

Formation et Analyse d'Images

James L. Crowley

ENSIMAG 3

Premier Bimestre 2008/2009

Séance 10

9 janvier 2009

Reconnaissance Bayésienne

Plan de la Séance :

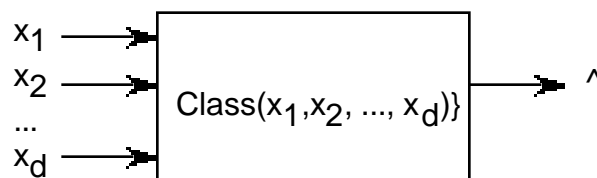
La Classification Bayésienne	2
Représentation de la Probabilité avec La Loi Normale	4
Estimations des moments d'une densité.....	5
Le premier moment : La Moyenne.....	5
Le deuxième moment (La variance).....	6
La Loi Normale pour $D = 1$	7
La Loi Normale pour $D > 1$	8
Forme en Algèbre Linéaire.....	11
Transformations Linéaire.....	12
Fonctions de Discrimination	13

La Classification Bayésienne

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

Soit les événements E décrit par une vecteur de caractéristiques $X : (E, X)$.
Soit K classes d'événements $\{T_k\} = \{T_1, T_2, \dots, T_K\}$

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



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

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

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

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

$$p(k) = p(E = k)$$

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

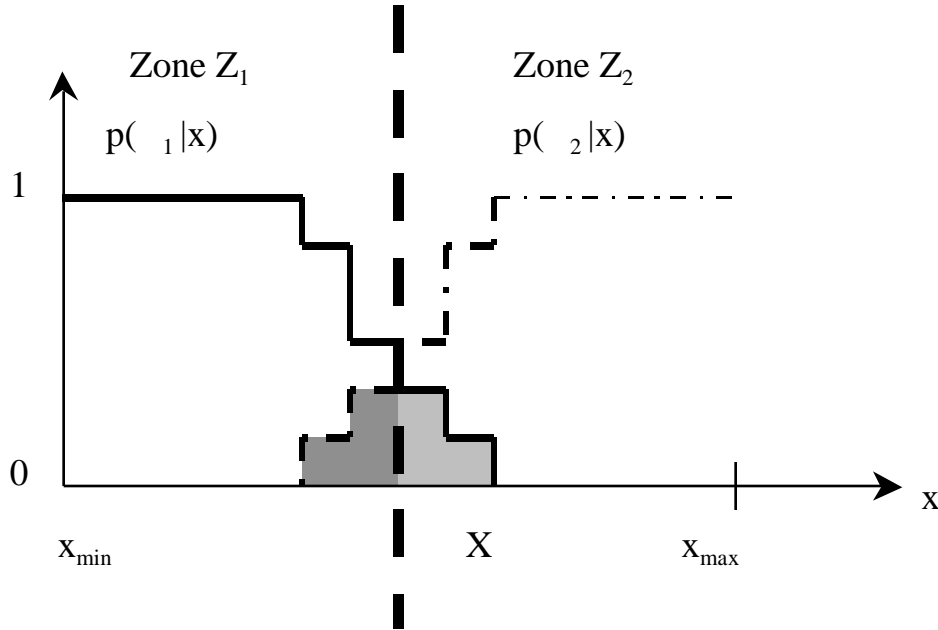
$$\hat{k} = \text{Decider}(E = k) = \arg\text{-max}_k \{p(k | X)\}$$

Pour estimer la probabilité nous utilisons les caractéristiques, X , de l'événement.

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

La classification est équivalente à un découpage du domaine d' X en deux zones : Z_1 et Z_2 .

$$\hat{c}_1 \text{ si } X \in Z_1 \text{ et } \hat{c}_2 \text{ si } X \in Z_2$$



La probabilité d'erreur est la somme des probabilités de $p(\hat{c}_2)$ en Z_1 et la somme de probabilité de $p(\hat{c}_1 | X)$ en zone 2.

$$p(\text{erreur}) = \int_{Z_1} p(\hat{c}_2 | X) + \int_{Z_2} p(\hat{c}_1 | X)$$

La minimum est atteint quand $d(g_k(X)) = \arg\text{-max}_k \{p(\hat{c}_k | X)\}$

Pour faire $g_k(X)$ nous allons utiliser :

1) La règle de Bayes :
$$p(\hat{c}_k | X) = \frac{P(X | \hat{c}_k) p(\hat{c}_k)}{P(X)}$$

2) La loi Normale :

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

Représentation de la Probabilité avec La Loi Normale

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

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

On écrit $p(X) = g(X)$.

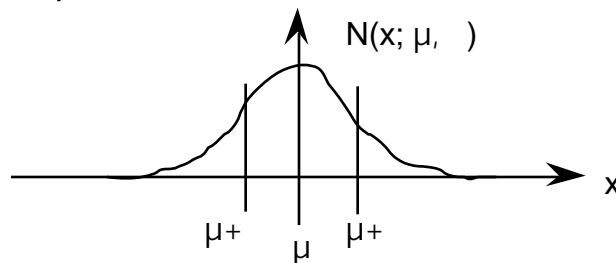
Quand les variables aléatoires sont issues d'une séquence d'événements aléatoires, leur densité de probabilité prend la forme de la loi normale, $\mathcal{N}(\mu, \sigma^2)$. Ceci est démontré par le théorème de la limite centrale. Il est un cas fréquent en nature.

Les paramètres de $\mathcal{N}(\mu, \sigma^2)$ sont les premiers et deuxième moments des exemples. Donc, on peut les estimer pour n'importe quel nombre d'exemples. 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 $p(x)$.

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



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

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

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

Le terme $d^2(x) = \frac{(x-\mu)^2}{2\sigma^2}$ est la distance x et μ normalisée par la variance.

La différence $(x - \mu)^2$ est la "distance" entre une caractéristique et la caractéristique "nominale" d'une classe. 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, $\{X_1, X_2, \dots, X_M\} = \{X_m\}$
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). Cette somme existe, même quand X_{\min} et X_{\max} n'existent pas et quand X est réelle.

On note que dans le cas X est un nombre entier, 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

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

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

Pour X continue :

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

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 !)

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 on a

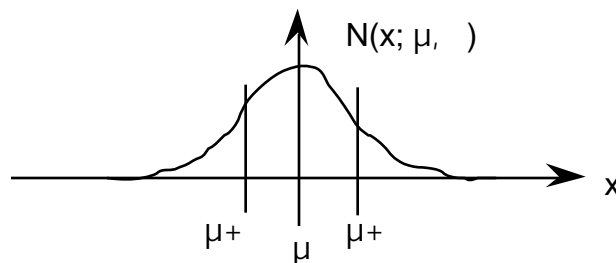
$$\sigma^2 = E\{(X_m - \mu)^2\} = \int p(x) \cdot (x - \mu)^2 \, dx$$

La Loi Normale pour $D = 1$

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

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

$\mathcal{N}(x; \mu, \sigma^2)$ a la forme :



La moyenne est le premier moment de la densité $p(x)$.

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

La variance 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 $D > 1$

Pour un vecteur de D propriétés

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

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

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 $\{X_m\}$

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

Et encore, pour éviter le bias, on peut utiliser :

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

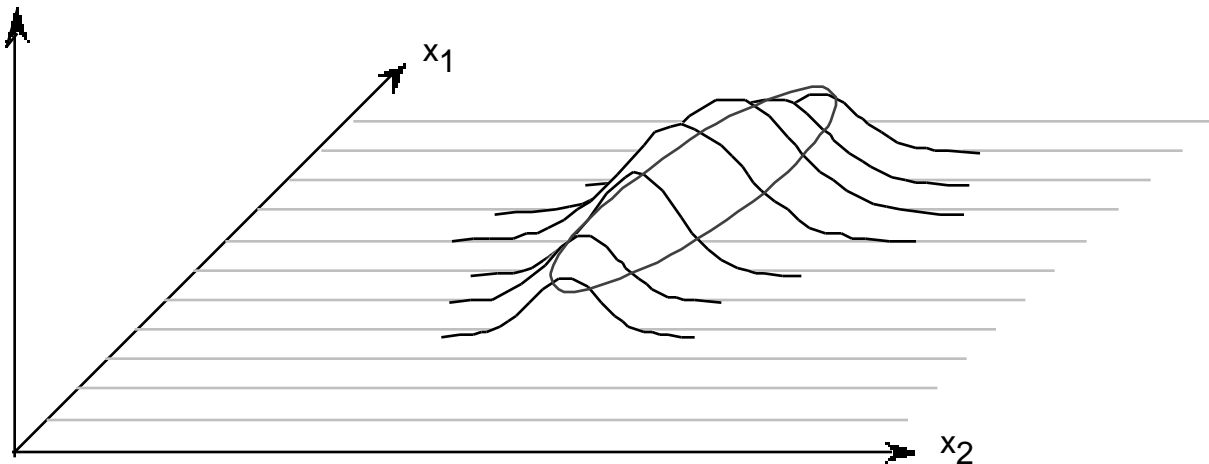
Ces coefficients composent une matrice de covariance. C

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

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

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

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



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

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

La déterminante, $\det(C)$ 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.)

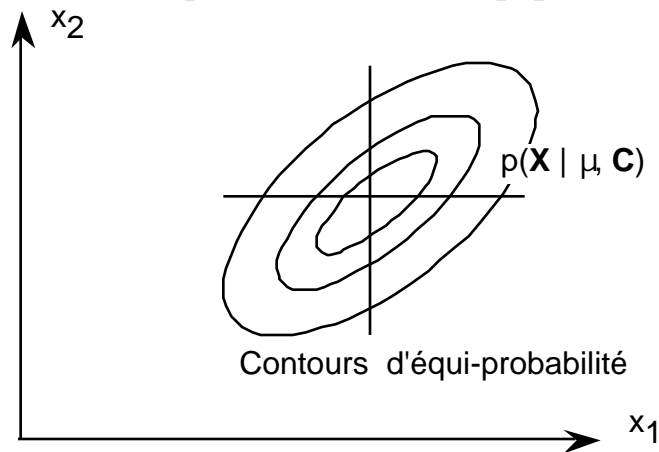
Cette valeur est connue comme la "distance de Mahalanobis".

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

Il s'agit d'une distance euclidienne, normalisé par la covariance 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 positive et semi-définie. Nous allons nous limiter au cas où C est positive et définie (C -à-d. $\det(C) = |C| > 0$)

Forme en Algèbre Linéaire

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 :

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

Soit V = X - E{X} = X - μ

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

Ceci peut etre exprimé en forme de matrice.

Soit V_m = X_m - μ

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

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

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

$$= \begin{matrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{matrix}$$

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 :

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

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, σ_y^2 dans la direction R .

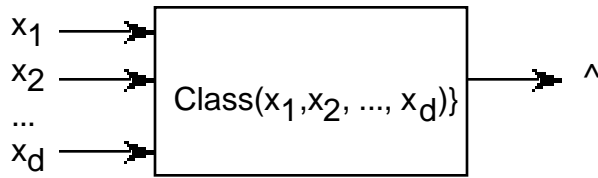
$$\mu_y = R^T \mu_X, \quad \sigma_y^2 = R^T C_X R$$

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

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

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

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

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

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

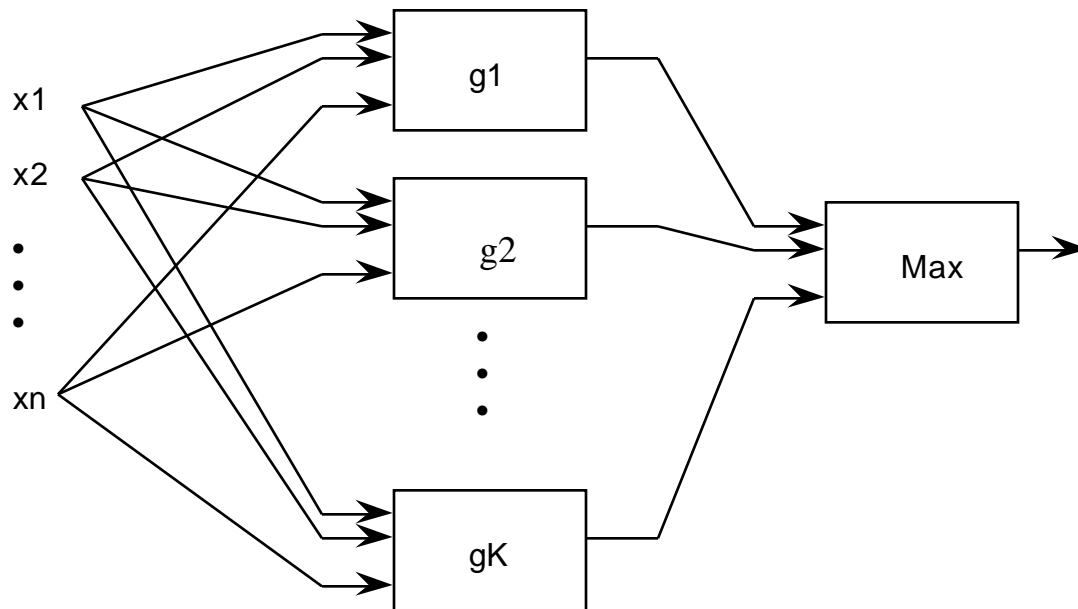
$$k = \arg\max_k \{g_k(X)\} = \arg\max_k \{p(k | X)\} = \arg\max_k \left\{ \frac{P(X|k)}{P(X)} p(k) \right\}$$

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

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

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

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



Fonctions classiques :

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

ou encore

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