

# Intelligent Systems: Reasoning and Recognition

James L. Crowley

ENSIMAG 2 / MoSIG M1

Second Semester 2011/2012

Lesson 17

18 april 2012

## Normal Probability Density Functions

Notation .....	2
Bayesian Classification.....	3
Quadratic Discrimination.....	4
Discrimination using Log Likelihood .....	6
Example for $K > 2$ and $D > 1$ .....	7
Canonical Form for the discrimination function .....	9
Noise and Discrimination .....	11
Decision Surfaces for different Noise assumptions .....	13
Two classes with equal means .....	15

Sources Bibliographiques :

"Pattern Recognition and Machine Learning", C. M. Bishop, Springer Verlag, 2006.

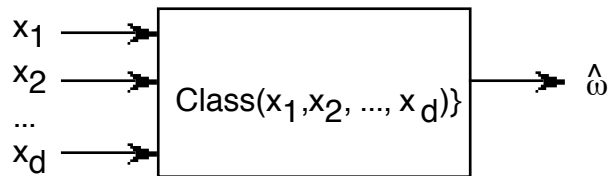
"Pattern Recognition and Scene Analysis", R. E. Duda and P. E. Hart, Wiley, 1973.

**Notation**

$\vec{x}$	A vector of $D$ variables.
$\vec{X}$	A vector of $D$ random variables.
$D$	The number of dimensions for the vector $\vec{x}$ or $\vec{X}$
$E$	An observation. An event.
$C_k$	The class $k$ .
$k$	Class index
$K$	Total number of classes
$\omega_k$	The statement (assertion) that $E \in C_k$
$p(\omega_k) = p(E \in C_k)$	Probability that the observation $E$ is a member of the class $k$ . Note that $p(\omega_k)$ is lower case.
$P(X)$	Probability density function for $X$
$P(\vec{X})$	Probability density function for $\vec{X}$
$P(\vec{X}   \omega_k)$	Probability density for $\vec{X}$ the class $k$ . $\omega_k = E \in T_k$ .

## Bayesian Classification

Our problem is to build a box that maps a set of features  $\vec{X}$  from an Observation, E into a class  $T_k$  from a set of K possible Classes.



Let  $\omega_k$  be the proposition that the event belongs to class k:  $\omega_k = E \in T_k$

$\omega_k$  Proposition that event  $E \in$  the class k

In order to minimize the number of mistakes, we will maximize the probability that  $\omega_k \equiv E \in T_k$

$$\hat{\omega}_k = \arg\max_k \{ \Pr(\omega_k | \vec{X}) \}$$

We will call on two tools for this:

1) Baye's Rule :

$$p(\omega_k | \vec{X}) = \frac{P(\vec{X} | \omega_k) p(\omega_k)}{P(\vec{X})}$$

2) Normal Density Functions

$$P(\vec{X} | \omega_k) = \frac{1}{(2\pi)^{\frac{D}{2}} \det(C_k)^{\frac{1}{2}}} e^{-\frac{1}{2}(\vec{X}-\vec{\mu}_k)^T C_k^{-1}(\vec{X}-\vec{\mu}_k)}$$

Last week we looked at Baye's rule. Today we concentrate on Normal Density Functions.

## Quadratic Discrimination

The classification function can be decomposed into two parts:  $d()$  and  $g_k()$ :

$$\hat{\omega}_k = d(g_k(\vec{X}))$$

$g(\vec{X})$ : A discriminant function:  $\mathbb{R}^D \rightarrow \mathbb{R}^K$   
 $d()$ : a decision function  $\mathbb{R}^K \rightarrow \{\omega_K\}$

The discriminant is a vector of functions:

$$\vec{g}(\vec{X}) = \begin{pmatrix} g_1(\vec{X}) \\ g_2(\vec{X}) \\ \vdots \\ g_K(\vec{X}) \end{pmatrix}$$

Quadratic discrimination functions can be derived directly from  $p(\omega_k | X)$

$$p(\omega_k | \vec{X}) = \frac{P(\vec{X} | \omega_k) p(\omega_k)}{P(\vec{X})}$$

To minimize the number of errors, we will choose  $k$  such that

$$\hat{\omega}_k = \arg\max_{\omega_k} \left\{ \frac{P(\vec{X} | \omega_k) p(\omega_k)}{P(\vec{X})} \right\}$$

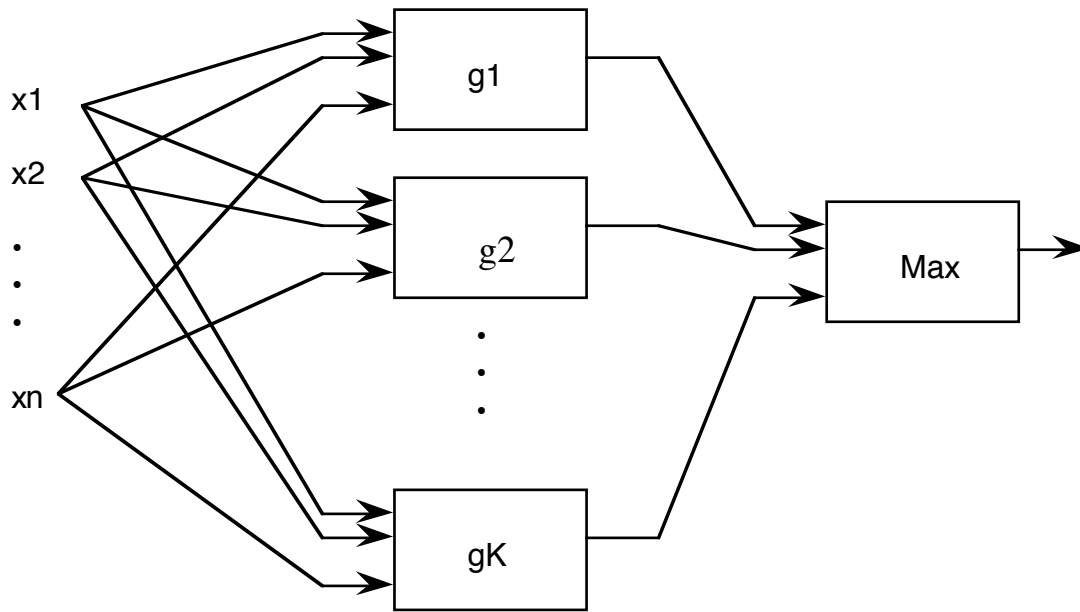
but because  $P(\vec{X})$  is constant for all  $k$ , it is common to use:

$$\hat{\omega}_k = \arg\max_{\omega_k} \{P(\vec{X} | \omega_k) p(\omega_k)\}$$

Remember that the confidence is

$$CF_{\hat{\omega}_k} = p(\hat{\omega}_k | \vec{X}) = \frac{P(\vec{X} | \hat{\omega}_k) p(\hat{\omega}_k)}{P(\vec{X})}$$

Thus the classifier can be decomposed to a selection among a set of parallel discriminant functions.



This is easily applied to the multivariate norm:

$$P(\vec{X} | \omega_k) = \mathcal{N}(\vec{X}; \vec{\mu}_k, \mathbf{C}_k)$$

**Discrimination using Log Likelihood**

As a simple example, let  $D=1$

$$P(X = x | \omega_k) = \mathcal{N}(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The discrimination function takes the form:

$$g_k(X) = P(X | \omega_k)P(\omega_k) = p(\omega_k) \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{(x-\mu_k)^2}{2\sigma_k^2}}$$

Note that  $k = \arg\max_k \{g_k(X)\} = \arg\max_k \{\text{Log}\{g_k(X)\}\}$

because  $\text{Log}\{\}$  is a monotonic function.

$$k = \arg\max_k \{\text{Log}\{p(\omega_k)\mathcal{N}(X; \mu_k, \sigma_k)\}\}$$

$$k = \arg\max_k \{\text{Log}\left\{\frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{(x-\mu_k)^2}{2\sigma_k^2}}\right\} + \text{Log}\{p(\omega_k)\}\}$$

$$k = \arg\max_k \{\text{Log}\left\{\frac{1}{\sqrt{2\pi}\sigma_k}\right\} + \text{Log}\left\{e^{-\frac{(x-\mu_k)^2}{2\sigma_k^2}}\right\} + \text{Log}\{p(\omega_k)\}\}$$

$$k = \arg\max_k \left\{-\text{Log}\{\sqrt{2\pi} \sigma_k\} - \frac{(x-\mu_k)^2}{2\sigma_k^2} + \text{Log}\{p(\omega_k)\}\right\}$$

$$k = \arg\max_k \left\{-\text{Log}\{\sigma_k\} - \frac{(x-\mu_k)^2}{2\sigma_k^2} + \text{Log}\{p(\omega_k)\}\right\}$$

**Example for  $K > 2$  and  $D > 1$**

In the general case, there are  $D$  characteristics.

$$g_k(\vec{X}) = p(\omega_k | \vec{X}) p(\omega_k)$$

Thus the classifier is a machine that calculates  $K$  functions  $g_k(\vec{X})$  Followed by a maximum selection.

The discrimination function is  $g_k(\vec{X}) = p(\vec{X} | \omega_k) p(\omega_k)$

Choose the class  $\omega_k$  for which  $\arg\text{-max}_k \{g_k(\vec{X})\}$

From Bayes rule:

$$\begin{aligned} \arg\text{-max}_k \{p(\omega_k | \vec{X})\} &= k = \arg\text{-max}_k \{p(\vec{X} | \omega_k) p(\omega_k)\} \\ &= \arg\text{-max}_k \{\text{Log}\{p(\vec{X} | \omega_k)\} + \text{Log}\{p(\omega_k)\}\} \end{aligned}$$

For a Gaussian (Normal) density function

$$p(\vec{X} | \omega_k) = \mathcal{N}(\vec{X}; \vec{\mu}_k, C_k)$$

$$\text{Log}(P(\vec{X} | \omega_k)) = \text{Log}\left\{ \frac{1}{(2\pi)^{\frac{D}{2}} \det(C_k)^{\frac{1}{2}}} e^{-\frac{1}{2}(\vec{X} - \vec{\mu}_k)^T C_k^{-1} (\vec{X} - \vec{\mu}_k)} \right\}$$

$$\text{Log}(P(\vec{X} | \omega_k)) = -\frac{D}{2} \text{Log}(2\pi) - \frac{1}{2} \text{Log}\{\text{Det}(C_x)\} - \frac{1}{2} (\vec{X} - \vec{\mu}_k)^T C_k^{-1} (\vec{X} - \vec{\mu}_k)$$

We can observe that  $-\frac{D}{2} \text{Log}(2\pi)$  can be ignored because it is constant for all  $k$ .

The discrimination function becomes:

$$g_k(\vec{X}) = -\frac{1}{2} \text{Log}\{\det(C_k)\} - \frac{1}{2} (\vec{X} - \vec{\mu}_k)^T C_k^{-1} (\vec{X} - \vec{\mu}_k) + \text{Log}\{p(\omega_k)\}$$

$$g_k(\vec{X}) = -\frac{1}{2} \text{Log}\{\det(C_k)\} - \frac{1}{2} (\vec{X} - \vec{\mu}_k)^T C_k^{-1} (\vec{X} - \vec{\mu}_k) + \text{Log}\{p(\omega_k)\}$$

Different families of Bayesian classifiers can be defined by variations of this formula. This becomes more evident if we reduce the equation to a quadratic polynomial.



**Canonical Form for the discrimination function**

The quadratic discriminant can be reduced to a standard (canonical) form.

$$g_k(\vec{X}) = -\frac{1}{2} \text{Log}\{\det(C_k)\} - \frac{1}{2} (\vec{X} - \vec{\mu}_k)^T C_k^{-1} (\vec{X} - \vec{\mu}_k) + \text{Log}\{p(\omega_k)\}$$

Let us start with the term  $(\vec{X} - \vec{\mu}_k)^T C_k^{-1} (\vec{X} - \vec{\mu}_k)$ .

This can be rewritten as :

$$(\vec{X} - \vec{\mu}_k)^T C_k^{-1} (\vec{X} - \vec{\mu}_k) = \vec{X}^T C_k^{-1} \vec{X} - \vec{X}^T C_k^{-1} \vec{\mu}_k - \vec{\mu}_k^T C_k^{-1} \vec{X} + \vec{\mu}_k^T C_k^{-1} \vec{\mu}_k$$

We note that  $\vec{X}^T C_k^{-1} \vec{\mu}_k = \vec{\mu}_k^T C_k^{-1} \vec{X}$

and thus :  $-\vec{X}^T C_k^{-1} \vec{\mu}_k - \vec{\mu}_k^T C_k^{-1} \vec{X} = -(2C_k^{-1} \vec{\mu}_k)^T \vec{X}$

we define:  $\vec{W}_k = -2C_k^{-1} \vec{\mu}_k$

to obtain  $-\vec{X}^T C_k^{-1} \vec{\mu}_k - \vec{\mu}_k^T C_k^{-1} \vec{X} = \vec{W}_k^T \vec{X}$

Let us also define  $D_k = -\frac{1}{2} C_k^{-1}$

The remaining terms are constant. Let us defined the constant

$$b_k = -\frac{1}{2} \vec{\mu}_k^T C_k^{-1} \vec{\mu}_k - \text{Log}\{\det(C_k)\} + \text{Log}\{p(\omega_k)\}$$

which gives a quadratic polynomial

$$g_k(\vec{X}) = \vec{X}^T D_k \vec{X} + \vec{W}_k^T \vec{X} + b_k$$

where:  $D_k = -\frac{1}{2} C_k^{-1}$

$$\vec{W}_k = -2C_k^{-1} \vec{\mu}_k$$

and  $b_k = -\frac{1}{2} \vec{\mu}_k^T C_k^{-1} \vec{\mu}_k - \text{Log}\{\det(C_k)\} + \text{Log}\{p(\omega_k)\}$

A set of  $K$  discrimination functions  $g_k(\vec{X})$  partitions the space  $\vec{X}$  into a disjoint set of regions with quadratic boundaries. The boundaries are points for which

$$g_i(\vec{X}) = g_j(\vec{X}) \geq g_k(\vec{X}) \quad \forall k \neq i, j$$

The boundaries are the functions  $g_i(\vec{X}) - g_j(\vec{X}) = 0$

## Noise and Discrimination

Under certain conditions, the quadratic discrimination function can be simplified by eliminating either the quadratic or the linear term.

If we could perfectly model the universe, then sensor reading would be a predictable value,  $\bar{x}$ . The normal density attempts to represent this with the "average" feature  $\bar{\mu}_k$ .

In reality, the features of a class are generally dispersed by un-modeled phenomena. These may be effects that are beyond the abilities of the available sensors, or they may be effects that we choose to ignore because they are "unimportant".

Although the true variation may not be additive, we will model it as an additive random term  $N_k$ . The term is random because we are unable to predict it.

Thus the observed feature is random:  $\vec{X} = \bar{x} + N_k$

For example, the color of your eyes could be predicted from your genetic code, but in the absence of a genetic decoder, this becomes random.

In addition, every observation system (or sensor) is subject to some form of sensor noise. This sensor Noise is modeled as an additive random term  $N_s$ . Sensor noise is generally independent of the class  $k$ .

Thus the sensor returns a random feature  $\vec{X} = \bar{x} + \vec{N}_k + \vec{N}_s$

The Normal density function represents these two forms of "noise" as a second moment of the class,  $C_k$ .

Thus  $\Sigma_k = E\{E\{(N_k + N_s)(N_k + N_s)^T\}$

Depending on the nature of  $\vec{N}_k$  and  $\vec{N}_s$  different simplifications are possible.

For example if  $\vec{N}_s \gg \vec{N}_k$  then the term  $\Sigma_k$  will be nearly constant for all  $k$ . In this case, the discrimination function can be reduced to a linear equation.

$$g_k(\vec{X}) = \vec{W}_k^T \vec{X} + b_k$$

This is very useful because there are simple powerful techniques to calculate the terms of such an equation.

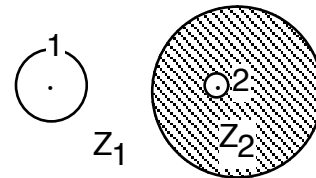
**Decision Surfaces for different Noise assumptions**

In the more general case we can not make any assumptions on  $\vec{N}_k$  and  $\vec{N}_s$   
 Depending on the nature  $\vec{N}_k$  we may find a variety of different second order decision surfaces :

For example (K=2, D=2)

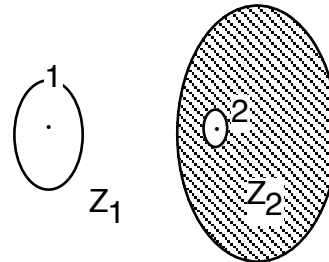
Hyper-sphere :

Let  $\Sigma_k = \sigma_k^2 I$   
 and  $\det\{\Sigma_1\} > \det\{\Sigma_2\}$



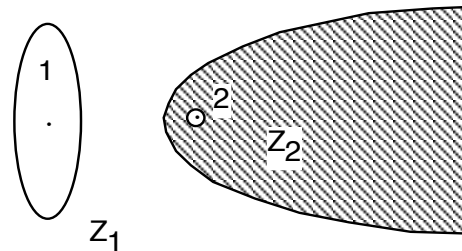
Hyper-ellipsoid :

For  $\sigma_{x1}^2 > \sigma_{x2}^2$   
 and  $\det\{\Sigma_1\} > \det\{\Sigma_2\}$

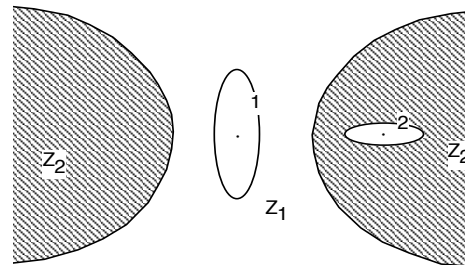


Hyper-paraboloid :

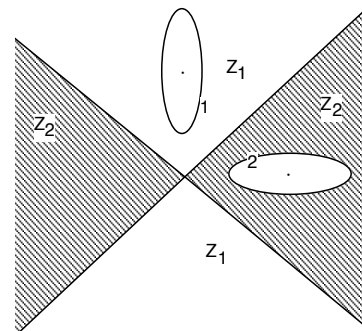
for  $\sigma_{x1k=1}^2 \gg \sigma_{x1k=2}^2$   
 et  $\sigma_{x2k=1}^2 > \sigma_{x2k=2}^2$



Hyper-hyperboloids :

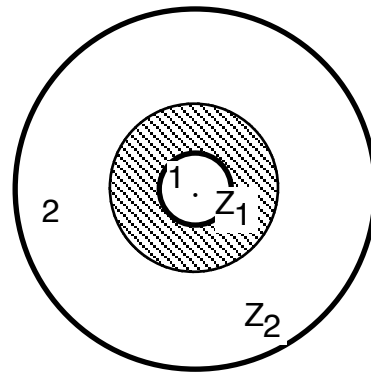


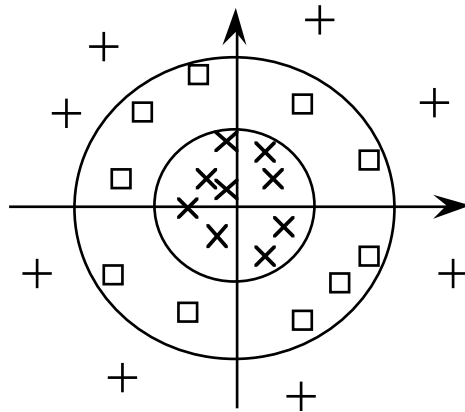
Hyperplanes.



$\vec{\mu}_1 = \vec{\mu}_2$  and  $\det\{\Sigma_1\} > \det\{\Sigma_2\}$   
with  $\sigma_{11} = \sigma_{22}$  et  $\sigma_{12} = \sigma_{21} = 0$ .

a hypersphere.



**Two classes with equal means**

Suppose tht we have 2 classes  $i, j$  such that

$$\vec{\mu}_i = \vec{\mu}_j \text{ and } \det\{\Sigma_1\} > \det\{\Sigma_2\}.$$

Is it possible to assign an observation to one of the classes?

$$g_i(\vec{X}) - g_j(\vec{X}) = 0$$

takes the form of a sphere with observations assigned to  $C_i$  outside the sphere and  $C_j$  on the inside.

$$g_k(\vec{X}) = \vec{X}^T D_k \vec{X} + \vec{W}_k^T \vec{X} + b_k$$