

# Vision par Ordinateur

James L. Crowley

DEA IVR

Premier Bimestre 2005/2006

Séance 5

9 novembre 2005

## Reconnaissance de formes

### Plan de la Séance :

La Classification des Formes .....	2
Les Formes.....	3
Les Observations.....	4
La Loi Normale.....	6
Estimations des moments d'une densité .....	7
La Loi Normale pour $D = 1$ .....	9
La Loi Normale pour $D > 1$ .....	10
La Classification.....	13
Classification entre deux Catégories ( $K=2, D=1$ ).....	13
La Probabilité d'erreur.....	14
Classification pour $K > 2$ .....	15
Classification pour $K > 2$ et $D > 1$ .....	16
Forme Canonique de la fonction de discrimination.....	17

### **Notations**

$x$	Un vecteur
$X$	Un vecteur aléatoire (non-prévisible).
$D$	Nombre de dimensions de $X$
$w_k$	La classe $k$
$k$	Indice d'une classe
$K$	Nombre de classes
$M_k$	Nombre d'exemples de la classe $k$ .
$M$	Nombre totale d'exemples de toutes les classes
$p(w_k)$	Probabilité a priori de rencontrer un membre de la classe $k$ .
$Y$	Une observation (un vecteur aléatoire).
$P(Y)$	Probabilité d'une observation $Y$

## La Classification des Formes

La classification est une capacité fondamentale de l'intelligence.

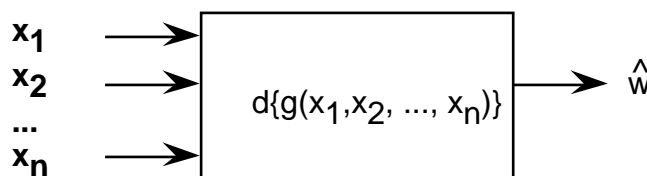
Reconnaissance : Le fait de reconnaître, d'identifier un objet, un être comme tel.

Identifier : Reconnaître un individu

Classer : Reconnaître un membre d'une catégorie, ou d'une classe.

Un ensemble est défini par un test d'appartenance.

La classification est un processus d'association d'une observation à une classe par un teste d'appartenance.



Pour un vecteur de caractéristique il sort une estimation de la classe,  $\hat{w}$

Les techniques de reconnaissance de formes statistiques fournissent une méthode pour induire des tests d'appartenance à partir d'un ensemble d'échantillons.

La classification se résume à une division de l'espace de caractéristique en partition disjoint. Cette division peut-être fait par estimation de fonctions paramétrique ou par une liste exhaustives des frontières.

Le critère est la probabilité d'appartenance.

Cette probabilité est fournie par la règle de Bayes.

$$p(w_k | X) = \frac{p(X | w_k) p(w_k)}{p(X)}$$

**Les Formes**

Forme n. f. : A. Apparence, aspect visible. 1) ... 2) apparence extérieure donnant à un objet ou à un être sa spécificité.

Les méthodes statistique de la reconnaissance de forme traite les observations sous forme de vecteur de caractéristiques.

Caractéristiques : (En anglais : Feature) Signes ou ensembles de signes distinctifs.  
Une ensemble de propriétés.  $\{ x_1, x_2 \dots x_D \}$ .

En notation vectorielle :

$$X = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_D \end{pmatrix}$$

La formation des vrais objets physiques est sujette aux influences aléatoires. Pour les objets d'une classe,  $w_k$ , les propriétés des objets individuels sont, les valeurs aléatoires. On peut resumer ceci par une somme d'une forme "intrinsèque"  $x$  plus ces influences aléatoires individuelles,  $B_i$ .

$$X = x + B_i$$

Les techniques probabiliste de RF suppose un bruit additif.

En notation vectorielle :

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \dots \\ X_n \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + \begin{pmatrix} B_1 \\ B_2 \\ \dots \\ B_n \end{pmatrix}$$

## Les Observations

Une observation : une constatation attentive des phénomènes.

Pour des machines, des observations sont fournies par les capteurs.

Ceci donne une observation (un phénomène) sous forme d'une ensemble de caractéristiques :  $\{ Y_1, Y_2 \dots Y_D \}$ .

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{pmatrix}$$

Les observations sont corrompues par un bruit,  $B_o$ .

$$Y = y + B_o$$

Le bruit est, par définition, imprévisible. Il est aléatoire.

Donc les caractéristiques observées sont des vecteurs aléatoires.

La corruption des observations par un bruit aléatoire est fondamentale aux capteurs physiques.

Pour chaque classe  $w_k$ , la probability d'observé  $Y$  est fournie par la règle de Bayes.

$$p(w_k | Y) = \frac{p(Y | w_k) p(w_k)}{p(Y)}$$

Si  $Y = X + B_o$  est issu de la classe  $w_k$  ayant caractéristique  $X = x + B_i$

Par exemple :

$$L = R+G+B \quad C_1 = \frac{R}{R+G+B} \quad C_2 = \frac{G}{R+G+B}$$

$R, G, B$  sont les entiers. Donc,  $C_1, C_2$  sont issu d'une ensemble finit de valeurs dans l'intervalle  $[0, 1]$ . On peut transformer  $C_1, C_2$  en entier entre  $[0, N-1]$ , par

$$C_1 = \text{Round} \left( N \cdot \frac{R}{R+G+B} \right). \quad C_2 = \text{Round} \left( N \cdot \frac{G}{R+G+B} \right).$$

On aura  $N^2$  cellules de chrominances dans l'histogramme.

Par exemple, pour  $N=32$ , on a  $32^2 = 1024$  cellules à remplir est il nous faut que  $M = 10$  K pixels d'exemples. (Une image = 256 K pixels).

Dans ce cas, pour M observations  $p(X) = \frac{1}{M} h(X)$

La probabilité à posteriori peut être calculé par la règle de Bayes.

$$p(w_k | X) = \frac{p(X | w_k) p(w_k)}{p(X)}$$

Dans le cas des valeurs de X discrètes tel que  $x \in [X_{\min}, X_{\max}]$ , on a

probabilité de la classe  $w_k$ :  $p(w_k) = \frac{M_k}{M}$

probabilité conditionnelle de X):  $p(X | w_k) = \frac{1}{M_k} h_k(X)$

Probabilité à priori de X :  $p(X) = \frac{1}{M} h(X)$

ce qui donne :

$$p(w_k | \vec{X}) = \frac{p(\vec{X} | w_k) p(w_k)}{p(\vec{X})} = \frac{\frac{M_k}{M} \frac{1}{M_k} h_k(X)}{\frac{1}{M} h(X)} = \frac{h_k(X)}{h(X)}$$

Que faire si la masse d'exemple est insuffisante :  $M < 10 (X_{\max} - X_{\min})$  ?

Que faire si x n'est pas entier ? Il faut une fonction paramétrique pour  $p(X)$ .

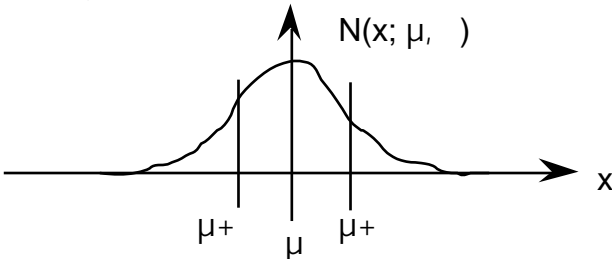
## La Loi Normale

Quand les variables aléatoires sont issues d'une séquence d'événements aléatoires, leur densité de probabilité prend la forme de la loi normale,  $\mathcal{N}(\mu, \sigma^2)$ . Ceci est démontré par le théorème de la limite centrale. Il est un cas fréquent en nature.

Les paramètres de  $\mathcal{N}(\mu, \sigma^2)$  sont les premiers et deuxième moments des exemples. Donc, on peut les estimer pour n'importe quel nombre d'exemples. On peut même estimer les moments quand il n'existe pas les bornes ( $X_{\max}-X_{\min}$ ) ou quand X est une variable continue.

Dans ce cas,  $p(\cdot)$  est une "densité" et il faut une fonction paramétrique pour  $p(\cdot)$ .

Dans la plupart des cas, on peut utiliser  $\mathcal{N}(\mu, \sigma^2)$  comme une fonction de densité pour  $p(x)$ .

$$p(x) \quad \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$


The graph shows a bell-shaped curve representing the normal distribution function  $N(x; \mu, \sigma^2)$ . The horizontal axis is labeled  $x$  and the vertical axis is labeled  $N(x; \mu, \sigma^2)$ . The curve is centered at  $\mu$ . Two vertical lines are drawn at  $\mu + \sigma$  and  $\mu - \sigma$ , indicating the standard deviation from the mean.

Le base "e" est :  $e = 2.718281828\dots$ . Il s'agit du fonction tel que  $\int e^x dx = e^x$

Le terme  $\frac{1}{\sqrt{2\pi}}$  sert à normaliser la fonction en sorte que sa surface est 1.

$$\int e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \sqrt{2\pi} \cdot \sigma$$

Le terme  $d^2(x) = \frac{(x-\mu)^2}{2\sigma^2}$  est la distance  $x$  et  $\mu$  normalisée par la variance.

La différence  $(x - \mu)^2$  est la "distance" entre une caractéristique et la caractéristique "nominale" d'une classe. La variance,  $\sigma^2$ , sert à "normaliser" cette distance.

La différence normalisée par la variance est connue sous le nom de "Distance de Mahalanobis". La Distance de Mahalanobis est un test naturel de similarité

**Estimations des moments d'une densité**Le premier moment : La Moyenne

Soit  $M$  observations d'un variable aléatoire,  $\{X_1, X_2, \dots, X_M\} = \{X_m\}$

La moyenne est l'espérance de  $\{X_m\}$ .

$$\mu = E\{x\} = \frac{1}{M} \sum_{m=1}^M X_m$$

Il s'agit d'une somme sur  $M$  (le nombre exemples). Cette somme existe, même quand  $X_{\min}$  et  $X_{\max}$  n'existent pas et quand  $X$  est réelle.

On note que dans le cas  $X$  est un nombre entier, on peut aussi estimer la moyenne par la table de fréquence. La masse d'un histogramme,  $h(x)$  est le nombre d'échantillons qui composent l'histogramme,  $M$ .

$$M = \sum_{x=X_{\min}}^{x_{\max}} h(x)$$

Pour  $X$  entier, tel que  $X \in [x_{\min}, x_{\max}]$  on peut démontrer que

$$\mu = E\{X\} = \frac{1}{M} \sum_{x=X_{\min}}^{x_{\max}} h(x) \cdot x = \sum_{x=X_{\min}}^{x_{\max}} p(x) \cdot x$$

$$\text{donc : } \mu = E\{X\} = \frac{1}{M} \sum_{m=1}^M X_m = \frac{1}{M} \sum_{x=X_{\min}}^{x_{\max}} h(x) \cdot x$$

Pour  $X$  continue :

$$\mu = E\{X\} = \int p(x) \cdot x \, dx$$

Le deuxième moment (La variance)

La variance  $\sigma^2$  est le deuxième moment de la densité de probabilité.

Pour un ensemble de  $M$  observations  $\{x_m\}$

$$\sigma^2 = E\{(X_m - \mu)^2\} = \frac{1}{M} \sum_{m=1}^M (X_m - \mu)^2$$

Mais l'usage de  $\mu$  estimé avec le même ensemble, introduit un biais dans  $\sigma^2$ .

Pour l'éviter, on peut utiliser une estimation sans biais.

$$\sigma^2 = \frac{1}{M-1} \sum_{m=1}^M (X_m - \mu)^2$$

Lequel est correct ? (les deux !)

Pour  $X$  entier, tel que  $X \in [X_{\min}, X_{\max}]$  on peut démontrer que

$$\sigma^2 = E\{(X_m - \mu)^2\} = \frac{1}{M} \sum_{x=X_{\min}}^{X_{\max}} h(x)(x - \mu)^2$$

Ceci est vrai par ce que la table  $h(x)$  est faite de  $\{X_m\}$ .

Donc :

$$\sigma^2 = \frac{1}{M} \sum_{m=1}^M (X_m - \mu)^2 = \frac{1}{M} \sum_{x=X_{\min}}^{x_{\max}} h(x)(x - \mu)^2$$

Pour  $X$  réel on a

$$\sigma^2 = E\{(X_m - \mu)^2\} = \int p(x) \cdot (x - \mu)^2 dx$$

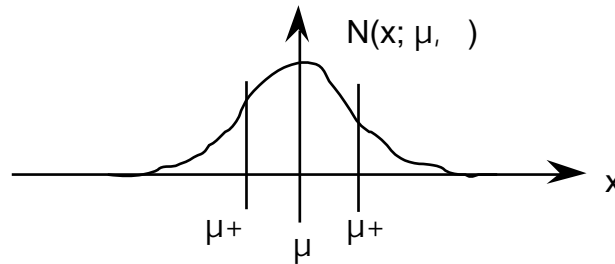


**La Loi Normale pour  $D = 1$** 

Avec  $\mu$  et  $\sigma^2$ , on peut estimer la densité  $p(x)$  par  $\mathcal{N}(x; \mu, \sigma^2)$

$$\text{pr}(X=x) = p(x) = \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\mathcal{N}(x; \mu, \sigma^2)$  a la forme :



La moyenne est le premier moment de la densité  $p(x)$ .

$$\mu = E\{X\} = \int p(x) \cdot x \, dx$$

La variance  $\sigma^2$  est le deuxième moment de  $p(x)$ .

$$\sigma^2 = E\{(X-\mu)^2\} = \int p(x) \cdot (x-\mu)^2 \, dx$$

**La Loi Normale pour  $D > 1$** 

Pour un vecteur de  $D$  propriétés

$$\mu = E\{\vec{X}\} = \frac{1}{M} \sum_{m=1}^M X_m = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_D \end{pmatrix} = \begin{pmatrix} E\{X_1\} \\ E\{X_2\} \\ \dots \\ E\{X_D\} \end{pmatrix}$$

Pour  $X$  entier, tel que pour chaque  $d \in [1, D]$ ,  $x_d \in [x_{dmin}, x_{dmax}]$  on peut démontrer que

$$\mu_d = E\{x_d\} = \frac{1}{M} \sum_{x_1=x_{1min}}^{x_{1max}} \dots \sum_{x_D=x_{Dmin}}^{x_{Dmax}} h(x) x_d$$

Pour  $x$  réel,  $\mu_d = E\{x_d\} = \dots \int p(x) \cdot x_d dX$

Dans tous les cas :

$$\mu = E\{\vec{X}\} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_n \end{pmatrix} = \begin{pmatrix} E\{x_1\} \\ E\{x_2\} \\ \dots \\ E\{x_n\} \end{pmatrix}$$

Pour  $D$  dimensions, la covariance entre les variables  $x_i$  et  $x_j$  est estimée à partir de  $M$  observations  $\{X_m\}$

$$\begin{aligned} \sigma_{ij}^2 &= E\{ (X_i - E\{X_i\})(X_j - E\{X_j\}) \} \\ &= \frac{1}{M} \sum_{m=1}^M (X_{im} - \mu_i)(X_{jm} - \mu_j) \end{aligned}$$

Et encore, pour éviter le bias, on peut utiliser :

$$\sigma_{ij}^2 = \frac{1}{M-1} \sum_{m=1}^M (X_{im} - \mu_i)(X_{jm} - \mu_j)$$

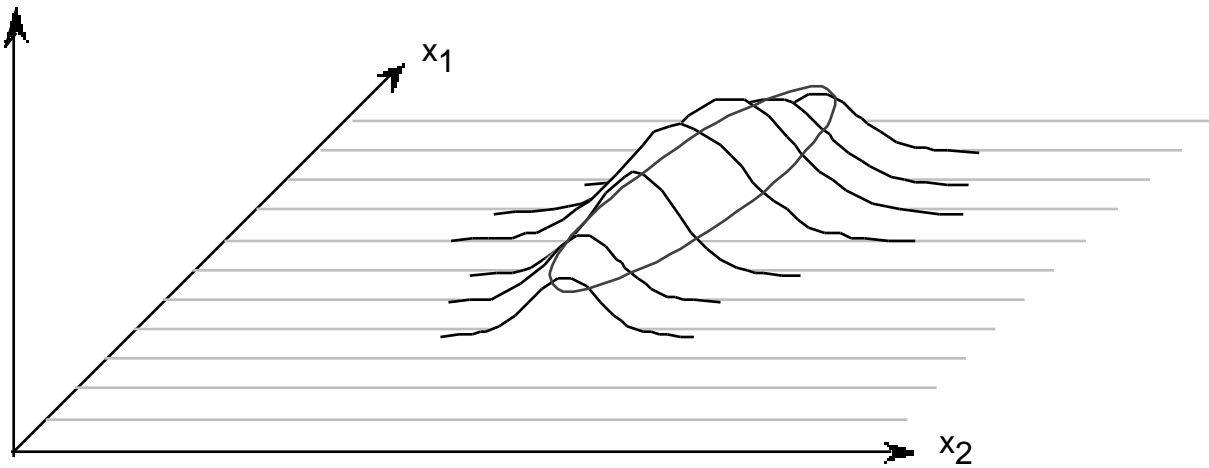
Ces coefficients composent une matrice de covariance.  $C$

$$C_x = E\{(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T\} = E\{(\mathbf{X} - E\{\mathbf{X}\})(\mathbf{X} - E\{\mathbf{X}\})^T\}$$

$$C_x = \begin{pmatrix} 11^2 & 12^2 & \dots & 1D^2 \\ 21^2 & 22^2 & \dots & 2D^2 \\ \dots & \dots & \dots & \dots \\ D1^2 & D2^2 & \dots & DD^2 \end{pmatrix}$$

Dans le cas d'un vecteur de propriétés, X, la loi normale prend la forme :

$$p(\mathbf{X}) = \mathcal{N}(\mathbf{X}; \boldsymbol{\mu}, C) = \frac{1}{(2\pi)^{\frac{D}{2}} \det(C)^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{X} - \boldsymbol{\mu})^T C_x^{-1} (\mathbf{X} - \boldsymbol{\mu})}$$



Le terme  $\frac{1}{(2\pi)^{\frac{D}{2}} \det(C)^{\frac{1}{2}}}$  est un facteur de normalisation.

$$\dots e^{-\frac{1}{2}(\mathbf{X} - \boldsymbol{\mu})^T C_x^{-1} (\mathbf{X} - \boldsymbol{\mu})} dX_1 dX_2 \dots dX_D = \frac{1}{(2\pi)^{\frac{D}{2}} \det(C)^{\frac{1}{2}}}$$

La déterminante,  $\det(C)$  est une opération qui donne la "énergie" de C.

Pour D=2  $\det \begin{pmatrix} a & b \\ c & d \end{pmatrix} = a \cdot d - b \cdot c$

Pour D=3

$$\det \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = a \cdot \det \begin{pmatrix} e & f \\ h & i \end{pmatrix} + b \cdot \det \begin{pmatrix} f & d \\ i & g \end{pmatrix} + c \cdot \det \begin{pmatrix} d & e \\ g & h \end{pmatrix}$$

$$= a(ei - fh) + b(fg - id) + c(dh - eg)$$

pour D > 3 on continue récursivement.

L'exposant est une valeur positive et quadrique.

(si  $X$  est en mètre,  $\frac{1}{2} (\mathbf{X} - \boldsymbol{\mu})^T \mathbf{C}_x^{-1} (\mathbf{X} - \boldsymbol{\mu})$  est en mètre<sup>2</sup>.)

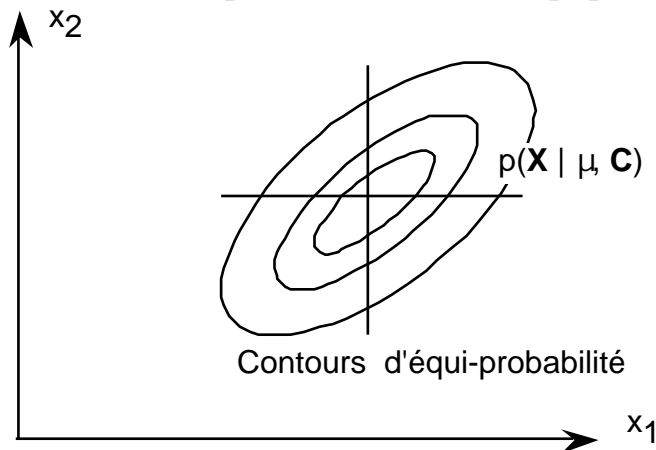
Cette valeur est connue comme la "distance de Mahalanobis".

$$d^2(\mathbf{X}) = \frac{1}{2} (\mathbf{X} - \boldsymbol{\mu})^T \mathbf{C}_x^{-1} (\mathbf{X} - \boldsymbol{\mu})$$

Il s'agit d'une distance euclidienne, normalisé par la covariance  $\mathbf{C}_x$ .

Cette distance est bien définie, même si les composants de  $X$  n'ont pas les mêmes unités. (Ceci est souvent le cas).

La loi Normale peut être visualisé par ses contours d'"équiprobabilité"



Ces contours sont les contours de constant  $d^2(\mathbf{X})$

La matrice  $C$  est positive et semi-définie. Nous allons nous limiter au cas où  $C$  est positive et définie ( $C$ -à-d.  $\det(C) = |C| > 0$ )

si  $x_i$  et  $x_j$  sont statistiquement indépendants,  $c_{ij} = 0$ .

si pour tout  $i \neq j$ ,  $c_{ij} = 0$  alors  $p(\mathbf{X}) = \prod_{i=1}^n p(x_i)$

## La Classification

### Classification entre deux Catégories (K=2, D=1)

Soit un événement A avec propriété X.  
Soit  $w_1$  une tribu d'événements.

On note que  $A \in w_1$  ou  $A \notin w_1$

soit la proposition p que X vaut la valeur "x".

$p : X = x$ . on écrit :  $p(X) = \Pr\{X=x\}$

Soit une classe  $w_1$ . Soit  $A \in w_1$  ou  $A \notin w_1$   
soit proposition q :  $A \in w_1$  donc  $\neg q : A \notin w_1$

on écrit :  $p(q) = \Pr\{A \in w_1\} = p(w_1)$

Bayes :  $p(q | p) p(p) = p(p | q) p(q)$

$\Pr(A \in w_1 | X=x) \Pr\{X=x\} = \Pr(X | A \in w_1) \Pr(A \in w_1) =$

ou bien :  $p(w_1 | X) p(X) = p(X | w_1) p(w_1)$

Ceci donne :

$$p(w_1 | X) = \frac{p(X | w_1) p(w_1)}{p(X)}$$

Soit une deuxième classe  $w_2$  pour tous les événements qui n'est pas de  $w_1$ .

$A \in w_1$     $A \in w_2$

Les classes  $w_1$  et  $w_2$  sont mutuellement exclusifs.

Soit x. Comment décider entre  $w_1$  et  $w_2$  ?

Une idée simple serait de chercher à minimiser la probabilité d'erreur.

$$p(\text{Erreur} | X) = \begin{cases} p(w_1 | X) = \Pr\{A \in w_1 | X=x\} & \text{si on décide } w_2 \\ p(w_2 | X) = \Pr\{A \in w_2 | X=x\} & \text{si on décide } w_1 \end{cases}$$

Donc pour tout  $X$  :  $p(\text{Erreur} | X)$  est minimale si on décide le  $w_k$  tel que

si  $p(w_1 | X) > p(w_2 | X)$  décide  $w_1$  sinon décide  $w_2$

Par règle de Bayes ceci est

si  $p(X | w_1) p(w_1) > p(X | w_2) p(w_2)$  décide  $w_1$  sinon décide  $w_2$

L'idée intuitive est de voir les valeurs de la propriété  $x$  comme un espace. Pour une scalaire  $x$  ( $D = 1$ ) l'espace est une droite. (Une droite est un ensemble ordonné de points.)

La frontière  $p(x | w_1) p(w_1) = p(x | w_2) p(w_2)$  partitionne l'espace  $x$  dans deux régions disjointes,  $z_1$  et  $z_2$ .

### La Probabilité d'erreur

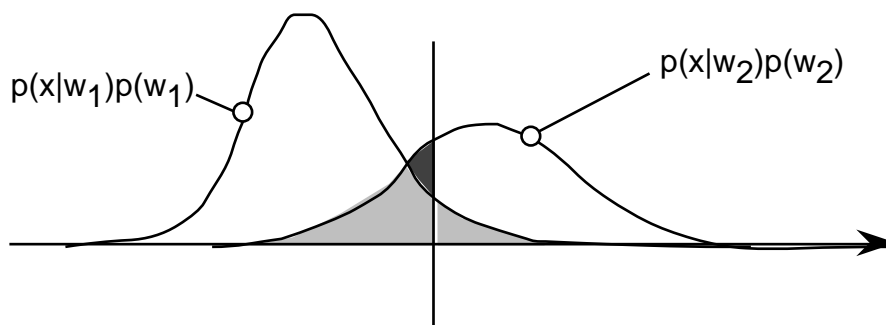
Notre règle

si  $p(X | w_1) p(w_1) > p(X | w_2) p(w_2)$  décide  $w_1$  sinon décide  $w_2$

Notre surface de décision est  $p(X | w_1) p(w_1) = p(X | w_2) p(w_2)$

On peut noter que

$$\begin{aligned} p(\text{Erreur} | X) &= p(X \in z_2, w_1) p(w_1) + p(X \in z_1, w_2) p(w_2) \\ &= \int_{z_2} p(x | w_1) p(w_1) dx + \int_{z_1} p(x | w_2) p(w_2) dx \end{aligned}$$



La valeur minimum de  $P(\text{erreur})$  est atteinte pour  $\frac{d P(\text{erreur})}{dx} = 0$

Donc quand  $p(x | w_1) p(w_1) - p(x | w_2) p(w_2) = 0$ .

**Classification pour  $K > 2$** 

Soit  $D=1$ . Dans le cas de  $K > 2$ , il y a plus de possibilité d'erreurs.

Il vaut mieux maximiser  $P(\text{vrai} | X) = p(\neg \text{Erreur} | X)$

Dans le cas général, on sélection la classe  $w_k$  pour laquelle

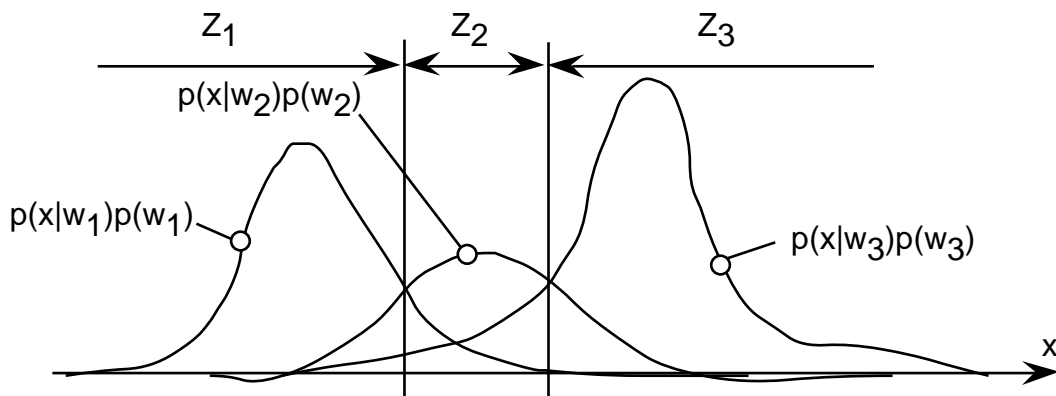
$$k = \underset{k}{\text{arg-max}} \{g_k(X)\} \quad \text{avec } g_k(X) = p(w_k | X)$$

par règle de Bayes :

$$\underset{k}{\text{arg-max}} \{p(w_k | X)\} = \underset{k}{\text{arg-max}} \left\{ \frac{p(X | w_k) p(w_k)}{p(X)} \right\}$$

Les frontières entre régions  $i$  et  $j$  sont les valeurs pour lesquelles

$$g_i(X) = g_j(X)$$



Une fonction de discrimination partition l'espace de caractéristique en régions disjointes  $Z_1, \dots, Z_k$  pour chaque classe.

Classification :

$$w_k = \underset{k}{\text{arg-max}} \{g_k(X)\}$$

**Classification pour  $K > 2$  et  $D > 1$ .**

Dans le cas général,  $g_k(X) = p(w_k | X) p(w_k)$

Soit  $g_k(X) = p(w_k | X) p(w_k)$ .

Quand l'observation est un vecteur ( $D > 1$ ) on aura

$$\text{si } i > j \quad g_i(X) > g_j(X) \text{ décide } w_i$$

Dans cette forme le classificateur est une machine qui calcule  $K$  fonctions  $g_k(x)$  suivie d'une sélection du maximum.

La fonction de discrimination est :  $g_k(X) = p(w_k | X) = \frac{p(X | w_k) p(w_k)}{p(X)}$

On sélectionne la classe  $w_k$  pour laquelle  $\arg\text{-max}_k \{g_k(X)\}$

par règle de Bayes :

$$\arg\text{-max}_k \{p(w_k | X)\} = \arg\text{-max}_k \left\{ \frac{p(X | w_k) p(w_k)}{p(X)} \right\}$$

=  $\arg\text{-max}_k \{p(X | w_k) p(w_k)\}$  parce que  $p(x)$  est indépendant de  $w_k$

=  $\arg\text{-max}_k \{ \text{Log}\{p(X | w_k)\} + \text{Log}\{p(w_k)\} \}$

parce que  $\text{Log}\{\}$  est une fonction monotone.

Si le bruit est d'une densité Normale :

$$p(X | w_k) = \mathcal{N}(\mu_k, C_k)$$

$$\text{Log}\{p(X | w_k)\} = \text{Log}\left\{ \frac{1}{(2\pi)^{\frac{D}{2}} \det(C_k)^{\frac{1}{2}}} e^{-\frac{1}{2}(X - \mu_k)^T C_k^{-1} (X - \mu_k)} \right\}$$

$$\text{Log}\{p(X | w_k)\} = -\frac{D}{2} \text{Log}\{2\pi\} - \frac{1}{2} \text{Log}\{\det(C_k)\} - \frac{1}{2}(X - \mu_k)^T C_k^{-1} (X - \mu_k)$$

On note que  $-\frac{D}{2} \text{Log}\{2\pi\}$  peut être éliminé parce qu'il est constant pour tout  $k$ .



La fonction de discrimination devient :

$$g_k(x) = -\frac{1}{2} \text{Log}\{\det(C_k)\} - \frac{1}{2}(X - \mu_k)^T C_k^{-1} (X - \mu_k) + \text{Log}\{p(w_k)\}$$

Les classifieurs Bayésiennes sont définies par les variations de ce formula.

### Forme Cannonique de la fonction de descrimination

La classe  $w_k$  est celle qui donne un maximum pour

$$g_k(X) = -\frac{1}{2}(X - \mu_k)^T C_k^{-1} (X - \mu_k) + \frac{1}{2} \text{Log}\{\det(C_k)\} + \text{Log}\{p(w_k)\}$$

On peut réécrire  $(X - \mu_k)^T C_k^{-1} (X - \mu_k)$  comme

$$X C_k^{-1} X - X C_k^{-1} \mu_k - \mu_k^T C_k^{-1} X + \mu_k^T C_k^{-1} \mu_k$$

On note que  $C_k^{-1}$  est symétrique, et donc  $X C_k^{-1} \mu_k = \mu_k^T C_k^{-1} X$

Donc  $-X C_k^{-1} \mu_k - \mu_k^T C_k^{-1} X = 2(\mu_k^T C_k^{-1})^T X = 2(C_k^{-1} \mu_k)^T X$

On peut réécrire  $g_k(X)$  comme

$$g_k(X) = -X^T \left(\frac{1}{2} C_k^{-1}\right) X + (C_k^{-1} \mu_k)^T X - \frac{1}{2} (\mu_k^T C_k^{-1} \mu_k) - \frac{1}{2} \text{Log}\{\det(C_k)\} + \text{Log}\{p(w_k)\}$$

ou bien

$$g_k(X) = X^T (D_k) X + d_k^T X + d_{k0}.$$

avec

$$D_k = \frac{1}{2} C_k^{-1} \quad d_k = C_k^{-1} \mu_k$$

$$d_{k0} = -\frac{1}{2} (\mu_k^T C_k^{-1} \mu_k) - \frac{1}{2} \text{Log}\{\det(C_k)\} + \text{Log}\{p(w_k)\}$$

Cette fonction est composé de trois termes :

	une terme quadrique	$X^T (D_k) X$ ,
	une terme linéaire :	$d_k^T X$
et	une terme constant :	$d_{k0}$