

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 (suite) | 3 |
| Forme en Algebre Linéaire..... | 4 |
| Transformations Linéaire..... | 5 |
| Fonctions de Discrimination..... | 6 |
| Discrimination..... | 6 |
| Simplification de la fonction de Discrimination..... | 8 |
| Classification pour $K > 2$ et $D > 1$ | 9 |
| Forme Canonique de la fonction de discrimination..... | 11 |
| Bruit et Choix de la Fonction de Discrimination..... | 12 |
| Formes Quadratique..... | 14 |
| Cas des classes avec moyennes égales. | 16 |

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 |
| E_k | L'affirmation que Evènement $E \in T_k$ |
| 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 entieres 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 $h(x)$. $Q = N^D$ |
| $p_k = p(E \in T_k)$ | 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 |
| $P(X k)$ | Densité de Probabilité pour X etant donné que k |
| | $P(X) = \sum_{k=1}^K p(X k) p_k$ |

La Loi Normale (suite)

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

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

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

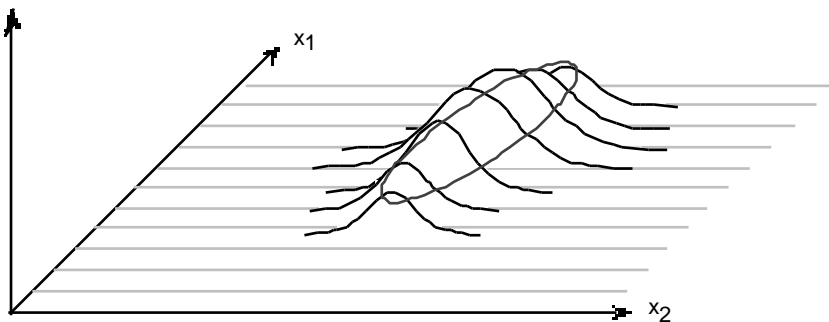
Ces coefficients composent une matrice de covariance. C

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

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

Dans le cas d'un vecteur de propriétés, 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.

Forme en Algèbre Linéaire

Soit les événements E décrit par un vecteur de caractéristiques X : (E,X).

Soit une ensemble aléatoire de M événements avec leurs caractéristiques.

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

Pour un vecteur de D caractéristiques :

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

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

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

Ceci peut être exprimé en forme de matrice.

Soit $V_m = X_m - \mu$

On peut faire une matrice V composé de M colonnes {V_m}

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

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

$C_x = V V^T$ est D x D. Note que $C_m = V^T V$ est de taille M x M.

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 :

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

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 .

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

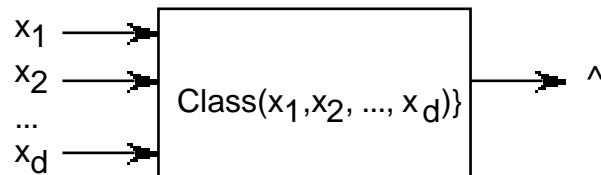
$$p(y) = \mathcal{N}(y; R^T \mu_X, R^T C_X R) = \mathcal{N}(y; \mu_y, \sigma_y^2)$$

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

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

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

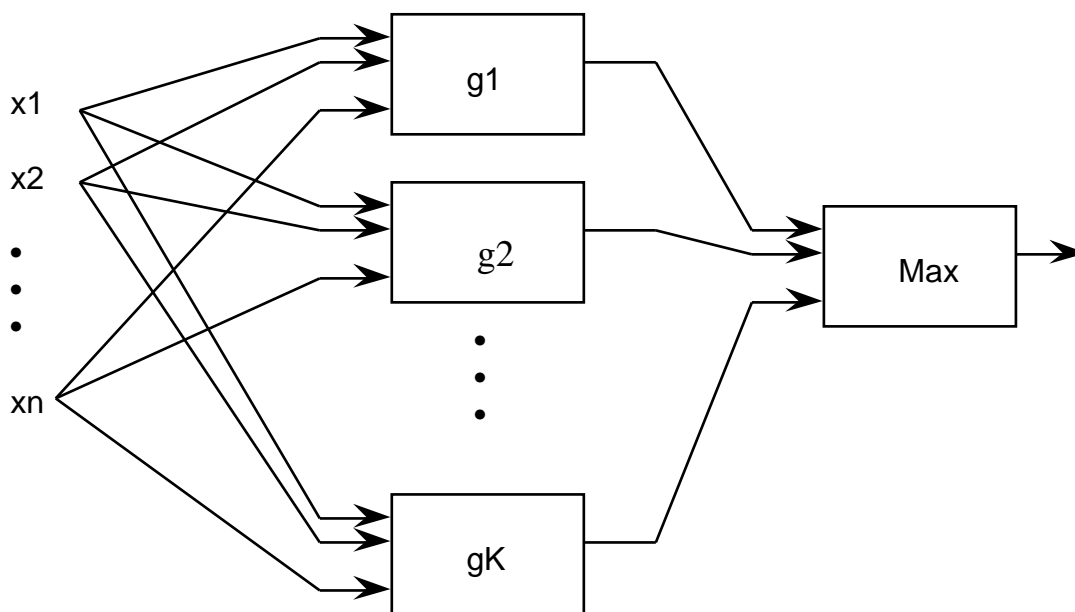
$$k = \arg\text{-max}_k \{g_k(\mathbf{X})\} = \arg\text{-max}_k \{p(k | \mathbf{X})\} = \arg\text{-max}_k \left\{ \frac{P(\mathbf{X} | k)}{P(\mathbf{X})} p(k) \right\}$$

mais, comme $P(\mathbf{X})$ est constant pour tous k ,

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

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

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



Fonctions classiques :

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

ou encore

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

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 la règle de décision est :

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

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

On sélectionne la classe ω_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 :

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

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

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

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

La fonction de discrimination devient :

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

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

Forme Canonique de la fonction de discrimination

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

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

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

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

Comment choisir la fonction pour $P(X|k)$?

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.

$$X = x_k + B_k$$

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

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

et

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

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

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

Donc $Y = x_k + B_k + B_o$

où

x_k Est la forme type (moyenne) de la classe 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 $p(k)$ on peut faire certaines simplifications.

par exemple :

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

Si $\int_{i,j} ij^2 = 2$: 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.

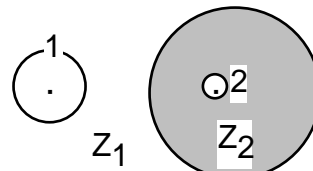
La surface de décision est une fonction quadrique en D dimensions.
 (une hyper quadriques)

Elle peut être les hyperplans, hyper-sphères, hyper-ellipsoïdes, hyper-paraboloïdes, ou les hyperhyperboloïdes.

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

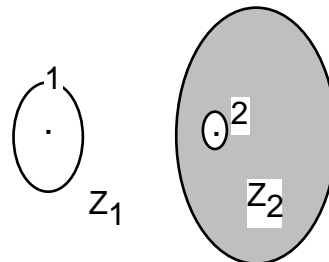
Hyper-sphère :

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



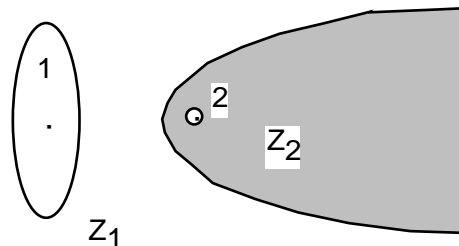
Hyper-ellipsoïde :

Pour $k = 1, 2$ $x_1^2 > x_2^2$
 et $\det\{C_1\} > \det\{C_2\}$

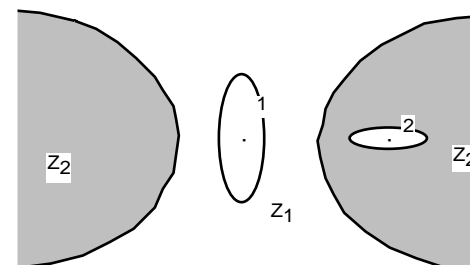


Hyper-paraboloïde :

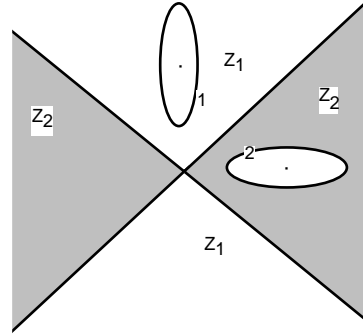
Pour $k = 1, 2$ $x_{1k=1}^2 \gg x_{1k=2}^2$
 et $x_{2k=1}^2 > x_{2k=2}^2$



Hyper-hyperboloïdes :

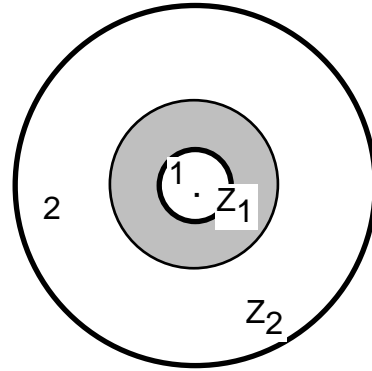


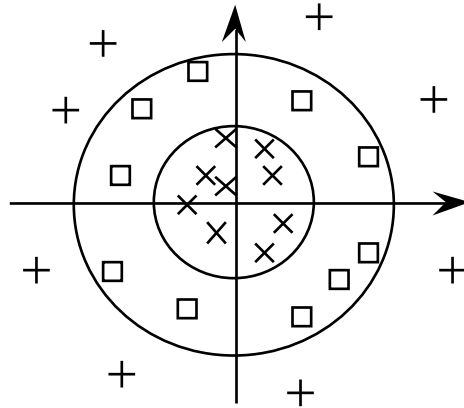
Hyperplanes.



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

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

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

Il s'agit d'un biais pour chaque classe

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

.

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.