

Systemes Intelligents : Raisonnement et Reconnaissance

James L. Crowley

Deuxième Année ENSIMAG

Deuxième Semestre 2005/2006

Séance 9

5 avril 2006

Reconnaissance Bayesienne

Notations	2
La Loi Normale.....	3
Estimations des moments d'une densité.....	4
Le premier moment : La Moyenne.....	4
Le deuxième moment (La variance).....	5
La Loi Normale pour $D = 1$	6
La Loi Normale pour $D > 1$	7
Forme en Algebre Linéaire.....	12
Transformations Linéaire.....	13
Fonctions de Discrimination	14
Discrimination.....	14
Bruit d'observation.....	17
Classification pour $K > 2$ et $D > 1$	18
Forme Canonique de la fonction de discrimination.....	20
Bruit et Choix de la Fonction de Discrimination.....	21
Formes Quadratique.....	22
Cas des classes avec moyennes égales.....	24

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.
A, B	des classes d'événements.
T_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
	$M = \sum_{k=1}^K M_k$
k	L'affirmation que l'événement E est dans la classe T_k
$h(x)$	Histogrammes des valeurs (x est entières avec range limité)
$h_k(x)$	Histogramme des valeurs pour la class k .
	$h(x) = \sum_{k=1}^K h_k(x)$
k	Proposition que l'événement E est dans la classe k
$p(k) = p(E = T_k)$	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
$P(X k)$	Densité de Probabilité pour X étant donné que k
	$P(X) = \sum_{k=1}^K p(X k) p(k)$

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.

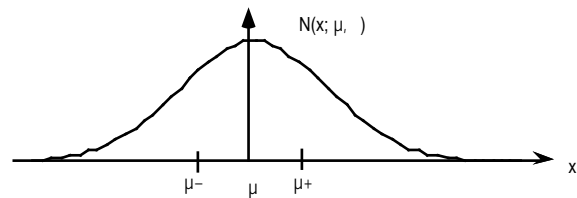
Soit M exemple d'observation d'un événement $E_m : X_m$

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 ($M > 0$)

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



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

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

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

Le terme $d^2(x) = \frac{(x-\mu)^2}{\sigma^2}$ est la différence entre x et μ 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, σ^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, $\{X_m\} : \{X_1, X_2, \dots, X_M\}$

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).

On note que dans le cas où il existe un histogramme pour X , 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

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

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

Pour X continue la moyenne peut être calculé par le 1^e moment du pdf.

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

Le deuxième moment (La variance)

La variance σ^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 μ estimé avec le même ensemble, introduit un biais dans σ^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 !) Ils ont les usages différents.

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, la variance est la deuxième moment de la pdf.

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

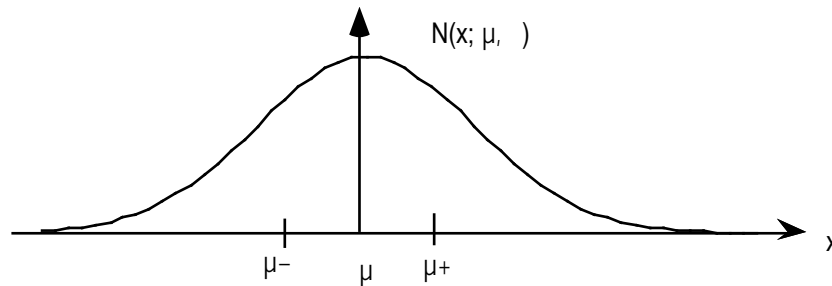
La Loi Normale pour D = 1

La cas le plus simple concerne une seule caractéristique.

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

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

$\mathcal{N}(x; \mu, \sigma^2)$ 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 σ^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

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

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

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

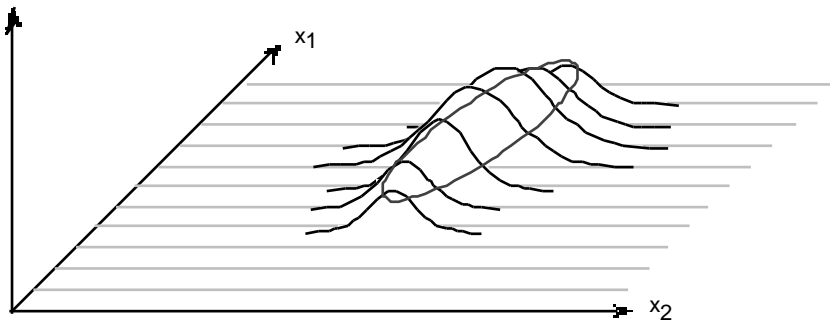
Ces coefficients composent une matrice de covariance. C

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

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

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

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



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 quadratique.

(si \mathbf{X} est en mètre, $\frac{1}{2}(\mathbf{X} - \boldsymbol{\mu})^T \mathbf{C}_X^{-1} (\mathbf{X} - \boldsymbol{\mu})$ est en mètre².)

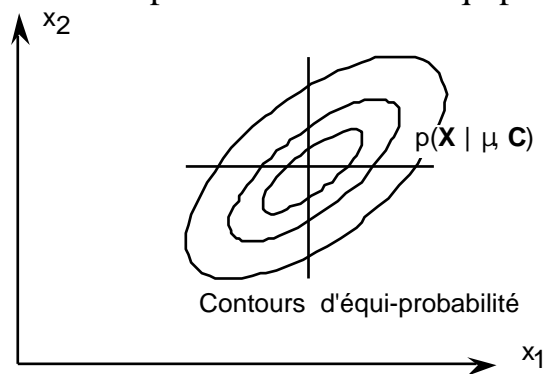
Cette valeur est connue comme la "distance de Mahalanobis".

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

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 $d^2(\mathbf{X})$

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

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

Soit les événements E décrit par une vecteur de caractéristiques $\mathbf{X} : (E, \mathbf{X})$.

Soit une ensemble aléatoire de M événements avec leurs caractéristiques.

Cet ensemble est dit l'ensemble d'entrainement (training set) $\{\mathbf{X}_m\}$

Pour un vecteur de D caractéristiques :

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

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

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

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

$$\text{Dans tous les cas : } \boldsymbol{\mu} = E\{\vec{\mathbf{X}}\} = \begin{matrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_n \end{matrix} = \begin{matrix} E\{x_1\} \\ E\{x_2\} \\ \dots \\ E\{x_n\} \end{matrix}$$

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

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

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

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

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

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

Forme en Algèbre Linéaire

Une expression en algèbre de Matrice est souvent utile.

Soit une ensemble aléatoire de M événements avec leurs caractéristiques. $\{X_m\}$

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

Soit $V_m = X_m - \mu$

On peut faire une matrix V composé de M colones $\{V_m\}$

$$V = \begin{pmatrix} V_{11} & V_{12} & \dots & V_{1M} \\ V_{21} & V_{22} & \dots & V_{2M} \\ \dots & \dots & \dots & \dots \\ V_{D1} & V_{D2} & \dots & V_{DM} \end{pmatrix}$$

$$C_x = V V^T = \begin{matrix} \begin{matrix} \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \end{matrix} & \begin{matrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{matrix} \\ \\ = \begin{matrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{matrix} \end{matrix}$$

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 μ_y et variance, σ_y^2 dans la direction R .

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

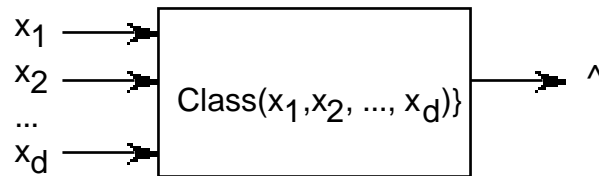
$$p(y) = \mathcal{N}(y; R^T \mu_X, R^T C_X R) = \mathcal{N}(y; \mu_y, \sigma_y^2)$$

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 $\mathbb{R}^D \rightarrow \mathbb{R}^K$

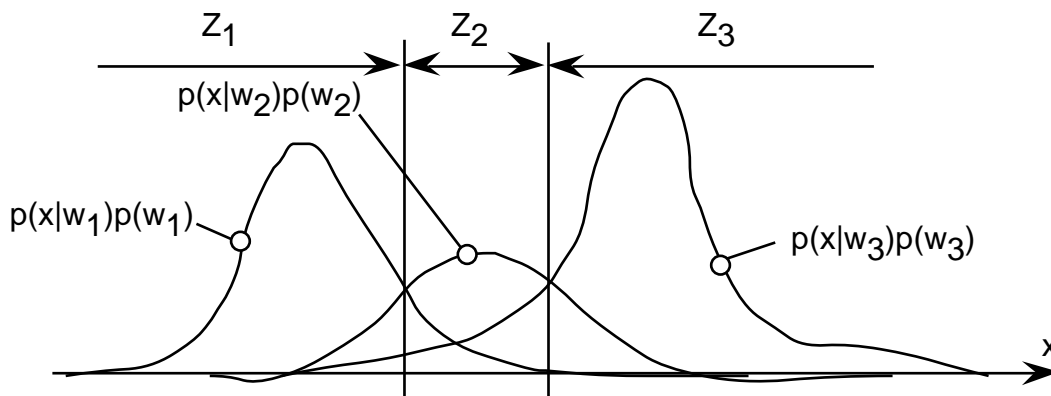
$$g(X) = \begin{pmatrix} g_1(X) \\ g_2(X) \\ \dots \\ g_K(X) \end{pmatrix}$$

Dans le cas général, $K > 2$, la nombre minimum d'erreur sont fait si k est choisi tel que :

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

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

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



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

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

Mais comment calculer $g_k(X)$?

Les caractéristiques X de l'événement E sont aléatoires avec une dispersion due aux variations naturelles de sa classe.

Ceci est modélisé par une variable aléatoire B_k autour d'une valeur "type" x_k . La valeur type est spécifique à la classe.

$$X = x_k + B_k$$

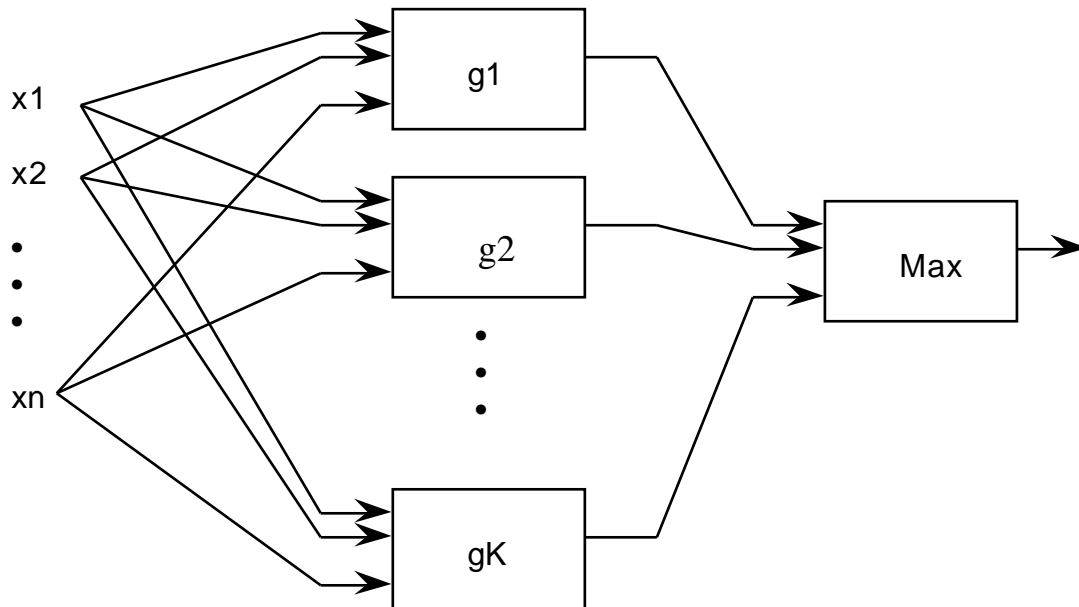
Si $D=1$, les membres de la classe k auront les caractéristiques X tel que :

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

Donc notre fonction de discrimination devient :

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

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



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

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

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

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

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

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

Bruit d'observation

Chaque observation d'un événement corrompu par un bruit d'observation.

$$Y = x_k + B_k + B_o$$

B_o est souvent Normale, avec moyenne 0 et variance σ^2 .

On dit que la variance σ^2 est la précision de la capteur.

Cette bruit ne dépend pas de la classe.

Dans ce cas, on observe Y avec

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

et

$$\begin{aligned} k &= \arg\text{-max}_k \left\{ -\text{Log} \left\{ \frac{1}{\sqrt{2\pi(\sigma_k^2 + \sigma^2)}} \right\} - \frac{(y-\mu_k)^2}{2(\sigma_k^2 + \sigma^2)} + \text{Log} \{ p(k) \} \right\} \\ &= \arg\text{-max}_k \left\{ -\text{Log} \{ \sigma_k \} - \frac{(y-\mu_k)^2}{2(\sigma_k^2 + \sigma^2)} + \text{Log} \{ p(k) \} \right\} \end{aligned}$$

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

Dans le cas général, il y a D caractéristique.

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

Et le règle de décision est :

$$\hat{k} : \text{si } j \neq i \text{ } g_i(\mathbf{X}) > g_j(\mathbf{X})$$

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

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

par règle de Bayes :

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

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

Si les caractéristiques suivent une densité Normale :

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

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

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

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

La fonction de discrimination devient :

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

Les classifieurs Bayésiennes 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}

Bruit et Choix de la Fonction de Discrimination

La conception d'un classifieur dépend de la choix de caractéristiques, x et du bruit observé sur ses caractéristiques :

Rappel qu'une observation $Y = x_k + B_k + B_o$

où

x_k Est la forme type (moyenne) de la classe w_k

B_k : Les variations aléatoires intra-classe.

Elle est spécifiques à chaque classe et chaque individus.

Elles n'est change pas entre observations.

B_o : Les variations aléatoires des observations.

Elles est indépendantes de la classe et de l'individu.

Elles changent avec les observations.

Selon la nature de B_k , B_o et de $p(x_k)$ on peut faire certaines simplifications.

par exemple :

Si $B_o \gg B_k$ on peut trouver que $w_k : C_k$ C : On a une classifieur linéaire.

Si $\sigma_{i,j}^2 = \sigma^2$ C : On a la détecteur optimale utilisé en communication hz.

Exemples de caractéristiques : x

- 1) Les échantillons d'un signal : $x(n)$ pour $n \in [1, N]$
- 2) Les caractéristiques d'un individu : [hauteur, poids, yeux, cheveux etc.]
- 3) Les caractéristiques géométriques d'un objet : Hauteur, largeur, nombre de faces, etc.

Formes Quadratique

Dans le cas le plus général, on ne fait aucune hypothèse sur B_k et B_0

Dans ce cas, C_k est arbitraire.

La surface de décision est une fonction quadrique en D dimensions.

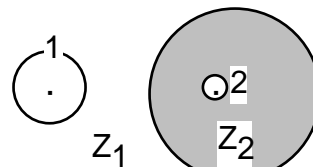
(une hyper quadriques)

Elle peut être les hyperplans, hyper-sphères, hyper-ellipsoïdes, hyper-paraboloides, ou les hyperhyperboloides.

Par exemples, en 2D (D=2) quand $K = 2$.

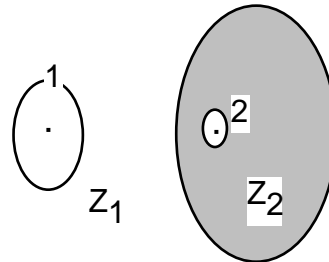
Hyper-sphère :

Pour $k = 1, 2$ $C_k = k^2 I$
et $\det\{C_1\} > \det\{C_2\}$



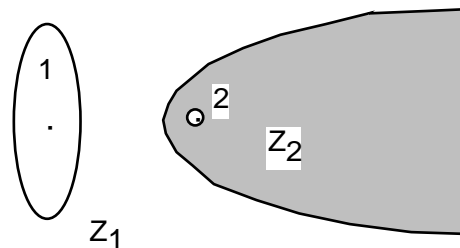
Hyper-ellipsoïde :

Pour $k = 1, 2$ $x_{1k}^2 > x_{2k}^2$
et $\det\{C_1\} > \det\{C_2\}$

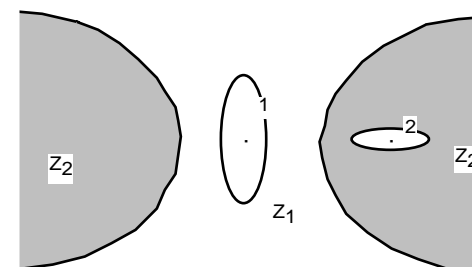


Hyper-paraboloïde :

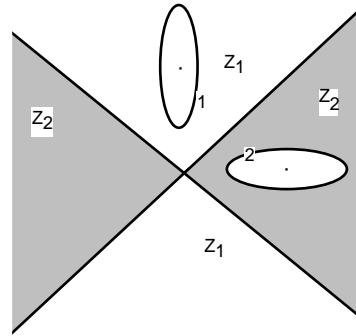
Pour $k = 1, 2$ $x_{1k=1}^2 \gg x_{1k=2}^2$
et $x_{2k=1}^2 > x_{2k=2}^2$



Hyper-hyperboloides :

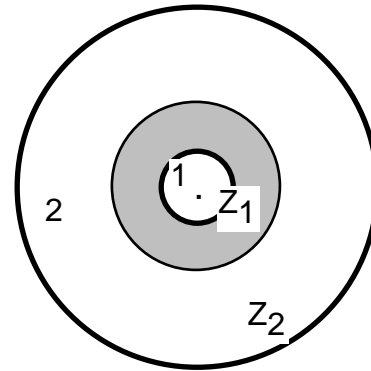


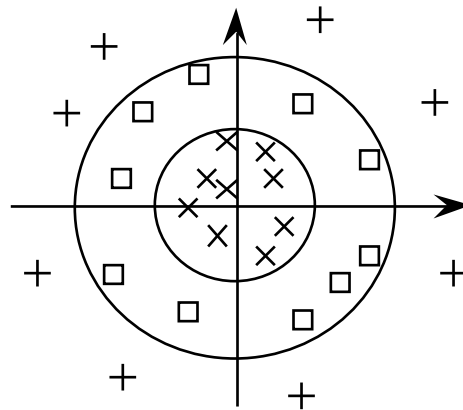
Hyperplanes.



$\mu_1 = \mu_2$ et $C_1 \ll C_2$
avec $\sigma_{11} = \sigma_{22}$ et $\sigma_{12} = \sigma_{21} = 0$.

Une hypersphere.



Cas des classes avec moyennes égales.

Supposons que nous avons K classes tel que

$$i, j: \mu_i = \mu_j \text{ et } \det(C_i) \neq \det(C_j).$$

Comment peut on décider la classe d'un événement (E, X) ?

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

- 1) $D_k = \frac{1}{2} C_k^{-1}$ est discriminant.
- 2) $d_k = C_k^{-1} \mu_k = C_k^{-1} \mu$ peut être éliminé.
- 3) $d_{k0} = -\frac{1}{2} \mu_k^T C_k^{-1} \mu_k + \text{Log}\{p(w_k)\}$ est réduit à

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

Il s'agit d'un biais pour chaque classe

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

Les surfaces de décisions entre classes i et j sont les hyper-surfaces

telles que $g_i(X) - g_j(X) = 0$ sont les hyper surfaces.