

# Systemes Intelligents : Raisonnement et Reconnaissance

James L. Crowley

Deuxième Année ENSIMAG

Deuxième Semestre 2008/2009

Séance 9

29 avril 2009

## Reconnaissance Bayesienne

Notations .....	2
La Classification.....	3
La Classification Bayesienne.....	4
Les Variations Aléatoires des Caractéristiques.....	6
La Loi Normale.....	9
La Loi Normale pour $D = 1$ .....	11
La Loi Normale pour $D > 1$ .....	12
Forme en Algebre Linéaire.....	18
Transformations Linéaire.....	19
Fonctions de Discrimination .....	20
Discrimination.....	20
Simplification de la fonction de Discrimination.....	22
Classification pour $K > 2$ et $D > 1$ .....	23
Forme Canonique de la fonction de discrimination.....	25

Sources Bibliographiques :

"Neural Networks for Pattern Recognition", C. M. Bishop, Oxford Univ. Press, 1995.

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

**Notations**

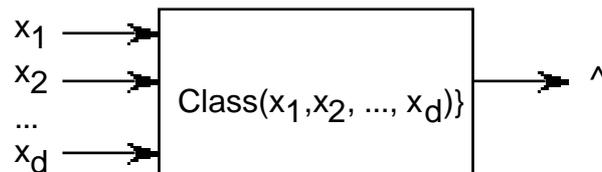
$x$	Une variable
$X$	Une valeur aléatoire (non-prévisible).
$N$	Le nombre de valeurs possible pour $x$ ou $X$
$x$	Un vecteur de $D$ variables
$X$	Un vecteur aléatoire (non-prévisible).
$D$	Nombre de dimensions de $x$ ou $X$
$E$	Une événement, un entité
$A, B$	des classes d'événements ou d'entités
$T_k$	La classe $k$ d'événements ou d'entités
$k$	Indice d'une classe
$K$	Nombre de classes
$\pi_k$	La proposition que l'entité $E \in T_k$
$M_k$	Nombre d'exemples de la classe $k$ .
$M$	Nombre totale d'exemples de toutes les classes
	$M = \sum_{k=1}^K M_k$
$h(x)$	Histogrammes des valeurs ( $x$ est entières avec range limité)
$h_k(x)$	Histogramme des valeurs pour la class $k$ .
	$h(x) = \sum_{k=1}^K h_k(x)$
$Q$	Nombre de Cellules dans $h(x)$ . $Q = N^D$
$p(\pi_k) = p(E \in T_k)$	Probabilité que $E$ est un membre de la classe $k$ .
$Y$	La valeur d'une observation (un vecteur aléatoire).
$P(X)$	Densité de Probabilité pour $X$
$p(X = x)$	Probabilité q'un vecteur $X$ prendre la valeur $x$
$P(X   \pi_k)$	Densité de Probabilité pour $X$ etant donné que $\pi_k$
	$P(X) = \sum_{k=1}^K p(X   \pi_k) p(\pi_k)$

## La Classification

Soit l'entité E décrit par un vecteur de caractéristiques  $X : (E, X)$ .

Soit K classes d'entité  $\{T_k\} = \{T_1, T_2, \dots, T_K\}$

La classification est un processus d'estimation de l'appartenance de l'entité E à une des classes  $T_k$  fondée sur les caractéristiques de l'événement, X.



$$\hat{k} = \text{Decider}(E, k)$$

$\hat{k}$  est la proposition que  $(E \in T_k)$ .

La fonction de classification est composée de deux parties  $d()$  et  $g_k()$ :

$$\hat{k} = d(g(X)).$$

$g(X)$  : Une fonction de discrimination :  $\mathbb{R}^D \rightarrow \mathbb{R}^K$

$d()$  : Une fonction de décision :  $\mathbb{R}^K \rightarrow \{1, \dots, K\}$

## La Classification Bayésienne

La technique Bayésienne de Classification repose sur une fonction de vérité probabiliste et le règle de Bayes.

Dans un système de vérité probabilisté, la valeur de vérité de la proposition une probabilité :

$$p(k) = p(E_k)$$

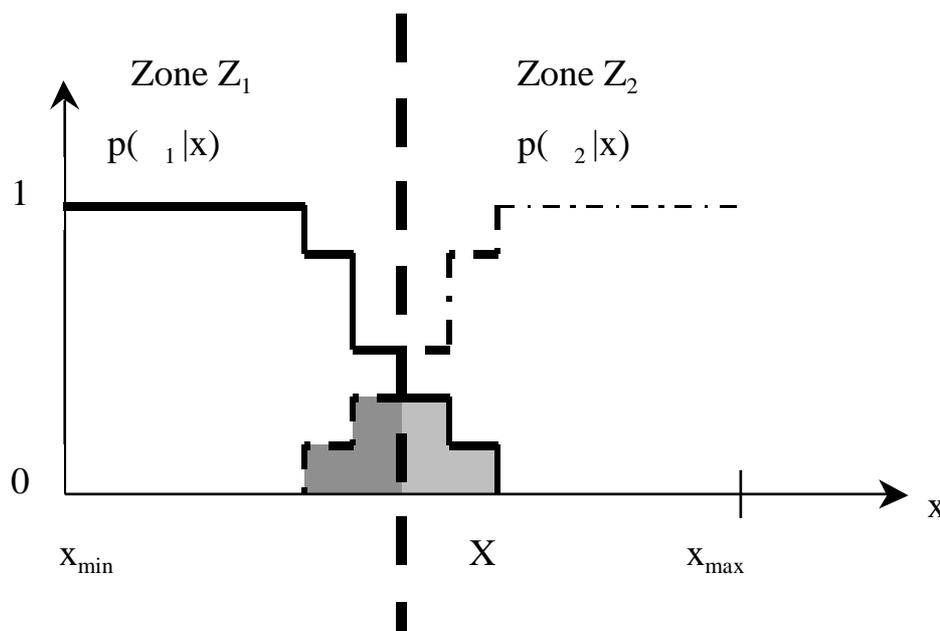
Le critère de décision est de minimiser le nombres d'erreurs. Dans un système probabiliste, ca revient de minimiser la probabilité d'erreur. Ceci est équivalent à choisir la classe le plus probable.

$$\hat{k} = \text{Decider}(E_k) = \arg\text{-max}_k \{p(k | X)\}$$

Considère le cas  $D = 1$  et  $K = 2$ . Dans ce cas, le domaine d' $X$  est le domaine.

La classification est équivalente à une decoupage du domaine d' $X$  en deux zones :  $Z_1$  et  $Z_2$ .

$$\hat{1} \text{ si } X \in Z_1 \text{ et } \hat{2} \text{ si } X \in Z_2$$



La probabilité d'erreur est la somme des probabilités de  $p(x_2)$  en  $Z_1$  et la somme de probabilité de  $p(x_1|x)$  en zone 2.

$$p(\text{erreur}) = \sum_{z_1} p(z_2 | X) + \sum_{z_2} p(z_1 | X)$$

La minimum est atteint quand :

$$\text{Donc } d(g_k(X)) = \arg\text{-max}_k \{p(z_k | X)\}$$

Dans ce cas, nous avons utilisé  $\arg\text{-max}_k \{p(z_k | X)\}$  en tant que fonction de décision

et  $g_k(X) = p(z_k | X)$  comme la fonction de discrimination

## Les Variations Aléatoires des Caractéristiques

Les vecteurs de caractéristique porte une composante "aléatoire" avec (au moins) deux origines :

- 1) Les variations des individus.
- 2) Le bruit des capteurs

### 1) Les variations entre individus d'une classe

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

$$X = x + B_i$$

Les plupart de techniques probabilistes de reconnaissance supposent un bruit additif. En notation vectorielle :

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

Attention. Le bruit peut être non-additif, et ceci devient une approximation.

### 2) Le bruit des capteurs

(def.) 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 :  $\{ Z_1, Z_2 \dots Z_D \}$ .

$$Z = \begin{pmatrix} Z_1 \\ Z_2 \\ \dots \\ Z_n \end{pmatrix}$$

Les observations sont corrompues par un bruit,  $B_0$ .

$$Z = X + B_0 = x + B_i + B_0$$

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  $k$ , la probabilité d'observé  $Z$  est fournie par la règle de Bayes.

$$p(k | Z) = \frac{p(Z | k) p(k)}{p(Z)}$$

Si  $Z = F(X + B_0)$  est issu de la classe  $k$  ayant caractéristique  $X = x + B_i$

Une Exemple - Le spectre observées par un satellite.

Une image satellite est composée de pixels  $s(i,j)$ . Chaque pixel compte le nombre de photons issus d'une surface carré de la terre (ex.  $10 \text{ m}^2$ ).

Les photons sont captés au travers des filtres spectraux. Ceci donne un vecteur de caractéristiques pour chaque pixel. Soit une région de végétation {blé, maize, etc}

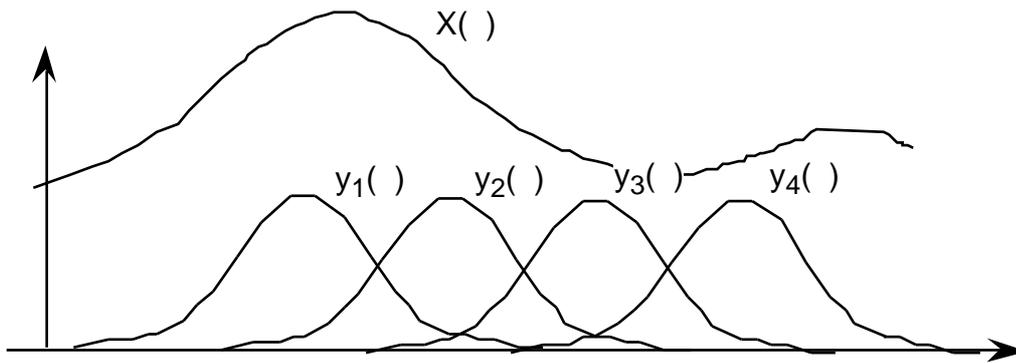
$x$  : Le spectre des pigments des feuilles pour une espèce.

$B_i$  : Les variations du spectre intrinsèque dû aux variations d'âge ou d'humidité.

$B_i$  est spécifique à un individu. Il ne change pas entre les observations.

$X = x + B_i$  : Le spectre des pigments des feuilles pour un individu

$B_0$  : Les variations d'observations dues à l'angle du soleil et les effets de filtrage de la lumière par l'atmosphère (humidité, pollution etc). Ils sont présents dans tous les pixels.



Une image est une table de pixels, avec les positions sur les lignes ( $j$ ) et colonne ( $i$ ). Un pixel,  $Z(i,j)$ , est un vecteur d'entiers,  $Z_1, \dots, Z_D$ . Chaque composant est l'intégral sur  $i, j$  et  $\lambda$  (longueurs d'onde) d'une produit d'un filtre spectral avec le spectre reçu sur la région  $c, r$ .

$$Z_d(i,j) = \text{Quant} \left\{ \int_i^{i+d_i} \int_j^{j+d_j} X(i,j, \lambda) \cdot f_d(\lambda) \, di \, dj \, d\lambda \right\}$$

C'est opération est réalisé par l'optique de la caméra. L'opération "Quant{ }" numérise les valeurs  $y_d$  avec un pas "q" sur une plage entre 0 et  $v_{\max}$ .

Une classification est une estimation de la classe de végétations dominant,  $k$ , pour la région observée par le pixel à partir de le vecteur d'observation  $\hat{Z}$

$$\hat{k} = \text{Class} \{ \hat{Z} \}$$

## La Loi Normale

La fonction paramétrique la plus utilisée est 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.

La loi Normale décrit une population d'exemples  $\{X_m\}$ .

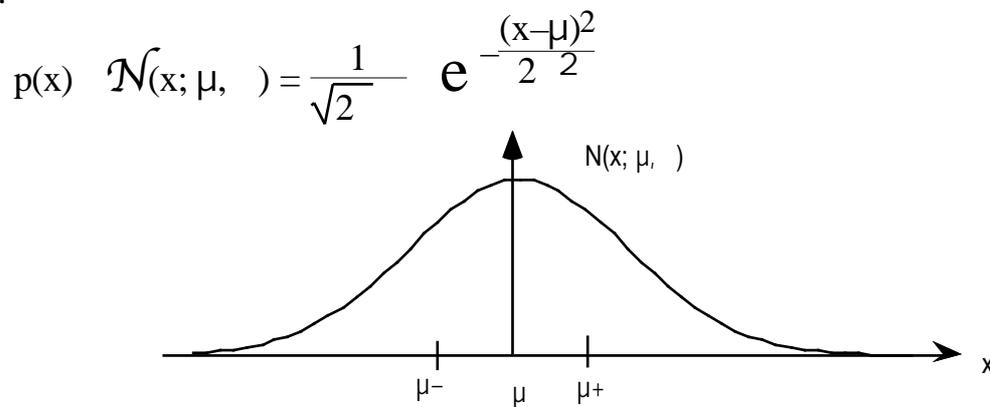
Les paramètres de  $\mathcal{N}(\mu, \sigma^2)$  sont les premiers et deuxième moments de la population.

On peut estimer les moments pour n'importe quel nombre d'exemples ( $M > 0$ )

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  $p(x)$ .



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.

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

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

La différence  $(x - \mu)^2$  est la "distance" entre une caractéristique et la moyenne de l'ensemble  $\{X_m\}$ . 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

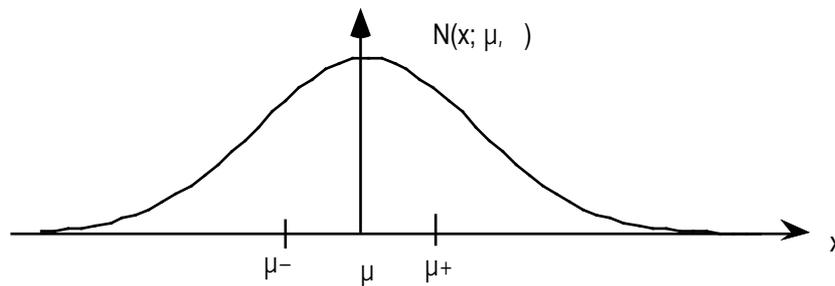
**La Loi Normale pour D = 1**

La cas le plus simple concerne une seule caractéristique.

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

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

$\mathcal{N}(x; \mu, \sigma^2)$  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  $D > 1$** 

Soit les événements  $E$  décrit par un vecteur de  $D$  caractéristiques  $X$

Soit un ensemble de  $M$  événements,  $\{E_m\}$  avec leurs caractéristiques.  $\{X_m\}$

Cet ensemble est dit l'ensemble d'entraînement (training set)

$$\mu_d = E\{x_d\} = \frac{1}{M} \sum_{m=1}^M X_{dm}$$

Pour le vecteur de  $D$  caractéristiques :

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

Pour  $M$  observations  $\{X_m\}$ , la covariance entre les variables  $x_i$  et  $x_j$  est

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

Ces coefficients composent une matrice de covariance.  $C_x$

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

En matrice on écrit :

$$\text{Soit } V = X - E\{X\} = X - \mu$$

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

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

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  $\{\mathbf{X}_m\}$

$$\hat{\sigma}_{ij}^2 = E\{(X_i - E\{X_i\})(X_j - E\{X_j\})\} = \frac{1}{M} \sum_{m=1}^M (X_{im} - \mu_i)(X_{jm} - \mu_j)$$

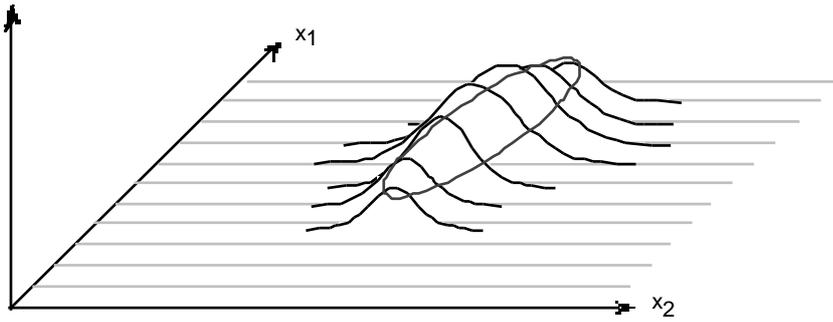
Ces coefficients composent une matrice de covariance.  $C$

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

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

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

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



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

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

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

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  $\mathbf{X}$  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>.)

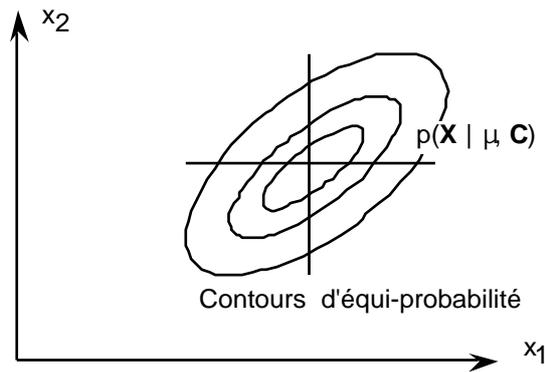
Cette valeur est connue comme la "distance de Mahalanobis".

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

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

Cette distance est bien définie, même si les composants de  $\mathbf{X}$  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  $d^2(X)$

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

si  $x_i$  et  $x_j$  sont statistiquement indépendants,  $\sigma_{ij}^2 = 0$ .

Soit les événements E décrit par une vecteur de caractéristiques X : (E,X).  
 Soit une ensemble aléatoire de M événements avec leurs caractéristiques.  
 Cet ensemble est dit l'ensemble d'entrainement (training set)  $\{X_m\}$

Pour un vecteur de D caractéristiques :

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

Pour X entier, tel que pour chaque d  $[1, D]$ ,  $X_d \in [x_{dmin}, x_{dmax}]$  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,  $\mu_d = E\{x_d\} = \int \dots p(x) \cdot x_d dX$

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  $\{\mathbf{X}_m\}$

$$\text{Soit } \mathbf{V} = \mathbf{X} - E\{\mathbf{X}\} = \mathbf{X} - \boldsymbol{\mu}$$

$$\mathbf{C}_x = E\{\mathbf{V} \mathbf{V}^T\} = E\{(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T\}$$

Ces coefficients composent une matrice de covariance.  $\mathbf{C}_x$

$$\mathbf{C}_x = \begin{matrix} & \begin{matrix} 11^2 & 12^2 & \dots & 1D^2 \\ 21^2 & 22^2 & \dots & 2D^2 \\ \dots & \dots & \dots & \dots \\ D1^2 & D2^2 & \dots & DD^2 \end{matrix} \end{matrix}$$

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

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

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

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

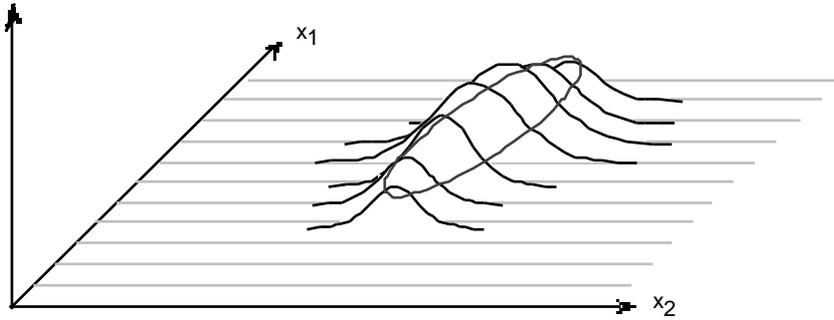
Ces coefficients composent une matrice de covariance.  $\mathbf{C}$

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

$$\mathbf{C}_x = \begin{matrix} & \begin{matrix} 11^2 & 12^2 & \dots & 1D^2 \\ 21^2 & 22^2 & \dots & 2D^2 \\ \dots & \dots & \dots & \dots \\ D1^2 & D2^2 & \dots & DD^2 \end{matrix} \end{matrix}$$

Dans le cas d'un vecteur de propriétés,  $\mathbf{X}$ , la loi normale prend la forme :

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



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



## Transformations Linéaire

La transformation linéaire d'une loi normale et une loi normale. Les moments d'une transformation linéaire sont les transformations linéaires des moments.

$$\text{Soit un vecteur unitaire } R = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} = \begin{pmatrix} \cos(\theta_1) \\ \cos(\theta_2) \\ \dots \\ \cos(\theta_D) \end{pmatrix} \quad \text{tel que } \|R\| = 1.$$

La projection (transformation linéaire) de  $X$  sur  $y$  est

$$y = R^T X.$$

Pour la covariance :

$$\begin{aligned} \sigma_y^2 &= E\{(R^T V)(R^T V)^T\} \\ &= E\{(R^T V)(V^T R)\} \quad \text{car } (R^T V)^T = (V^T R) \\ &= E\{R^T (V V^T) R\} \\ &= R^T E\{V V^T\} R = R^T C_X R \end{aligned}$$

La projection de la covariance est la covariance de la projection.

La projection de la moyenne et la covariance sur un axe,  $R$  donne une moyenne  $\mu_y$  et variance,  $\sigma_y^2$  dans la direction  $R$ .

$$\mu_y = R^T \mu_x, \quad \sigma_y^2 = R^T C_X R$$

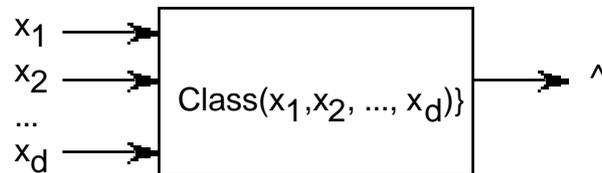
$$p(y) = \mathcal{N}(y; R^T \mu_x, R^T C_X R) = \mathcal{N}(y; \mu_y, \sigma_y^2)$$

Les moments d'une projection sont les projections des moments.

$$\mu_y = E\{p(y)\} = R^T \mu_x \quad \sigma_y^2 = E\{(p(y) - \mu_y)(p(y) - \mu_y)^T\} = R^T C_X R$$

## Fonctions de Discrimination

La classification est un processus d'estimation de l'appartenance d'un événement à une des classes  $A_k$  fondée sur les caractéristiques de l'événement,  $X$ .



$$\hat{k} = \text{Classer}(E) = \text{Decider}(E \quad k)$$

$\hat{k}$  est la proposition que  $(E \quad k)$ .

La fonction de classification est composée de deux parties  $d()$  et  $g_k()$ :

$$\hat{k} = d(g(X)).$$

$g(X)$  : Une fonction de discrimination :  $\mathbb{R}^D \rightarrow \mathbb{R}^K$

$d()$  : Une fonction de décision :  $\mathbb{R}^K \rightarrow \{K\}$

### Discrimination

$g(X)$  : Une fonction de discrimination est une fonction  $\mathbb{R}^D \rightarrow \mathbb{R}^K$

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

Etant donnée  $X$ , pour chaque  $k$  il existe une valeur de probabilité  $p(k | X)$

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

Dans le cas général la nombre minimum d'erreur est fait si  $k$  est choisi tel que

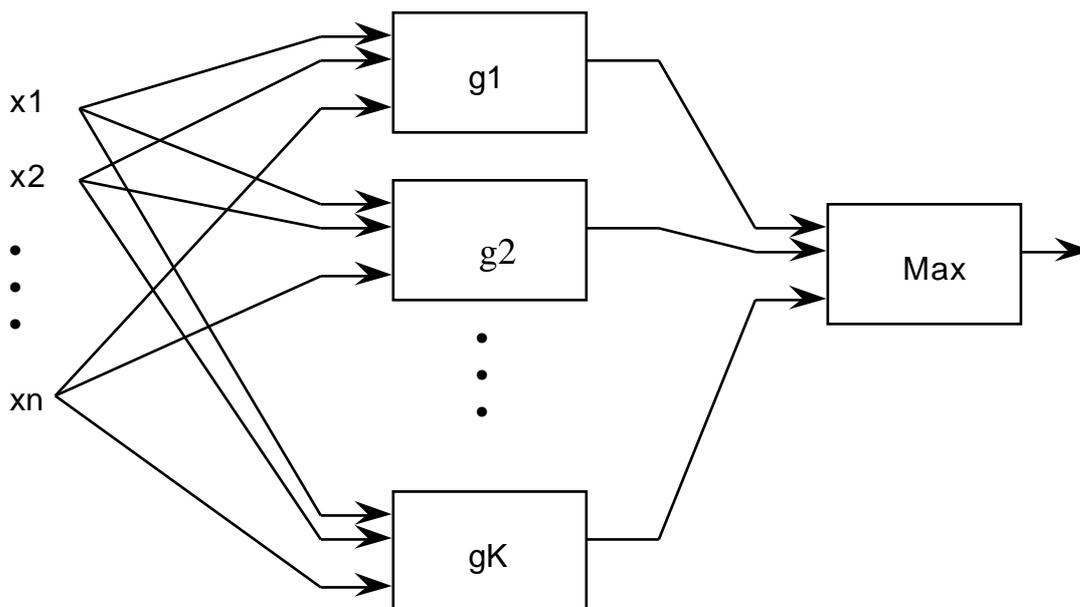
$$k = \arg\text{-max}_k \{g_k(X)\} = \arg\text{-max}_k \{p(k | X)\} = \arg\text{-max}_k \left\{ \frac{P(X|k)}{P(X)} p(k) \right\}$$

mais, comme  $P(X)$  est constant pour tous  $k$ ,

$$k = \arg\text{-max}_k \{ P(X|k) p(k) \}$$

Il suffit de l'évaluer  $P(X|k)$ , pour  $X=x$

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



Fonctions classiques :

$$P(X|k) = \mathcal{N}(X; \mu_k, C_k)$$

ou encore

$$P(X|k) = \prod_{n=1}^N \mathcal{N}(X_n; \mu_{kn}, C_{kn})$$

**Simplification de la fonction de Discrimination**

Soit  $D=1$ , avec

$$p(X=x | k) = \mathcal{N}(x; \mu_k, \sigma_k^2) = \frac{1}{\sqrt{2\pi} \sigma_k} e^{-\frac{(x-\mu_k)^2}{2\sigma_k^2}}$$

Donc notre fonction de discrimination devient :

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

On peut noter que  $k = \arg\text{-max}_k \{g_k(X)\} = \arg\text{-max}_k \{\text{Log}\{g_k(X)\}\}$

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

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

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

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

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

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

Dans le cas général, il y a  $D$  caractéristique.

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

Et le règle de décision est :

$$\hat{\omega}_i : \text{si } \omega_j \text{ si } g_i(\mathbf{X}) > g_j(\mathbf{X})$$

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

On sélection la classe  $\omega_k$  pour laquelle  $\arg\text{-max}_k \{g_k(\mathbf{X})\}$

par règle de Bayes :

$$\arg\text{-max}_k \{p(\omega_k | \mathbf{X})\} = k = \arg\text{-max}_k \{p(\mathbf{X} | \omega_k) p(\omega_k)\}$$

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

Si les caractéristiques suivent une densité Normale :

$$p(\mathbf{X} | \omega_k) = \mathcal{N}(\mathbf{X}, \boldsymbol{\mu}_k, \mathbf{C}_k)$$

$$\text{Log}\{p(\mathbf{X} | \omega_k)\} = \text{Log}\left\{ \frac{1}{(2\pi)^{\frac{D}{2}} \det(\mathbf{C}_k)^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{X} - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1} (\mathbf{X} - \boldsymbol{\mu}_k)} \right\}$$

$$\text{Log}\{p(\mathbf{X} | \omega_k)\} = -\frac{D}{2} \text{Log}\{2\pi\} - \frac{1}{2} \text{Log}\{\det(\mathbf{C}_k)\} - \frac{1}{2}(\mathbf{X} - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1} (\mathbf{X} - \boldsymbol{\mu}_k)$$

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

La fonction de discrimination devient :

$$g_k(x) = -\frac{1}{2} \log\{\det(C_k)\} - \frac{1}{2}(x - \mu_k)^T C_k^{-1} (x - \mu_k) + \log\{p(x|k)\}$$

Les classifieurs Bayésiens sont définies par les variations de cette formule.

## Forme Canonique de la fonction de discrimination

La décision  $w_k$  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

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

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

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

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

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

$$d_k = C_k^{-1} \mu_k$$

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

Cette fonction est composée de trois termes :

une terme quadratique  $X^T (D_k) X,$

une terme linéaire :  $d_k^T X$

et une terme constant :  $d_{k0}$