

Vision par Ordinateur

James L. Crowley

DEA IVR

Premier Bimestre 2006/2007

Séance 4

24 octobre 2006

Reconnaissance Probabilistique

Plan de la Séance :

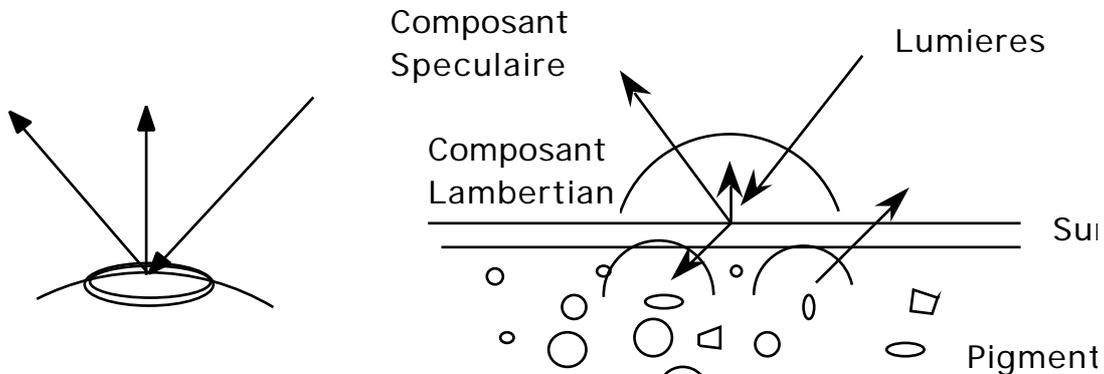
Notations.....	3
Analyse Probabiliste au Niveau Pixel	4
Invariants au niveau pixels.....	4
Histogrammes	5
Détection par ratio d'histogramme	7
Caractérisation par moments.....	8
Composantes principales.....	9
La Classification des Formes	10
La Forme	11
Les Observations.....	12
La Probabilité	14
Définition Fréquentielle.	14
Définition Axiomatique.....	14
La probabilité de la valeur d'une variable aléatoire.....	14
Fonctions de Discrimination.....	16
La Loi Normale :	19
Représentation Paramétrique de la Probabilité.....	19
La Loi Normale pour $D = 1$	20
Estimations des moments d'une densité	21
La Loi Normale pour $D > 1$	23
Forme en Algèbre Linéaire	28
Transformations Linéaire.....	29
Bruit d'observation.....	30
Classification pour $K > 2$ et $D > 1$	31
Forme Canonique de la fonction de discrimination.....	32
Bruit et Choix de la Fonction de Discrimination	33

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 $E = T_k$
$h(x)$	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 $E = T_k$ la classe k
$p(k) = p(E = T_k)$	Probabilité que E est un membre de la classe k .
Y	Une observation (un vecteur aléatoire).
$P(X)$	Densité de Probabilité pour X
$p(X = x)$	Valeur de probabilité q'un vecteur X prendre la valeur x
$P(X k)$	Densité de Probabilité pour X etant donné que k
	$P(X) = \sum_{k=1}^K p(X k) p(k)$

Analyse Probabiliste au niveau Pixel

Invariants au niveau pixels



Rappelons que l'albédo d'un objet non-métallique peut être approximé par la composition d'une réflexion "spéculaire" et d'une réflexion "lambertienne".

$$R(i, e, g, \lambda) = \alpha_s R_s(i, e, g, \lambda) + (1 - \alpha_s) R_l(i, e, g, \lambda)$$

Hors reflets spéculaire, pour une réflexion lambertienne :

$$R_l(i, e, g, \lambda) = p(\lambda) \cos(i)$$

Ceci peut être vu comme un produit de luminance terme $\cos(i)$ qui module la luminance et $p(\lambda)$ modulant la chrominance.

Le composant "chrominance" $p(\lambda)$ est déterminé par la composition du spectre de la source et le spectre d'absorption des pigments de la surface. Si le spectre de la source est constant, la chrominance indique l'identité de l'objet.

Ceci donne un indice pour la détection d'objet au niveau pixel.

Description de la chrominance

L'axe luminance, L , peut être défini par

$$L = R + V + B$$

La chrominance $C = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix}$ est une signature pour l'objet.

La chrominance peut être définie par plusieurs codages.

Par exemple, pour la détection de la peau, il est fréquent de voir la normalisation par la luminance laisser deux axes chromatiques : r, v

$$c_1 = r = \frac{R}{R+V+B} \quad c_2 = v = \frac{V}{R+V+B}$$

Une autre codage fréquente est un codage en "couleur opposée", par exemple :

$$L = \frac{R+G+B}{3}, \quad c_1 = \frac{R-G}{2} \quad \text{et} \quad c_2 = \frac{R+G}{2} - B.$$

$$\begin{array}{r} L \\ C_1 \\ C_2 \end{array} = \begin{array}{ccc} 0.33 & 0.33 & 0.33 \\ 0.5 & -0.5 & 0 \\ 0.5 & 0.5 & -1 \end{array} \begin{array}{l} R \\ G \\ B \end{array}$$

En tout cas, on peut apprendre et détecter la signature de la chrominance avec d'une manière statistique.

Histogrammes

Un histogramme est une table de fréquence. Il peut fournir une estimation d'une densité de probabilités.

On peut utiliser les histogrammes pour calculer la probabilité qu'un pixel est une projection d'un objet.

Par exemple, construisons un histogramme pour le vecteur de chrominance (r, v) .

Chaque pixel est un vecteur (r, v) : $C(i, j) = \begin{pmatrix} r \\ v \end{pmatrix} (i, j)$

On alloue un tableau 2D de taille N_h (exemple $32 \times 32 = 1024$ cellules) : $h(r, v)$.

Pour chaque pixel $C = C(i, j)$ dans l'image, on incrémente la cellule de l'histogramme qui correspond à $C = (r, v)$

$$h(r, v) := h(r, v) + 1 \quad \text{c'a-dire} \quad h(C) := h(C) + 1$$

Soit M Pixels dans l'image. Un histogramme des chrominances, $h(C)$, des M pixels dans une image donne leurs fréquences d'occurrence.

$$P(C) = \frac{1}{M} h(C)$$

Pour que l'estimation soit "raisonnable", il faut assurément que $M \gg N_h$

Considère une région W de M_o pixels du même image correspondance à l'objet O .

$$(i,j) \in W : h_o(C(i,j)) := h_o(C(i,j)) + 1$$

Ensuite: pour tout pixel $C(i, j) = \frac{r}{v}(i, j) : p(C| \text{objet}) = \frac{1}{M_o} h_o(C)$

Parce que W est dans l'image, la probabilité de rencontrer un pixel de W ,

$$P(W) = \frac{M_o}{M}$$

Détection par ratio d'histogramme

L histogramme permet d'utiliser la règle de Bayes afin de calculer la probabilité qu'un pixel corresponde à un objet.

Pour chaque pixel $C(i, j)$
$$p(\text{objet} | C) = p(C | \text{objet}) \frac{p(\text{objet})}{p(C)}$$

Soit M images de $I \times J$ pixels. Ceci fait $N = I \times J \times M$ Pixels.

Soit $h(r, v)$, l'histogramme de tous les N pixels.

Soit $h_o(r, v)$, l'histogramme des N_o pixels de l'objet "o".

$$p(\text{objet}) = \frac{M_o}{M}$$

$$p(C) = \frac{1}{M} h(C)$$

$$p(C | \text{objet}) = \frac{1}{M_o} h_o(C)$$

$$\text{Donc } p(\text{objet} | C) = p(C | \text{objet}) \frac{p(\text{objet})}{p(C)} = \frac{1}{M_o} h_o(C) \frac{\frac{M_o}{M}}{\frac{1}{M} h(C)}$$

$$p(\text{objet} | C) = \frac{h_o(C)}{h(C)}$$

Par exemple, voici une image de la probabilité de peau fait par ratio d'histogramme de r, v



Caractérisation par moments

Les ensemble connexes de pixels s'appelles les "blobs".

On peut décrire une blob par une vecteur de caractéristiques "invariantes" à l'orientation grâce aux "moments"

Les moments sont invariants aux transformations affines.

Pour une febre (imagerie) $w(i, j)$ de taille $N \times M$

$$\text{Somme des Pixels :} \quad S = \sum_{i=1}^M \sum_{j=1}^N w(i, j)$$

Premiers moments :

$$\mu_i = \frac{1}{S} \sum_{i=1}^M \sum_{j=1}^N w(i, j) \cdot i \quad \mu_j = \frac{1}{S} \sum_{i=1}^M \sum_{j=1}^N w(i, j) \cdot j$$

Le premier moment est le centre de gravité de la forme :

Deuxième moment :

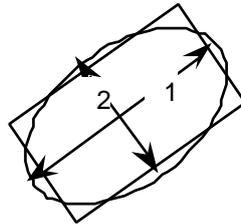
$$i^2 = \frac{1}{S} \sum_{i=1}^M \sum_{j=1}^N (w(i, j)) \cdot (i - \mu_i)^2$$

$$j^2 = \frac{1}{S} \sum_{i=1}^M \sum_{j=1}^N w(i, j) \cdot (j - \mu_j)^2$$

$$ji^2 = \frac{1}{S} \sum_{i=1}^M \sum_{j=1}^N w(i, j) \cdot (i - \mu_i)(j - \mu_j)$$

Ceci permet de définir les "axes", majeur, μ_1 et mineur, μ_2 , de la forme par analyse des composantes principales de la deuxième moment

$$C_o \cong \begin{pmatrix} i^2 & ij^2 \\ ij^2 & j^2 \end{pmatrix}$$

Composantes principales

Les deuxièmes moments sont "invariants" à l'orientation

Les axes sont calculés par une analyse en composantes principales de la matrice C . Il s'agit de trouver une rotation, Φ , dans l'espace de caractéristiques $\Phi^T C_P \Phi = \Lambda$ telles que Λ soit diagonale.

$$\Lambda = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \quad \text{tel que } 1 > 2 \quad \Phi = \begin{pmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{pmatrix}$$

$$\Phi^T C_P \Phi = \Lambda = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \quad \Phi^T \Phi = \mathbf{I} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Phi^T C_P \Phi = \Phi^T C_P \Phi = \Lambda \quad \Phi = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \quad \begin{pmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{pmatrix}$$

Les lignes du Φ sont des vecteurs propres du C .

La longueur des axes majeurs et mineur est les valeurs propres de la matrice C .

θ est l'orientation de l'axe "majeur" et $1/2$ est le rapport entre la longueur et la largeur.

$1/2$ est une caractéristique invariante de la taille et de l'orientation.

La Classification des Formes

La classification est une capacité fondamentale de l'intelligence.

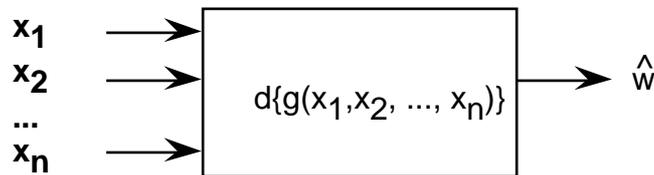
Classer : Reconnaître un membre d'une catégorie, ou d'une classe.

On peut distinguer "reconnaissance" et "identification".

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

Identifier : Reconnaître un individu

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.

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

Les Caractéristiques

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 :

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

Par exemple, $X = \begin{pmatrix} \mu_i \\ \mu_j \\ 1 \\ 2 \end{pmatrix}$ est une vecteur de caractéristique pour les "blobs".

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 resume 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 k , la probability d'observé Y est fournie par la règle de Bayes.

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

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 $M = 10 \text{ K}$ 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.

soit k , la proposition que l'observation X la classe k

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

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

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

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

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

ce qui donne :

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

La Probabilité

Définition Fréquentielle.

Une définition "Fréquentielle" de la probabilité sera suffisante pour la plupart des techniques vues dans ce cours.

Soit M observations des événement aléatoire dont M_k appartient à la classe A_k . La Probabilité d'observer un événement E de la classe A_k est

$$p(E = A_k) = \lim_M \left\{ \frac{M_k}{M} \right\}$$

Pour le cas pratique où M est fini, $p(E = A_k) \approx \frac{M_k}{M}$

La validité (ou précision) de l'approximation dépend du nombre d'échantillons M .

Définition Axiomatique.

Une définition axiomatique permet d'éviter certaine difficulté dans l'analyse de systèmes probabilistes. Trois postulats sont suffisants :

Postulat 1 : $A_k \in S : p(E = A_k) \geq 0$

Postulat 2 : $p(E \in S) = 1$

Postulat 3 :

$A_i, A_j \in S$ tel que $A_i \cap A_j = \emptyset : p(E \in A_i \cup A_j) = p(E \in A_i) + p(E \in A_j)$

La probabilité de la valeur d'une variable aléatoire

Pour x entier, tel que $x \in [x_{\min}, x_{\max}]$, on peut traiter chacun des valeurs possibles comme une classe d'événement.

Si les valeurs de x sont entières, tel que $x \in [x_{\min}, x_{\max}]$ on peut estimer la probabilité à partir de M observations de la valeur, $\{X_m\}$.

Pour estimer la probabilité d'une valeur on peut compter le nombre d'observation de chaque valeur, x , dans une table, $h(x)$.

L'existence des ordinateurs avec des centaines de megabytes rendre des tables de fréquence très pratique pour la mise en œuvre en temps réel des algorithmes de reconnaissance. Dans certains domaines, comme l'analyse d'images, par abus de langage, un tel table s'appelle une histogramme. Proprement dit, l'histogramme est une représentation graphique de $h(x)$

Ainsi la probabilité d'une valeur de $X \in [X_{\min}, X_{\max}]$ est la fréquence d'occurrence de la valeur. Avec M observations de la valeur, X , on peut faire une table, $h(x)$, de fréquence pour chacun des valeurs possibles. On observe M exemples de X , $\{X_m\}$.

Pour chaque observation on ajoute "1" à son entrée dans la table.

$$m=1, M : h(X_m) := h(X_m) + 1; M := M+1;$$

$h(x)$ est une table de fréquence pour chaque $x \in [x_{\min}, x_{\max}]$.

Ainsi, on peut définir la probabilité d'une valeur x par sa fréquence :

$$p(X_m=x) = \lim_M \left\{ \frac{1}{M} h(x) \right\}$$

Quand M est fini, on peut faire appel à l'approximation.

$$P(X=x) \approx \frac{1}{M} h(x)$$

La validité de l'approximation dépend du nombre de valeurs possible et de M . En règle générale, on dit qu'il faut 10 exemples par cellule de l'histogramme.

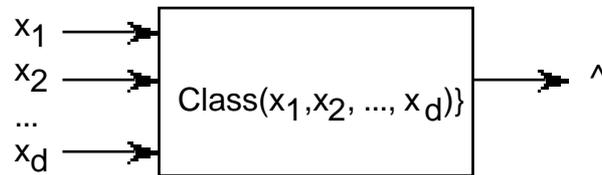
Que faire si la masse d'exemple est insuffisante : $M < 10 (X_{\max} - X_{\min})$?

Que faire si x n'est pas entier ?

Dans ces cas, on peut faire appel à une fonction paramétrique pour $p(X)$.

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 .



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

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

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

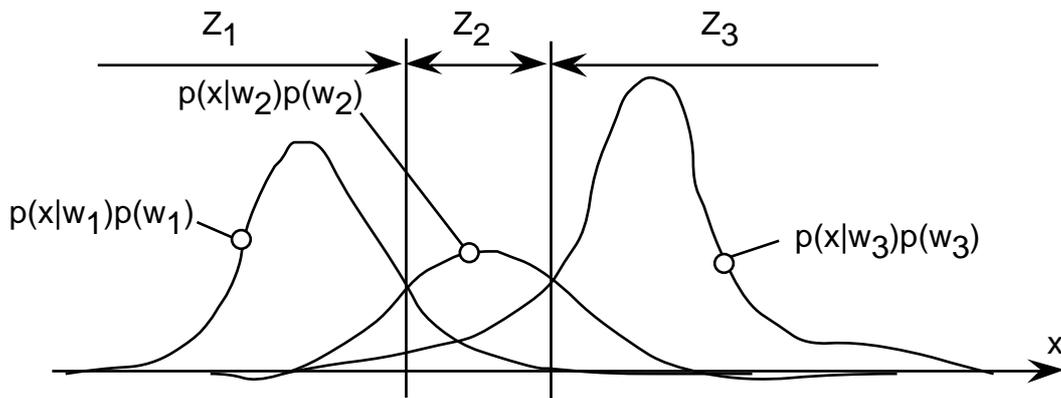
en générale $g(X) = \begin{pmatrix} g_1(X) \\ g_2(X) \\ \dots \\ g_K(X) \end{pmatrix}$ est un vecteur de K fonctions $g_k(X)$

Dans le cas général, $K > 2$, le nombre minimum d'erreurs est 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 partition l'espace de caractéristique en régions disjointes Z_1, \dots, Z_k pour chaque classe.

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

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 :

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

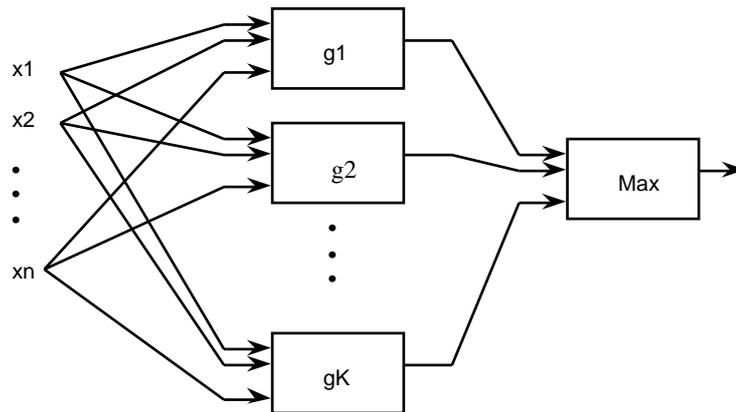
Donc notre fonction de discrimination devient :

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

pour $D > 1$ il faut la fonction normales multi-variante

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

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



Soit $D = 1$. (une seule caractéristique).

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.

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

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

Dans le cas générale $D > 1$ La fonction de discrimination devient :

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

Ceci peut être traduit dans une forme canonique :

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

avec une terme 2^{ieme} ordre : $D_k = \frac{1}{2} C_k^{-1}$

une terme 1^{iere} ordre : $d_k = C_k^{-1} \mu_k$

Constant : $d_{ko} = -\frac{1}{2}(\mu_k^T C_k^{-1} \mu_k) - \frac{1}{2} \text{Log}\{\det(C_k)\} + \text{Log}\{p(k)\}$

La Loi Normale :

Représentation Paramétrique de la Probabilité.

Une "Densité" de probabilité (écrit "pdf" pour "Probability Density Function") est une fonction, $P(X)$ représentant la probabilité pour une variable aléatoire, $X \in [-\infty, \infty]$ tel que

$$\int_{-\infty}^{\infty} P(x) dx = 1$$

Quand les variables aléatoires sont issues d'une séquence d'événements aléatoires, leur densité de probabilité prend la forme de la loi normale, $\mathcal{N}(\mu, \sigma)$. 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 $\mathcal{N}(\mu, \sigma)$ sont les premiers et deuxième moments de la population.

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

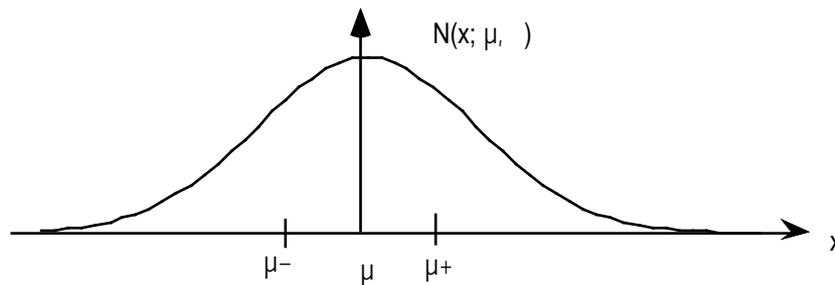
La Loi Normale pour D = 1

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 :



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.

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

Le terme $d^2(x) = \frac{(x-\mu)^2}{\sigma^2}$ est la difference 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

Les parametres du $\mathcal{N}(x; \mu, \sigma^2)$ sont le premier moment μ et le deuxieme moment

Estimations des moments d'une densitéLe 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 $\{X_m\}$.

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

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

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

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

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

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

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 $D > 1$

Soit les événements E décrit par une vecteur de D caractéristiques X

Soit une 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 en 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

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

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 $\{X_m\}$

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

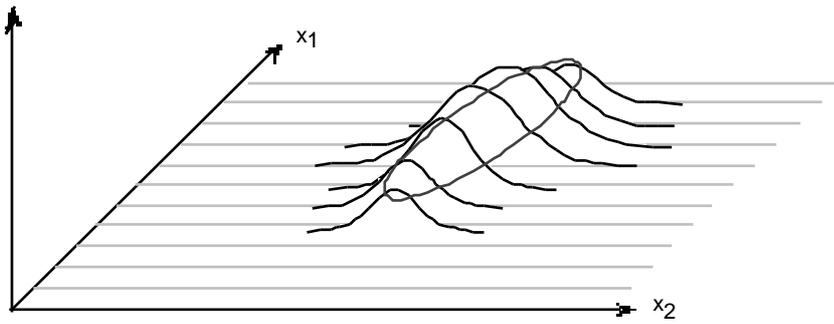
Ces coefficients composent une matrice de covariance. C

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

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

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.

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

La déterminante, $\det(\mathbf{C})$ 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 mètre².)

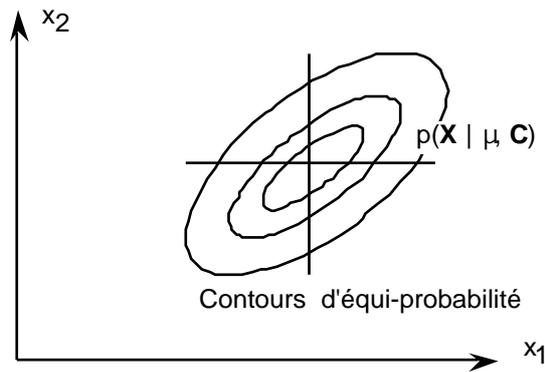
Cette valeur est connue comme la "distance de Mahalanobis".

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

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

Cette distance est bien définie, même si les composants de 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 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, $c_{ij} = 0$.

Soit les événements E décrit par un vecteur de caractéristiques X : (E, X).
Soit un ensemble aléatoire de M événements avec leurs caractéristiques.

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

Pour un vecteur de D caractéristiques :

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

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

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

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

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

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

Ces coefficients composent une matrice de covariance. C_x

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

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

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 :

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

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, y^2 dans la direction R .

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

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

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

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

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

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

et

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

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 :

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

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 :

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

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

Si les caractéristiques suivent une densité Normale :

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

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

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

One 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 :

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

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

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

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

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

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

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 $w_k : C_k \rightarrow C$: On a une classifieur linéaire.

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