

Analyse et Reconnaissance d'Images

James L. Crowley

M2R IVR

Premier Semestre 2007/2008

Séance 6

26 novembre 2007

Reconnaissance Probabiliste

Plan de la Séance :

La Reconnaissance.....	2
La Probabilité d'un Evénement.....	4
Définition Fréquentielle.....	4
Définition Axiomatique.....	4
La probabilité de la valeur d'une variable aléatoire.....	5
La Règle de Bayes.....	6
Classification des Pixels par Ratio d'Histogramme	8
Histogrammes.....	8
Exemple : Détection d'object par ratio d'histogramme de couleur.....	8
Histogrammes de Champs Réceptifs.....	10
Densités de Probabilité.....	12
La Loi Normale.....	13
La Loi Normale pour $D = 1$	14
La Loi Normale pour $D > 1$	15
Exemple : Analyse des Images Terrestre.....	19
Exemple : Caractérisation d'un région par moments.....	22
Composantes principales.....	23
Fonctions de Discrimination	24

La Reconnaissance

La reconnaissance est une capacité fondamentale de l'intelligence et même de la vie. Pour la survie, il faut savoir reconnaître les amis, les ennemies et la nourriture.

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

Identifier : Reconnaître un entité comme un individu

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

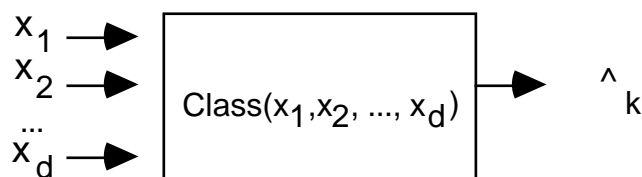
Un ensemble est défini par un test d'appartenance.

La classification est un processus d'association un entité (un événement) à une classe. L'entité est décrit par un vecteur des caractéristiques, produit par une observation. L'affectation de l'entité à une classe est fait par un test, calculer sur le vecteur de caractéristiques.

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

En notation vectorielle : $X = \begin{matrix} x_1 \\ x_2 \\ \dots \\ x_D \end{matrix}$

Pour un vecteur de caractéristique, X , un processus de classification propose une estimation de la classe, \hat{k}



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

Les classes peuvent êtres définis par
 extension : une liste complète des membres
 intention : une conjonction de prédicats.

par extension : Une comparaison d'une observation avec des membres connus de la classe (des prototypes). Ceci correspond (grosso modo) à des méthode dite "générative" de reconnaissance.

Dans ce cas, la classification peut être fait par comparaison avec les membre de la classe.

$$\hat{k} = \underset{k}{\text{arg-max}} \{ \text{Sim}(Y, X_m^k) \} \quad \text{pour tout } k, m \quad \text{ou bien}$$

$$\hat{k} = \underset{k}{\text{arg-min}} \{ \|Y, X_m^k\| \} \quad \text{pour tout } k, m$$

par intention : Conjonction de prédicats définis sur les propriétés observées. Ceci correspond (grosso modo) à des méthode dite "discriminative" de reconnaissance.

Génératives : Les techniques fondaient sur un modèle.

Discriminatives : Les techniques fondaient sur des tests quelconques.

La teste d'appartenance est une forme de partition de l'espace de caracteristiques.

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 conditionnelle d'appartenance.

$$p(k) = \Pr(E = T_k) \quad \begin{array}{l} \text{Proposition que l'évènement } E \text{ est la classe } k \\ \text{Probabilité que } E \text{ est un membre de la classe } k. \end{array}$$

Ayant une observation, X , le critère de partition est la probabilité.

$$p(k | X) = \Pr(E = T_k \text{ étant donnée l'observation de } X)$$

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

Cette probabilité est fournie par la règle de Bayes.

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

La Probabilité d'un Evénement.

La sémantique (ou "sens") de la probabilité d'un événement peut être fourni par sa fréquence d'occurrence ou par un système d'axiomes. L'approche fréquentielle a l'avantage d'être facile à comprendre. Par contre, elle peut entraîner les difficultés dans l'analyse. La définition axiomatique favorise les analyses mathématiques.

Dans le deux cas, la probabilité est une fonction numérique, $\Pr() \in [0, 1]$.
Le domaine de la fonction $\Pr()$ est les événements, E .

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 appartiennent à la classe A_k .
La Probabilité d'observer un événement E de la classe A_k est

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

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

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

Définition Axiomatique.

Une définition axiomatique permet d'appliquer certaines techniques d'analyse de systèmes probabilistes. Trois postulats sont suffisants :

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

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

La Règle de Bayes

Soit un événement "E". Soit deux tribus d'événements A et B tel que certains événements sont communs à A et à B.

E peut appartenir à A ∩ B ou à $\bar{A} \cap B$ ou à $A \cap \bar{B}$ ou à $\bar{A} \cap \bar{B}$

Soit deux propositions p et q.

donc $P(p) = \Pr\{E \in A\}$ et $P(q) = \Pr\{E \in B\}$.

Par axiome 2 de la définition des systèmes de probabilités :

$$P(q) + P(\bar{q}) = 1.$$

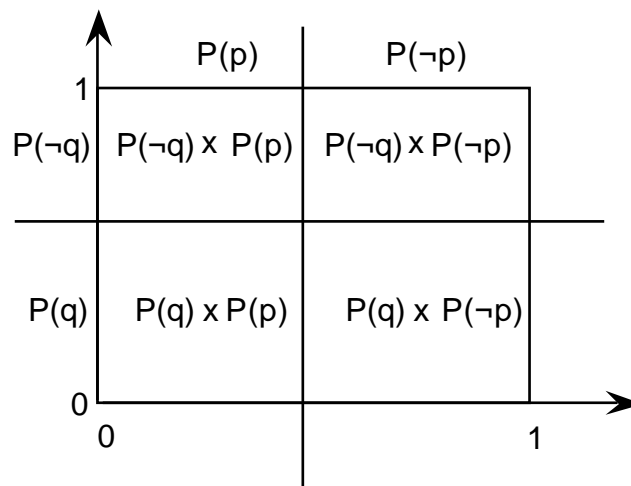
$P(p \cap q)$ est la probabilité "conjointe" de p et q.

Si p et q sont indépendantes

$$P(p \cap q) = P(p) \cdot P(q),$$

$$P(p \cup q) = P(p) + P(q).$$

On peut voir ça d'une manière graphique :



$$P(p \cap q) + P(p \cap \bar{q}) + P(\bar{p} \cap q) + P(\bar{p} \cap \bar{q}) = 1$$

Dans ce cas, les probabilités marginales sont

$$P(p) = P(p \cap q) + P(p \cap \bar{q})$$

$$P(q) = P(p \cap q) + P(\bar{p} \cap q)$$

La probabilité conditionnelle de q étant donnée p s'écrit $P(q | p)$

$$P(q | p) = \frac{P(p \cap q)}{P(p)} = \frac{P(p \cap q)}{P(p \cap q) + P(p \cap \neg q)}$$

de la même manière :

$$P(p | q) = \frac{P(p \cap q)}{P(q)} = \frac{P(p \cap q)}{P(p \cap q) + P(\neg p \cap q)}$$

Par algèbre on déduit :

$$P(q | p) P(p) = P(p \cap q) = P(p | q) P(q)$$

d'où

$$P(q | p) P(p) = P(p | q) P(q)$$

Ceci est une forme de règle de Bayes. On peut écrire :

$$P(q | p) = \frac{P(p | q) P(q)}{P(p)}$$

$P(q | p)$ est la probabilité "conditionnelle" ou "postérieur"

Classification des Pixels par Ratio d'Histogramme

Histogrammes

La probabilité d'une valeur x est sa fréquence d'occurrence.

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

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 ratio entre le nombre de Cellules, $Q = N^d$, de $h(x)$ et le nombre d'échantillons, M .

L'erreur moyenne entre $\frac{1}{M} h(C)$ et $P(C)$ est $E_{ms} \sim O\left(\frac{Q}{M}\right)$

Pour que l'estimation soit "raisonnable", il faut assuré que $M \gg Q = N^d$
En règle générale, on dit qu'il faut 10 exemples par cellule de l'histogramme.

Exemple : Détection d'objet par ratio d'histogramme de couleur

On peut utiliser les histogrammes avec la règle de Bayes pour détecter les objets.

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

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 du peau, il est fréquent de voir

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

Supposons qu'on code c_1 et c_2 avec les entiers entre 0 et $N-1$

$$c_1 = \text{Round} \left((N-1) \cdot \frac{R}{R+G+B} \right) \quad c_2 = \text{Round} \left((N-1) \cdot \frac{G}{R+G+B} \right)$$

On alloue un tableau 2D, $h(c_1, c_2)$, de taille $N \times N$ cellules.
(exemple $Q = 32 \times 32 = 1024$ cellules)

Pour chaque pixel $C = C(i, j)$ dans l'image, on incrémente la cellule de

l'histogramme qui correspond à C : $h(C) := h(C) + 1$
 c'a-dire $h(c_1, c_2) := h(c_1, c_2) + 1$

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

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

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 $C(i, j) = \begin{matrix} r \\ v \end{matrix} (i, j) : p(C | objet) = \frac{1}{M_o} h_o(C)$

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

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

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

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

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".

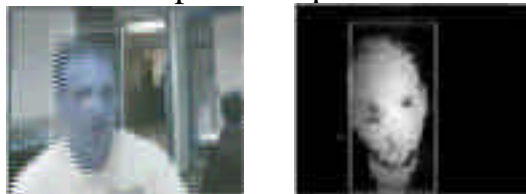
$$p(objet) = \frac{M_o}{M}$$

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

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

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

Voici une image de la probabilité de peau fait par ratio d'histogramme de r, v



Histogrammes de Champs Réceptifs.

Cette méthode peut être généralisé en remplaçant la chrominance par un vecteur de champs réceptifs. Mais il faut bien gérer la ration Q/M !

Soit une image $p(i,j)$ et un vecteur de "d" champs réceptifs G

$V(i,j) = \langle G, p(i,j) \rangle$ est un vecteur de caractéristiques de d dimensions.

$h(V)$ aura $Q = N^d$

N \ d	1	2	3	4	5	6
2	2^1	2^2	2^3	2^4	2^5	2^6
4	2^2	2^4	2^6	2^8	$2^{10} = 1 \text{ Kilo}$	$2^{12} = 2 \text{ Kil}$
8	2^3	2^6	2^9	2^{12}	2^{15}	2^{18}
16	2^4	2^8	2^{12}	2^{16}	$2^{20} = 1 \text{ Meg}$	$2^{24} = 4 \text{ M}$
32	2^5	$2^{10} = 1 \text{ Kilo}$	2^{15}	$2^{20} = 1 \text{ Meg}$	2^{25}	$2^{30} = 1 \text{ Gi}$
64	2^6	2^{12}	2^{18}	2^{24}	$2^{30} = 1 \text{ Gig}$	2^{36}
128	2^7	2^{14}	$2^{21} = 2 \text{ Meg}$	2^{28}	2^{35}	$2^{42} = 2 \text{ Ter}$
256	2^8	2^{16}	2^{24}	$2^{32} = 2 \text{ Gig}$	$2^{40} = 1 \text{ Tera}$	2^{48}

Soit les champs réceptifs achromatique

$$G = (G_x, G_y, G_{xx}, G_{xy}, G_{yy})$$

$$d = 5$$

ou chromatique avec normalisation de l'orientation et échelle :

$$G_c = (G_x^L, G_x^{C1}, G_x^{C2}, G_{xx}^L, G_{xy}^L, G_{xx}^{C1}, G_{xx}^{C2})$$

$$d = 7.$$

On peut faire

$$p(\text{objet}(i,j) | V(i,j)) = \frac{p(V(i,j) | \text{objet}(i,j)) p(\text{objet}(i,j))}{p(V(i,j))} \quad \frac{h_o(V(i,j))}{h_{\text{tot}}(V(i,j))}$$

sur condition de gérer M et Q.

Rappel :

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

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.

On peut démontrer que l'écart type moyenne de l'erreur est en proportion avec la racine de la taille de la cellule de $h(x)$, N , sur le nombre d'échantillons, M .

MSE —

Que faire si la masse d'exemple est insuffisante : $M \ll N$?

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

Densités de Probabilité.

Une fonction de densité de probabilité $P(X)$, est une fonction tel que

$P(X)$ est Real et positive pour tout X .

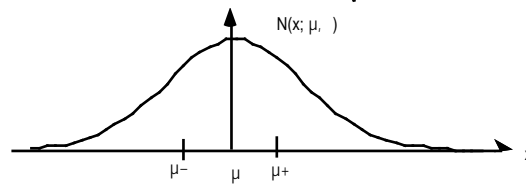
X est réel entre $[- ,]$

tel que

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

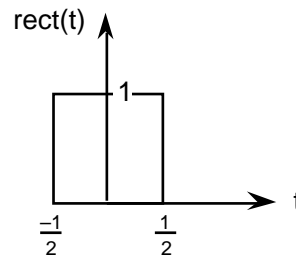
exemples :

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



Mélange de Normales $P(X) = \sum_{n=1}^N \mathcal{N}(x; \mu_n, \sigma_n)$

rect : $P(X) = \text{rect}(X)$.



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.

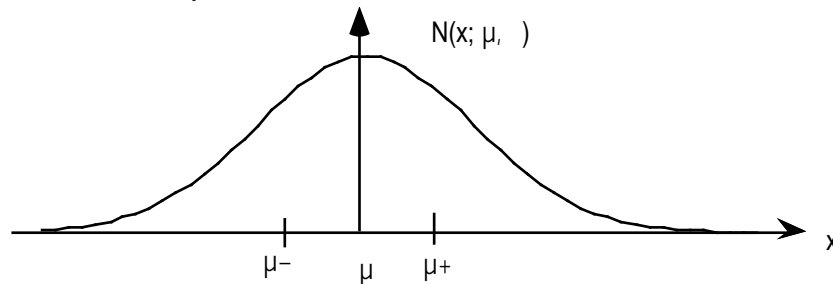
On peut estimer les moments pour n'importe quel nombre d'exemples ($M > 0$)

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 $p(x)$.

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



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

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

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

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

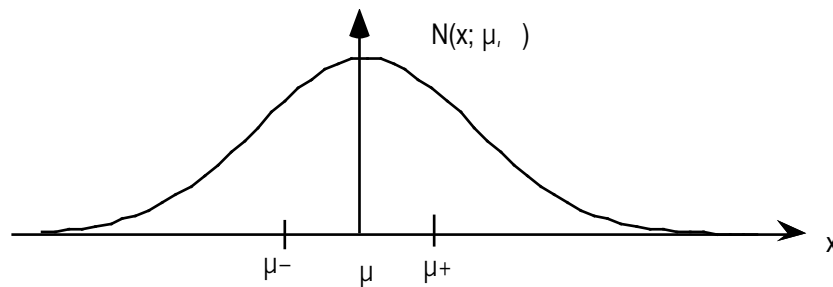
La Loi Normale pour D = 1

La cas le plus simple concerne une seule caractéristique.

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

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

$\mathcal{N}(x; \mu, \sigma^2)$ 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

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

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

Dans tous les cas :

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

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

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

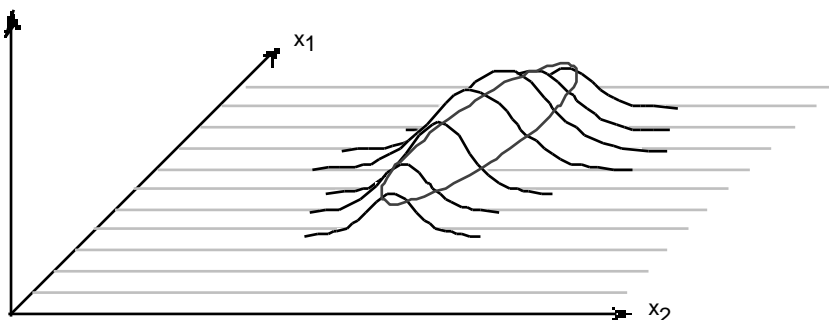
Ces coefficients composent une matrice de covariance. C

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

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

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

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



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

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

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.)

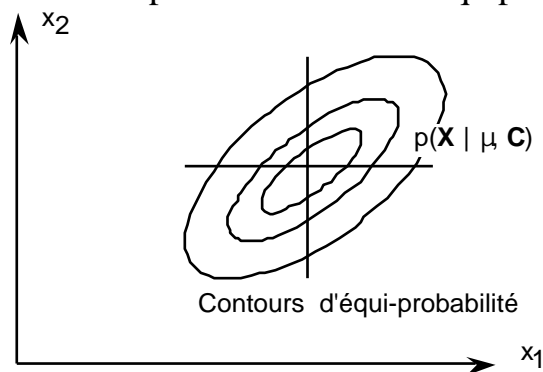
Cette valeur est connue comme la "distance de Mahalanobis".

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

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-définie. Nous allons nous limiter au cas où C est positif et définie (C.-à-d. $\det(\mathbf{C}) = |\mathbf{C}| > 0$)

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

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

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 $\{X_m\}$

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

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

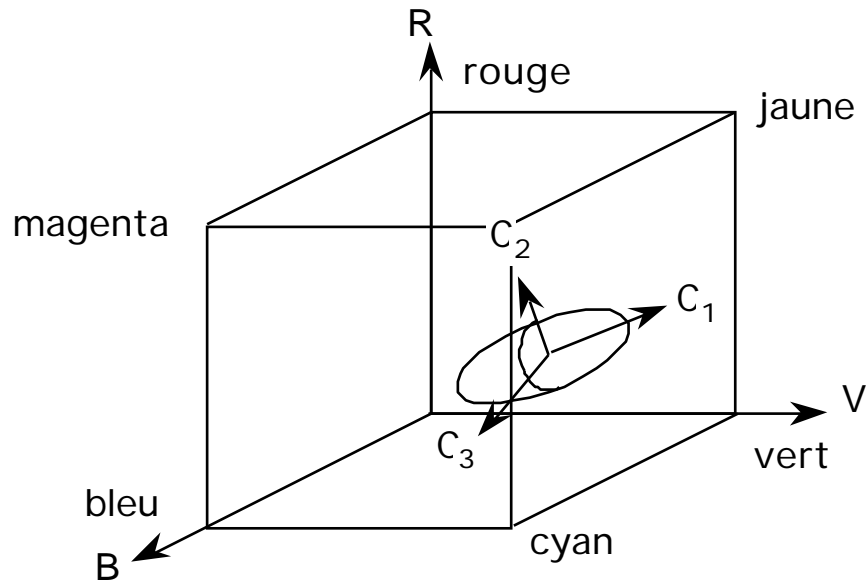
ou $\rho_{ij}^2 = \frac{1}{M} \sum_{m=1}^M (X_{im} - \mu_i)(X_{jm} - \mu_j)$

Exemple : Analyse des Images Terrestre

Les statistiques de la couleur sont utilisées couramment afin d'analyser les images pris par les satellites.

On peut utiliser la moyenne et covariance afin de caractériser la couleur d'une région. Soit un échantillon R d'une couleur dans une scène, nous pouvons calculer la distribution normale qui représente la couleur de R par la Moyenne,

$$\hat{X}_R = [r, v, b]^T \text{ et Covariance, } C_x.$$



La somme de l'histogramme :
$$S = \sum_{r,v,b=0}^{255} h(r,v,b)$$

Moyenne :
$$\hat{X} = \begin{bmatrix} \mu_r \\ \mu_v \\ \mu_b \end{bmatrix} \quad \text{où} \quad \mu_r = \frac{1}{S} \sum_{r,v,b=0}^{255} h(r,v,b)r$$

$$\mu_v = \frac{1}{S} \sum_{r,v,b=0}^{255} h(r,v,b)v \quad \mu_b = \frac{1}{S} \sum_{r,v,b=0}^{255} h(r,v,b)b$$

Covariance :
$$C_x \hat{=} \begin{bmatrix} r^2 & rv & rb \\ br & b^2 & bv \\ vr & vb & b^2 \end{bmatrix} \quad \text{où}$$

$$r^2 = \frac{1}{S} \sum_{r,v,b=0}^{255} h(r,v,b)(r - \mu_r)^2 \quad v^2 = \frac{1}{S} \sum_{r,v,b=0}^{255} h(r,v,b)(v - \mu_v)^2$$

$$b^2 = \frac{1}{S} \sum_{r,v,b=0}^{255} h(r,v,b)(b - \mu_b)^2$$

et

$$rv \quad vr \quad \frac{1}{S} \sum_{r,v,b=0}^{255} h(r,v,b)(r - \mu_r)(v - \mu_v)$$

$$rb \quad vb \quad \frac{1}{S} \sum_{r,v,b=0}^{255} h(r,v,b)(r - \mu_r)(b - \mu_b)$$

$$vb \quad vb \quad \frac{1}{S} \sum_{r,v,b=0}^{255} h(r,v,b)(v - \mu_v)(b - \mu_b)$$

Avec ces statistiques, on peut estimer la probabilité d'avoir observé une instance d'une classe, étant donné l'observation d'un pixel $Y = (r, v, b)$.

Soit K "classes", $\{w_k\}$, chacun décrit par \hat{X}_k, C_k^{-1}

Pour un observation Y

$$P(Y | \hat{X}_k) = \frac{1}{(2\pi)^{3/2} \det(C_k)^{1/2}} e^{-\frac{1}{2}(Y - \hat{X}_k)^T C_k^{-1} (Y - \hat{X}_k)}$$

La probabilité d'obtenir w_k ayant observé Y

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

La classe le plus "probable" est indépendant de $p(Y)$

$$w_k = \text{Max}_k \{ p(w_k | Y) \} = \text{Max}_k \{ p(Y | w_k) p(w_k) \}$$

$$= \text{Max}_k \{ \text{Log} \{ p(Y | w_k) p(w_k) \} \}$$

$$= \text{Max}_k \left\{ -\frac{1}{2} (Y - \mu_k)^T C_k^{-1} (Y - \mu_k) - \text{Log} \left(\frac{1}{\sqrt{2\pi}} \det(C_k)^{1/2} \right) + \text{Log} \{ p(Cw_k) \} \right\}$$

Si la bruit est indépendant de la classe, $j, k : \det(C_j) = \det(C_k)$

$$\text{Max}_k \{ p(w_k | Y) \} = \text{Max}_k \left\{ -\frac{1}{2} (\mathbf{Y} - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1} (\mathbf{Y} - \boldsymbol{\mu}_k) + \text{Log}\{p(w_k)\} \right\}$$

Si les classes sont equi-probable, $j, k : \text{Log}\{p(w_j)\} = \text{Log}\{p(w_k)\}$

$$\text{Max}_k \{ p(w_k | Y) \} = \text{Max}_k \left\{ -\frac{1}{2} (\mathbf{Y} - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1} (\mathbf{Y} - \boldsymbol{\mu}_k) \right\}$$

Pour trouver les pixels "blé" dans les images de satellite.

Déterminer la couleur "caractéristique" du blé $\hat{\mathbf{X}}_b = (R_b, V_b, B_b)$ et sa covariance. \mathbf{C}_b .

Pour un vecteur de propriétés de 3 Dimensions :

$$P(\mathbf{Y} | \hat{\mathbf{X}}_b) = \frac{1}{(2\pi)^{3/2} \det(\mathbf{C}_b)^{1/2}} e^{-\frac{1}{2} (\mathbf{Y} - \hat{\mathbf{X}}_b)^T \mathbf{C}_b^{-1} (\mathbf{Y} - \hat{\mathbf{X}}_b)}$$

Pour chaque pixel de l'image, $Y(i,j)$, on peut calculer la distance normalisée par la covariance.

$$D_b^2(i,j) = \frac{1}{2} \{ (\mathbf{Y}(i,j) - \hat{\mathbf{X}}_b)^T \mathbf{C}_b^{-1} (\mathbf{Y}(i,j) - \hat{\mathbf{X}}_b) \}$$

Une image de la "confiance" ou CF du blé :

$$p(i, j) = e^{-\frac{1}{2} (\mathbf{Y}(i,j) - \hat{\mathbf{X}}_k)^T \mathbf{C}_k^{-1} (\mathbf{Y}(i,j) - \hat{\mathbf{X}}_k)}$$

Les pixels avec une couleur proche de celle de l'échantillon apparaissent intenses.

Exemple : Caractérisation d'un région 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 fenêtre (image) $w(i, j)$ 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 :

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

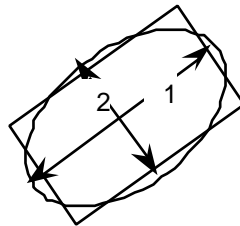
$$j^2 = \frac{1}{S} \sum_{i=1}^M \sum_{j=1}^N w(i, j) \cdot (j - \mu_j)^2$$

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

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

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

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 $\Phi C_P \Phi^T = \Lambda$ telles que Λ soit diagonale.

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

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

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

Les lignes du Φ sont des vecteurs propres du C .

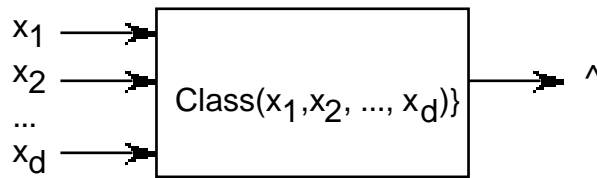
La longueur des axes majeur et mineur est les valeurs propres de la matrice C .

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.

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 .



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

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

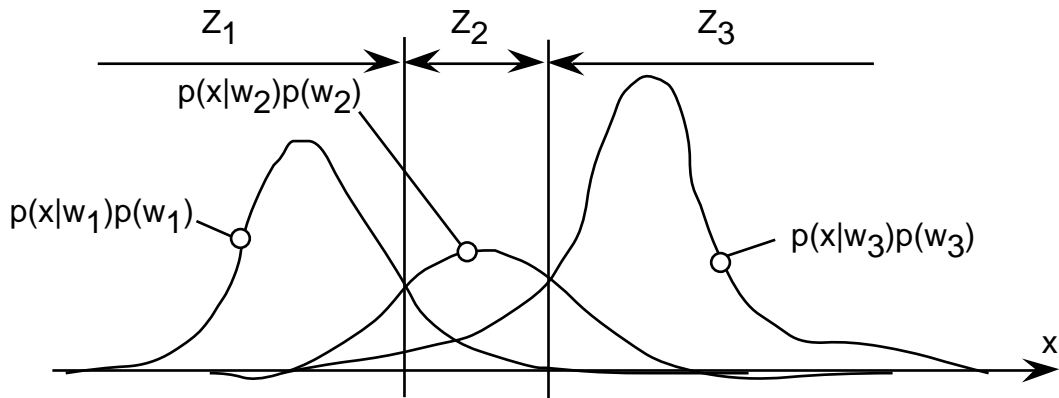
en générale $g(X) = \begin{pmatrix} g_1(X) \\ g_2(X) \\ \dots \\ g_K(X) \end{pmatrix}$ est une vecteur de K fonctions $g_k(X)$

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 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 :

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

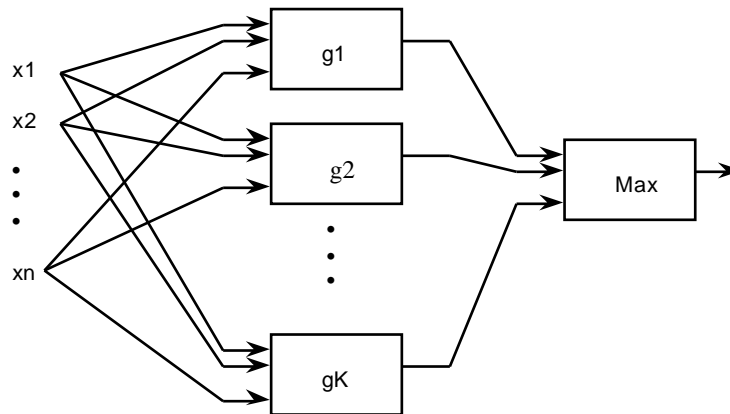
Donc notre fonction de discrimination devient :

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

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

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

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



Soit $D = 1$. (une seul 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.

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

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 :

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

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