

Systemes Intelligents : Raisonnement et Reconnaissance

James L. Crowley

Deuxième Année ENSIMAG

Deuxième Semestre 2005/2006

Séance 8

29 mar 2006

Reconnaissance Bayesienne

Notations	2
La Règle de Bayes.....	3
La règle de Bayes avec une ratio d'histogrammes.....	4
Exemple : Les statistiques de couleurs de la peau.....	5
Détection par ratio d'histogramme.....	7
La Classification.....	8
La Classification Bayesienne.....	9
Les Variations Aléatoires des Caractéristiques.....	10
La Loi Normale.....	13
La Loi Normale pour $D = 1$	15
La Loi Normale pour $D > 1$	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
M_k	Nombre d'exemples de la classe k .
M	Nombre totale d'exemples de toutes les classes
	$M = \sum_{k=1}^K M_k$
k	L'affirmation que l'événement E est dans la classe T_k
$h(x)$	Histogrammes des valeurs (x est entières avec range limité)
$h_k(x)$	Histogramme des valeurs pour la class k .
	$h(x) = \sum_{k=1}^K h_k(x)$
k	Proposition que l'événement E est dans la classe k
$p(k) = p(E = T_k)$	Probabilité que E est un membre de la classe k .
Y	La valeur d'une observation (un vecteur aléatoire).
$P(X)$	Densité de Probabilité pour X
$p(X = x)$	Probabilité q'un vecteur X prendre la valeur x
$P(X k)$	Densité de Probabilité pour X étant donné que k
	$P(X) = \sum_{k=1}^K p(X k) p(k)$

La Règle de Bayes

Soit un univers d'événement S , avec deux classes A et B . A et B ne sont pas disjoints. Soit deux propositions q et r .

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

$$p(q \wedge r) = (E=A) \wedge (E=B)$$

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

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

de la même manière :

$$p(r|q) = \frac{p(r \wedge q)}{p(q)} = \frac{p(r \wedge q)}{p(r \wedge q) + p(\neg r \wedge q)}$$

Par algèbre on déduit :

$$p(q|r) P(r) = p(r \wedge q) = p(r|q) P(q)$$

d'où

$$p(q|r) p(r) = p(r|q) p(q)$$

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

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

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

La règle de Bayes avec une ratio d'histogrammes.

On peut utiliser la règle de Bayes pour calculer la probabilité d'appartenance d'une classe. Soit un vecteur de caractéristiques, X discrètes tel que $x \in [X_{\min}, X_{\max}]$, et une ensemble de classe k .

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

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 :

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

Cette technique peut également marcher pour les vecteurs de caractéristiques.

La histogramme est une table à D dimensions : $h(x)$

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 :

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

Cette technique s'avère très utile dans les cas où il y a suffisamment d'échantillons pour faire un histogramme valable. Par exemple quand on traite des images ou les signaux.

Exemple : Les statistiques de couleurs de la peau

Une image est une table de pixels.

Chaque pixel est une observation d'une scène, et donc, une variable aléatoire.

Il y a beaucoup des pixels dans les images ($512 \times 512 = 2^{18} = 256 \text{ K}$ pixels)

Les pixels d'une image couleur sont représenté par 3 octets R, G et B avec (8 bits par octets). Dans ce cas, chaque pixel est une vecteur aléatoire.

$$X = (R, G, B)^T$$

ou R, G, et B sont issue du $[0, 255]$.

Pour un vecteur de caractéristique, on peut calculer une table à 3 dimensions.

Pour un image couleur, composé de (R, G, B), avec 8 bits par pixel, $h(X)$ contient $256^3 = 2^{24}$ valeurs. Mais chaque image contient $512^2 = 2^{18}$ pixels.

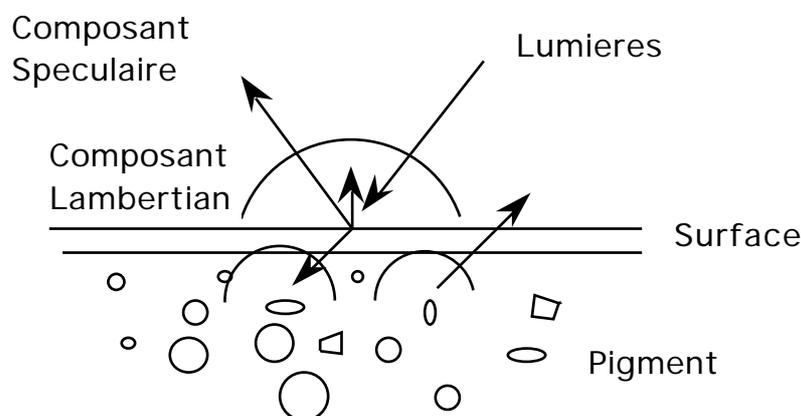
Si on suppose qu'il faut 10 exemples par cellule, Il faut 10×2^6 images = 640 images pour une estimation valable de $p(X) = \frac{1}{M} h(X)$.

A 25 images per second, ceci est 25.6 seconds de vidéo.

Si on souhaite, on peut reduire le nombre de dimensions, en normalisant la luminance.

On peut transformer le vecteur $(R, G, B)^T$ en luminance et chrominance.

La luminance, ou intensité, L, est en proportion de $\cos(i)$ où i est l'angle entre la source et la normale de la surface. La chrominance, C_1, C_2 est une signature pour la reconnaissance.



La composant "luminant" est déterminé par l'orientation de la surface.

La composant "chrominant" est déterminé par la composition de la spectre de la source et le spectre d'absorption des pigments de la surfaces. Si la spectre de la source est constante, la chrominance indique l'identité de l'objet

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

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

Donc pour chaque pixel (i,j) , $Y = (C_1, C_2)$

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 $M = 10$ K pixels d'exemples. (Une image = 256 K pixels).

Dans ce cas, pour M observations $p(X=x) = \frac{1}{M} h(x)$

Un histogramme de couleurs, $h(X)$, de les M pixels dans une l'image donne une approximation de la probabilité de chaque couleur dans l'image.

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

Un histogramme des de couleurs d'un entité, A , $h_A(X)$, de les M_A pixels dans une région d'une image de l'objet, $w(i, j)$, donne une approximation de la probabilité de chaque couleur de l'objet.

$$p(X | A) = p(X= x | A) = \frac{1}{M_A} h_A(x)$$

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 de chrominance $X(i, j)$: $p(\text{objet} | x) = \frac{p(x | \text{objet}) p(\text{objet})}{p(x)}$

Soit N images de M pixels. Ceci fait M Pixels.

Soit $h(c_1, c_2)$, l'histogramme de tous les M pixels.

Soit $h_A(c_1, c_2)$, l'histogramme des M_A pixels de l'objet "A".

$$p(c(i,j) = A) = \frac{M_A}{M} \quad p(x) = \frac{1}{M} h(x)$$

$$p(x | A) = \frac{1}{M_A} h_A(x)$$

$$\text{Donc } p(A | x) = \frac{p(x | A) p(A)}{p(x)} = \frac{1}{M_A} h_A(x) \frac{\frac{M_A}{M}}{\frac{1}{M} h(x)}$$

$$p(A | Y) = \frac{h_A(x)}{h(x)}$$

Il faut assurer que $h(x) \neq 0$!! Pour cela il faut que $w_A(i,j) = w(i,j)$.

Ainsi, on peut créer une image de probabilité de l'objet.

En peut ensuite déterminer un seuil pour l'image de probabilité.

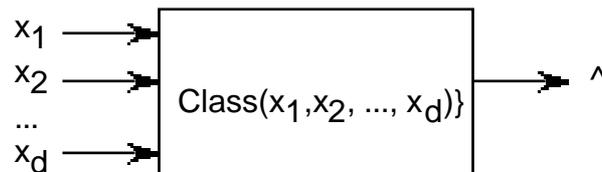


La Classification

Soit les événements E décrit par un 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 .



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

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

La Classification Bayesienne

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

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

$$p(k) = p(E_k)$$

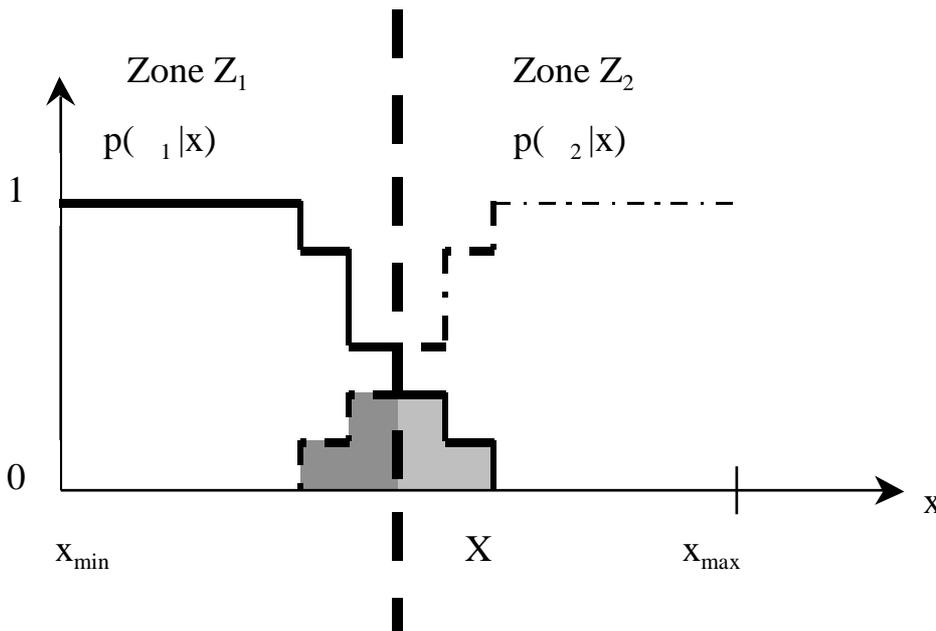
Le critère de décision est de minimiser le nombre d'erreur. Dans un système probabiliste, ca revient de minimiser la probabilité d'erreur. Ceci est équivalent à choisir la classe le plus probable.

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

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 une decoupage du domaine d' X en deux zones : Z_1 et Z_2 .

$$\hat{1} \text{ si } X \in Z_1 \text{ et } \hat{2} \text{ si } X \in Z_2$$



La probabilité d'erreur est la somme des probabilités de $p(2)$ en Z_1 et

la somme de probabilité de $p(z_1 | X)$ en zone 2.

$$p(\text{erreur}) = \int_{Z_1} p(z_2 | X) + \int_{Z_2} p(z_1 | X)$$

Le minimum est atteint quand :

$$\text{Donc } d(g_k(X)) = \arg\text{-max}_k \{p(z_k | X)\}$$

Dans ce cas, nous avons utilisé $\arg\text{-max}_k \{p(z_k | X)\}$ en tant que fonction de décision

et $g_k(X) = p(z_k | X)$ comme la fonction de discrimination

Les Variations Aléatoires des Caractéristiques

Les vecteurs de caractéristique porte une composante "aléatoire" avec (au moins) deux origines :

- 1) Les variations des individus.
- 2) Le bruit des capteurs

1) Les variations des individus.

La formation des vrais objets physiques est sujette aux influences aléatoires.

Pour les objets d'une classe, k , les propriétés des objets individuels sont,

les valeurs aléatoires. On peut résumer ceci par une somme d'une forme "intrinsèque"

x plus ces influences aléatoires individuelles, B_i .

$$X = x + B_i$$

Les plupart de techniques probabilistes de reconnaissance supposent un bruit additif.

En notation vectorielle :

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \dots \\ X_n \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + \begin{pmatrix} B_1 \\ B_2 \\ \dots \\ B_n \end{pmatrix}$$

Ceci n'est pas strictement nécessaire. Ils existe de méthodes pour la bruit non-additif

2) Le bruit des capteurs

(def.) 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'un 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 .

$$Y = X + B_o = x + B_i + B_o$$

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 probabilité d'observé Y est fournie par la règle de Bayes.

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

Si $Y = F(X + B_o)$ est issu de la classe k ayant caractéristique $X = x + B_i$

Une Exemple - Le spectre observées par un satellite.

Une image satellite est composée de pixels $s(x, y)$. Chaque pixel compte le nombre de photons issus d'une surface carré de la terre (ex. 10 m^2).

Les photons sont captés au travers des filtres spectraux. Ceci donne un vecteur de caractéristiques pour chaque pixel. Soit une région de végétation {blé, maize, etc}

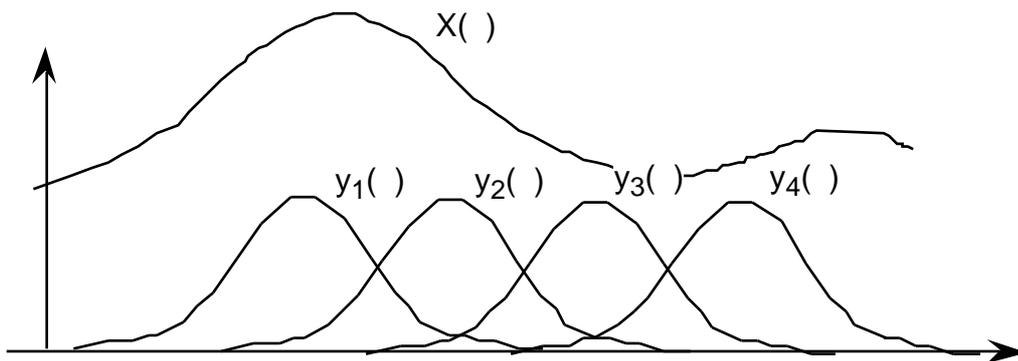
x : Le spectre des pigments des feuilles pour une espèce.

B_i : Les variations du spectre intrinsèque dû aux variations d'âge ou d'humidité.

B_i est spécifique à un individu. Il ne change pas entre les observations.

$X = x + B_i$: Le spectre des pigments des feuilles pour un individu

B_0 : Les variations d'observations dues à l'angle du soleil et les effets de filtrage de la lumière par l'atmosphère (humidité, pollution etc). Ils sont présents dans tous les pixels.



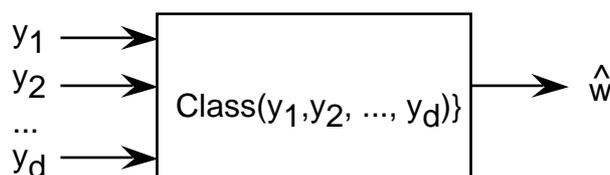
Une image est une table de pixels, avec les positions sur les lignes (r pour "row") et colonne (c). Un pixel, $Y(c,r)$, est un vecteur d'entiers, y_1, \dots, y_D . Chaque composant est l'intégral sur c, r et λ (longueurs d'onde) d'une produit d'un filtre spectral avec le spectre reçu sur la région c,r .

$$y_d(c,r) = \text{Quant} \left\{ \int_{c+dc}^{c+dc+dr} \int_{r+dr}^{r+dr+d} X(c,r, \lambda) \cdot f_d(\lambda) \, dc \, dr \, d\lambda \right\}$$

Cette opération est réalisée par l'optique de la caméra. L'opération "Quant{ }" numérise les valeurs y_d avec un pas "q" sur une plage entre 0 et v_{max} .

Une classification est une estimation de la classe de végétations dominante, k , pour la région observée par le pixel à partir de le vecteur d'observation \hat{Y}

$$\hat{k} = \text{Class} \{ \hat{Y} \}$$



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 $\{X_m\}$.

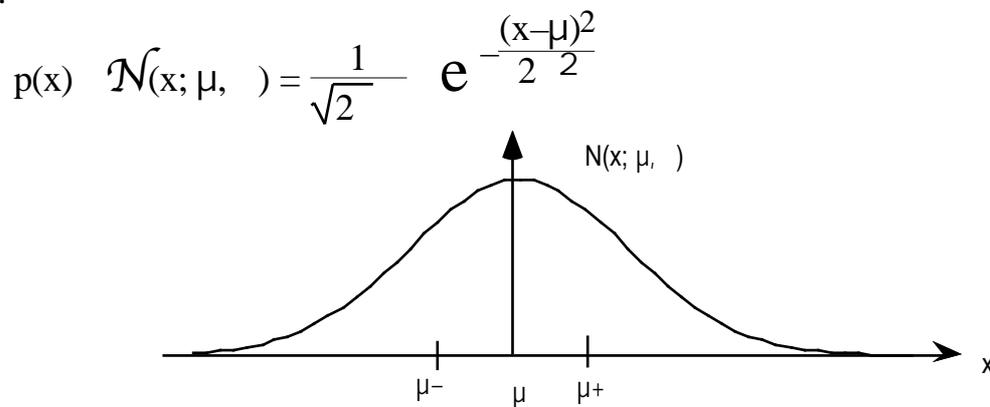
Les paramètres de $\mathcal{N}(\mu, \sigma^2)$ sont les premiers et deuxième moments de la population.

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 $p(x)$.



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.

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

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

La différence $(x - \mu)^2$ est la "distance" entre une caractéristique et la moyenne de l'ensemble $\{X_m\}$. La variance, σ^2 , sert à "normaliser" cette distance.

La différence normalisée par la variance est connue sous le nom de "Distance de Mahalanobis". La Distance de Mahalanobis est un test naturel de similarité

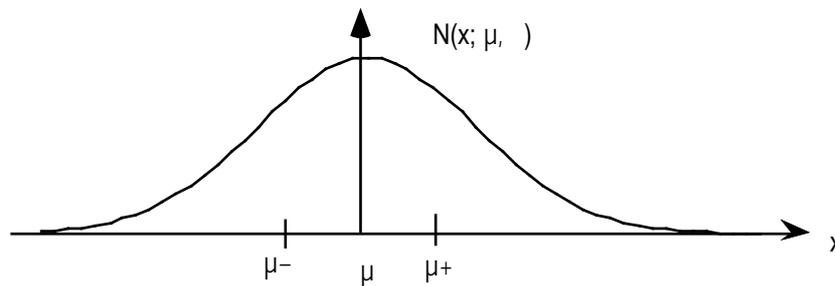
La Loi Normale pour D = 1

La cas le plus simple concerne une seule caractéristique.

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

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

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

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

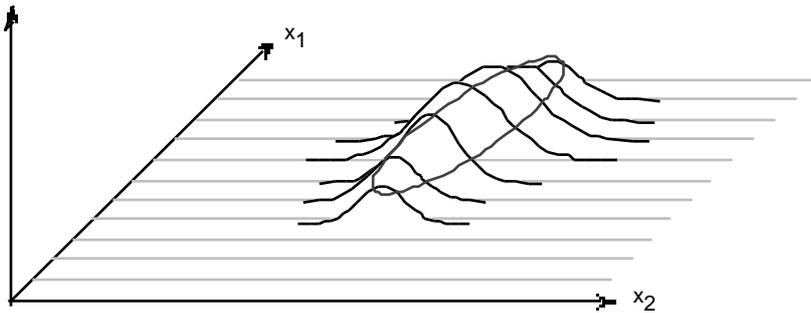
La Loi Normale pour D > 1

Pour un vecteur de D caractéristiques :

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

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

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



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

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

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 quadratique.

(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, normalisée 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).