neq j\}$ and our decision boundary $\beta_{ij} = \{\underline{x} | f_i(\underline{x}) = f_j(\underline{x})\}$ where the decision Regions are defined to be convex. (For the convex definition, see lecture.)

4.7 Linear (affine) discriminant functions

We define our discriminant function with the learnable parameters \underline{w}_j and w_{j0} as $f_j(\underline{x}) = \underline{w}_j^T \underline{x} + w_{j0}$ for the class ω_j with \underline{w}_j as the weight vector and w_{j0} as the offset. Example for two classes with classification capability:

$$\begin{aligned} f(\underline{x}) &= f_1(\underline{x}) - f_2(\underline{x}) = (\underline{w}_1 - \underline{w}_2)^T \underline{x} + (w_{10} - w_{20}) \\ &= \underline{w}^T \underline{x} + w_0 \stackrel{\omega_1}{\leq} 0 \end{aligned}$$

and the decision boundary $f(\underline{x}) = 0$ as an hyperplane. The vector \underline{w} is the normal vector to the hyperplane.

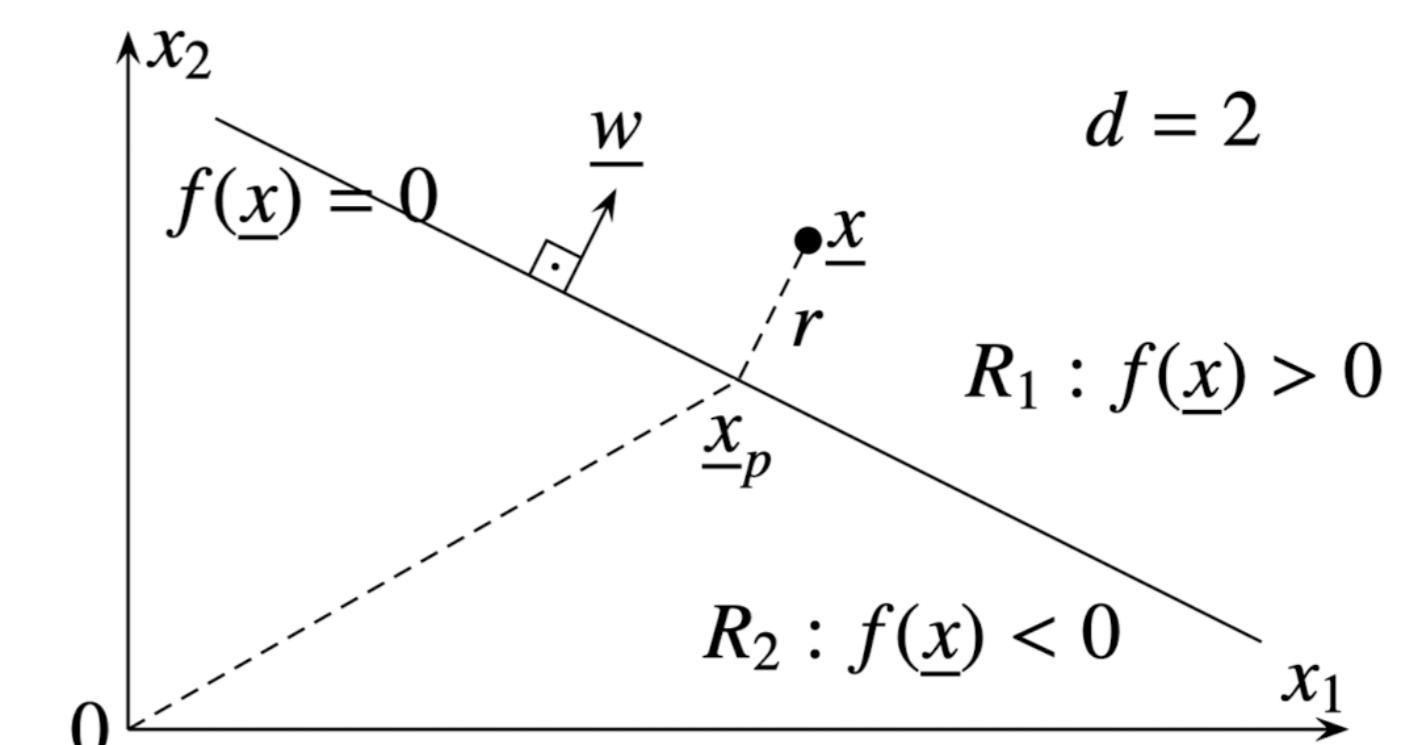


Figure 18: Hyperplane and the two decision regions.

However this approach with a linear affine function as discriminant function only works well for linearly separable data and is also not feasible in most cases because of the $N(c-1)$ equalities and infinite many solutions for linearly separable data. For non linearly separable data, there is no solution. To find a solution for the set S of infinite many solutions with the linear discriminant function: Training by optimization.

$$\begin{aligned} \underline{\theta} &= \{(\underline{w}_j, w_{j0}) \text{ from } S\} \\$$

$$J(\underline{\theta}) = \sum_{n=1}^N |\underline{w}^T \underline{x} + w_0 - y_n|$$

- c) **Maximum Margin:** see Support vector machine.
d) **Cross entropy:** see Deep Learning.

4.8 Feature Mapping

Is an approach to make non linearl separable data linear and non convex regions convex. See the examples on the lecture.

It is powerful but, d features grow and become unmanageable. Huge computational complexity. Might lead to overfitting.

4.9 Support vector machine

5. 5

Referências

- [1] J. Flusser and T. Suk. Pattern recognition by affine moment invariants. *Pattern Recognition*, 26:167–174, 1993.

- [2] M. K. Hu. Visual pattern recognition by moment invariants. *IRE Transactions on Information Theory*, IT-8:179–187, 1962.
[3] P. Maragos. Pattern spectrum and multiscale shape representation. *IEEE Trans. Patt. Anal. Mach. Intell.*, 11:701–715, 1989.
[4] A. Meijster and M. H. F. Wilkinson. A comparison of algorithms for connected set openings and closings. *IEEE Trans. Patt. Anal. Mach. Intell.*, 24(4):484–494, 2002.
[5] P. F. M. Nacken. *Image Analysis Methods Based on Hierarchies of Graphs and Multi-Scale Mathematical Morphology*. PhD thesis, University of Amsterdam, Amsterdam, The Netherlands, 1994.