

Systemes Intelligents : Raisonnement et Reconnaissance

James L. Crowley

Deuxième Année ENSIMAG

Deuxième Semestre 2006/2007

Séance 10

25 avril 2007

Mélange de Gaussiennes et L'Algorithme EM (Expectation-Maximization)

L'Analyse en Composantes Principales	2
Exemple :	6
Reconstruction.....	8
Reconnaissance de visages avec PCA.....	9
Mélange de Gaussiens.	10
Ebauche de l'Algorithme EM.....	11
Maximisation de la Vraisemblance	12
Maximum de vraisemblance pour le cas univariate.....	13
Le cas multi-variate.	15
L'Algorithme EM (Expectation-Maximization).....	16
L'etape E.....	17
L'etape M	18

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.

L'Analyse en Composantes Principales

L'analyse en composantes principales est une méthode de déterminer un sous-espace "optimale" pour la reconstruction. Il est souvent utilisé pour la discrimination dans la cas où le vecteur X est composé d'un grand nombre de caractéristiques.

Voici son application pour la reconnaissance d'images de visages.

Soit un ensemble de M images, toutes classes confondues, $W_m(i, j)$, composé de N pixels :

$W_m(i, j)$ pour $m \in \{1, M\}$, tel que $i \in [1, I], j \in [1, J], I \times J = N$

On les exprime sous forme de vecteur :

Ensemble d'images $X_m(n) = W_m(i, j)$ ou $n = j \cdot I + i$.

On cherche une base orthogonale $(n) = \{d(n)\}$ $d = 1, 2, \dots, D$ pour représenter les $X(n)$. telles que $D \ll M$.

Image moyenne $\mu(n) = \frac{1}{M} \sum_{m=1}^M X_m(n)$

Image moyenne zéro $V_m(n) = X_m(n) - \mu(n)$

Base orthogonale $(n) = \{d(n)\}$ $d = 1, 2, \dots, D$
avec $D \ll N$.

Vecteur $= \langle V(n), (n) \rangle$

Image reconstruite : $\hat{X}(n) = \mu(n) + \sum_{n=1}^D v_n(n)$

Image de résidu : $R(n) = X(n) - \hat{X}(n)$

Energie de résidu : $r^2 = \sum_{n=1}^D R^2(n)$

La covariance, C , est composée de $N \times N = N^2$ termes

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

(pour une image de taille $D = 2^N$, il y a 2^{2N} pixels et 2^{4N} termes dans la covariance.)

$$c_{ij}^2 = \frac{1}{M} \sum_{m=1}^M (V_{m(i)} \cdot V_{m(j)})$$

On peut faire une matrix V composé de M colonnes $\{V_m\}$

$$V = \begin{matrix} V_{11} & V_{12} & \dots & V_{1M} \\ V_{21} & V_{22} & \dots & V_{2M} \\ \dots & \dots & \dots & \dots \\ V_{N1} & V_{N2} & \dots & V_{NM} \end{matrix}$$

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

V est N par M . Chaque colonne de V est une image. Chaque coefficient, c_{ij}^2 , est une covariance de la pixel i et j pour l'ensemble de M images

Pour une image de 32×32 , le matrice $C = V V^T$ est de taille 1024×1024 . Il y a un coefficient par pair de pixels. Chaque terme est la covariance d'une paire de pixels.

$$C = V V^T \text{ est de taille } N \times N$$

On cherche une ensemble orthogonales $(n) = \{d(n) \mid d = 1, 2, \dots, N\}$ tels que : $V^T V V^T = c$

c est une matrice diagonale des valeurs principales de C .

Chaque colonne de V est un vecteur directeur $d(n)$. Les colonnes d sont orthogonales. Nous allons choisir une sous ensemble D des N colonnes de V .

Une telle matrice de rotation est fournie par une procédure d'analyse en composants principales.

$$(d(n), c) = \text{PCA}(C).$$

Pour une image de 32×32 , le matrice $C = V V^T$ est de taille $2^{10} \times 2^{10} = 2^{20}$.

Pour une image de 512×512 , le matrice $C = V V^T$ est de taille $2^{18} \times 2^{18} = 2^{36}$.

Heureusement, il y a une astuce pour éviter la matrice de covariance de $N \times N$ coefficients. Le rank de $V V^T$ est M , ou M est le nombre d'images. $M \ll N$.

Noter que $B = V^T V$ est de taille M^2 , M est le nombre d'images.

Chaque coefficient est une produit de deux images !

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

Les coefficients sont les covariances des images (et non pas les pixels!).

$$b_{ij}^2 = \sum_{n=1}^N V_{i(n)} V_{j(n)}$$

Pour $M < 512$ on peut facilement calculer une matrice R de rotation tels que chaque colonne est un vecteur directeur orthogonal.

Soit R les composantes principales de $\tilde{X}^T \tilde{X}$

$$R^T V^T V R = b$$

On multiplie les deux cotés par R :

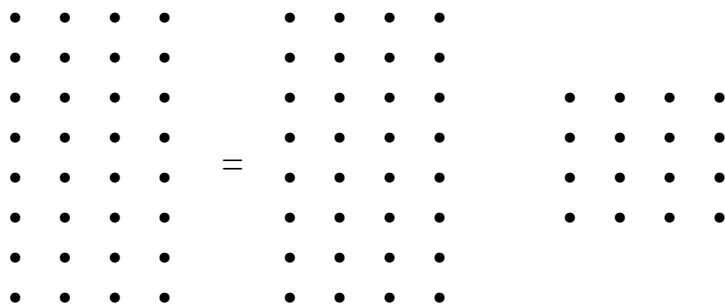
$$R R^T V^T V R = V^T V R = R \quad b$$

$R^T R = I$. Maintenant on multiplie par A

$$\begin{aligned} V V^T V R &= V R \quad b \\ &= (V V^T)(V R) = (V R) \quad b \\ &= (V V^T) \quad = \quad c \end{aligned}$$

Donc, par inspection : $\quad = VR$.

$$\quad = V \quad R$$



Raison : Les mêmes images ont été utilisé pour \quad et pour R.

Les premiers M vecteurs propres de $V^T V$ sont aussi les premiers M vecteurs propres de $V V^T$, le Rank de \quad est le même que R, et $\quad b$ sont les premiers M valeurs propre de $\quad c$ et Donc, les vecteurs propres : $\quad d(n) = V R$ triés par $\quad b$

$$\quad = VR$$

Chaque colonne est un vecteur dans une base orthogonale.

$$d(n) = \sum_{m=1}^M V_{m(n)} R(m,n)$$

$d(n)$ fournit une base "ortho-normal" pour $X_m(n)$.

Une image (normalisée) peut être exprimé par

$$d = \langle X(n), d(n) \rangle = \sum_{d=1}^M X(n) d(n)$$

Les valeurs d sont un "code" qui représente $X(n)$ pour la reconnaissance ou la transmission.

Exemple :

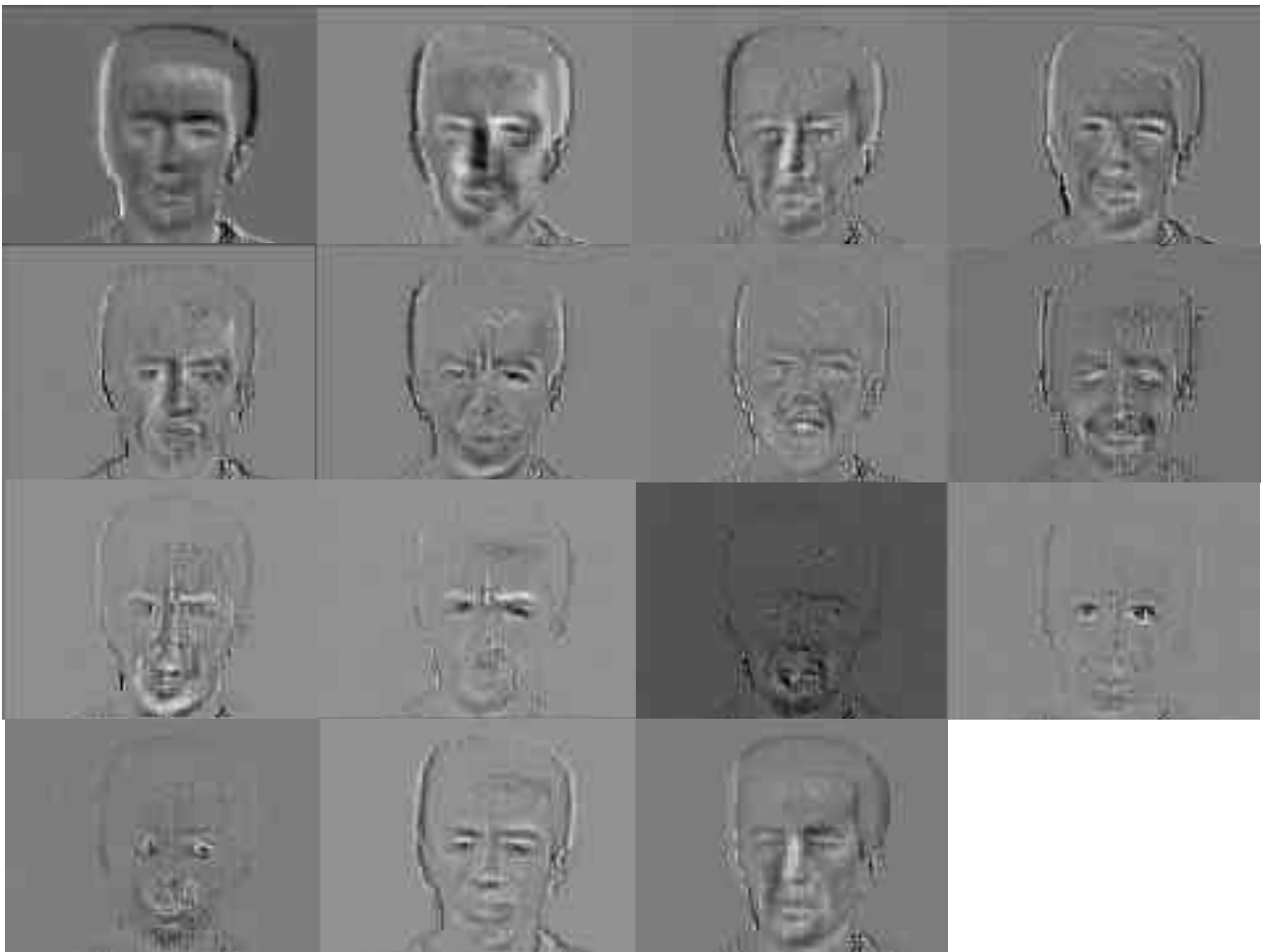
16 images pris au hasard dans une séquence de 2 minutes. (F. Bérard, 1995).



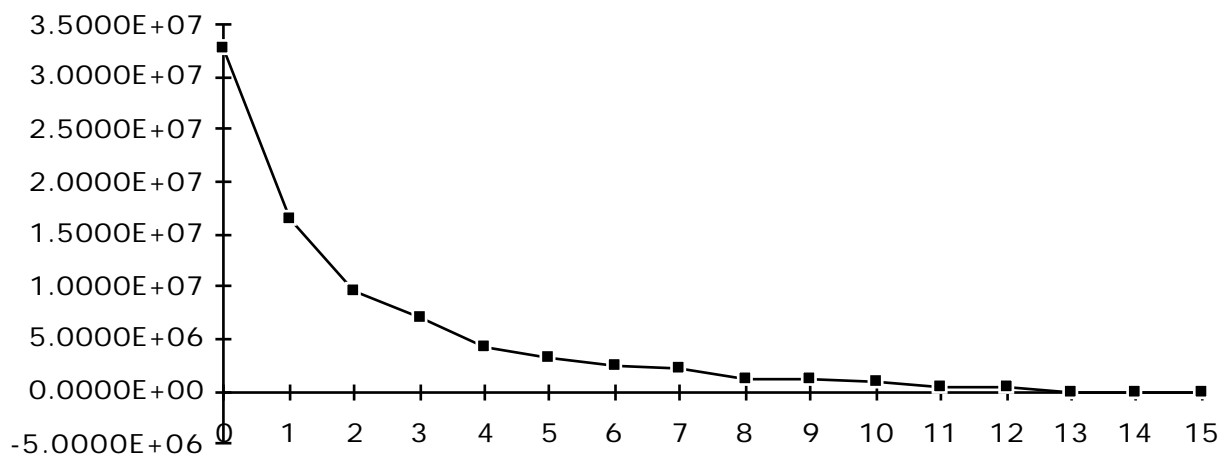
Average Image



Components Principales :



Eigen Values.



Reconstruction

Image Reconstructed image (120 bytes) Error Image.



Reconstruction (120 bytes)



Image Error

Reconnaissance de visages avec PCA

Dans l'espace PCA, pour chaque individu, on fait une ensemble de M_k images.

Ensemble d'images de la classe k : $X_{mk}(n) = W_{mk}(i, j)$ ou $n = j*32 + i$.

Image moyenne de la classe T_k : $\mu_k(n) = \frac{1}{M} \sum_{m=1}^M X_{mk}(n)$

Image moyenne zéro $V_{mk}(n) = X_{mk}(n) - \mu_k(n)$

Projection sur : $Z_{mk} = \langle V_{mk}, \rangle$

La covariance, C_k , est composée de $D_y \times D_y$ termes

$$C_k = E\{ Z_{mk} Z_{mk}^T \}$$

(pour une image de taille 2^n , il y a 2^{2n} pixels et 2^{4n} termes dans la covariance.)

$$ij^2 = \frac{1}{M} \sum_{m=1}^M (Z_{mk}(i) Z_{mk}(j))$$

ou bien :

$$ij^2 = \frac{1}{M} \sum_{m=1}^M (Z_{mk}(i) - \mu_k(i))(Z_{mk}(j) - \mu_k(j))$$

$$k = \arg\text{-max}_k \{ g_k(Z) \}$$

ou $g_k(Z) = Z^T (B_k) Z + b_k^T Z + b_{k0}$.

avec $B_k = \frac{1}{2} C_k^{-1}$

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

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

Quelle sont les sources de bruit? $\det(B_0) \gg \det(B_k)$ est dominé par

- 1) L'orientation du visage
- 2) Les ombres, et
- 3) Les expressions.

Ceci donne plusieurs densité Gaussiennes. Comment les estimer?

Mélange de Gaussiens.

Si les événements sont issus d'une composition de "N" phénomènes, la densité $p(X)$ prendra la forme d'une composition de lois Normales.

Dans un tel cas, on peut approximer par une somme pondérée de densités Normales.

$$p(x) = \sum_{n=1}^N \pi_n \mathcal{N}(x; \mu_n, \sigma_n)$$

Un tel somme est connu par la terme "mélange de Gaussiens" pour chaque Gaussien il faut estimer trois paramètres :

$$\theta_n = (\pi_n, \mu_n, \sigma_n)$$

En total il y a $3 \cdot N$ paramètres à estimer.

$$\theta = (\pi_1, \mu_1, \sigma_1, \pi_2, \mu_2, \sigma_2, \dots, \pi_N, \mu_N, \sigma_N)$$

Si toutes les μ et σ était fixe (et égale), on pouvait calculer les π_n directement.

Mais parce qu'ils sont libres il faut les estimer par un processus itérative.

Un tel processus est composé de deux étapes.

Ces étapes fournissent une estimation des variables cachées.

Pour un mélange de Gaussiens, les variables cachées sont les "sources" pour des événements.

On suppose chaque événement est issu d'un des N sources.

Nous allons construire une table de probabilités. $h(m, n)$

$$h(m, n) = \Pr\{\text{l'événement } E_m \text{ est issu de la source } N\}$$

les probabilités, $h(m, n)$, nous donnera les facteurs de Mélange, π_n , ainsi que μ_n, σ_n . L'algorithme d'estimation s'appelle "Expectation-Maximization" ou "EM".

Ebauche de l'Algorithme EM

Soit un ensemble ("training set") de M observations $\mathbf{T} = \{X_m\}$.

Fait une première estimation des paramètres $^{(0)}$ et puis on alterne "E" et "M".

E: Faire une estimation des valeurs manquantes, $h(m, n)$, pour les événements.

$h(m, n)^{(i)} = p(h_m | X_1, X_2, \dots, X_M, \quad ^{(i)})$ pour chaque terme "n".

$$h(m, n)^{(i)} = \frac{n^{(i)} \mathcal{N}(X_m; \mu_n^{(i)}, \sigma_n^{(i)})}{\sum_{j=1}^N j^{(i)} \mathcal{N}(X_m; \mu_j^{(i)}, \sigma_j^{(i)})}$$

M: Recalculer $^{(i+1)}$ avec $p(h_m | X_1, X_2, \dots, X_M, \quad ^{(i)})$

$$S_n^{(i+1)} = \sum_{m=1}^M h(m, n)^{(i)}$$

$$\hat{n}^{(i+1)} = \frac{1}{M} S_n^{(i+1)} = \frac{1}{M} \sum_{m=1}^M h(m, n)^{(i)}$$

$$\hat{\mu}_n^{(i+1)} = \frac{1}{S_n^{(i+1)}} \sum_{m=1}^M h(m, n)^{(i)} X_m$$

$$\hat{\sigma}_n^{(i+1)} = \frac{1}{S_n^{(i+1)}} \sum_{m=1}^M h(m, n)^{(i)} (X_m - \hat{\mu}_n^{(i+1)})^2$$

Pour la dérivation l'algorithme EM, il faut introduire le concept de "likelihood" (vraisemblance).

Maximisation de la Vraisemblance

Nous commençons par le cas où $D = 1$. (une seule caractéristique)

Soit un ensemble de M exemples de caractéristiques $\{X_m\}$.

Les exemples sont supposés d'être les échantillons independents.

On souhaite estimer une forme paramétrique pour $p(x)$.

Supposons que $p(x) = \mathcal{N}(x; \mu, \sigma)$.

Pour une loi normale avec $D = 1$ les paramètres sont

$$= (\mu, \sigma)$$

On souhaite trouver une estimation $\hat{\theta} = (\hat{\mu}, \hat{\sigma})$ qui minimise la probabilité d'erreur.

Pour ceci on définit le Likelihood $L(\theta | X_1, X_2, \dots, X_M)$ de θ étant donnée $\{X_m\}$

Si les X_m sont indépendants,

$$p(X_1, X_2 | \theta) = p(X_1 | \theta) \cdot p(X_2 | \theta)$$

En général pour M événements :

$$p(\{X_m\} | \theta) = p(X_1, X_2, \dots, X_M | \theta) = \prod_{m=1}^M p(X_m | \theta)$$

Nous définissons le "Likelihood" (vraisemblance) de θ étant $\{X_m\}$ comme

$$\begin{aligned} L(\theta | \{X_m\}) &= L(\theta | X_1, X_2, \dots, X_M) = p(X_1, X_2, \dots, X_M | \theta) \\ &= \prod_{m=1}^M p(X_m | \theta) \end{aligned}$$

Notre objectif est d'estimer les paramètres $\hat{\theta}$ pour maximiser $L(\theta | S)$.

$$\hat{\theta} = \arg\text{-max} \{ L(\theta | \{X_m\}) \} = \arg\text{-max} \left\{ \prod_{m=1}^M p(X_m | \theta) \right\}$$

Pour simplifier l'analyse, on travaille avec la log : $\ell(\theta) = \text{Log}\{L(\theta | S)\}$.

$$\ell(\theta) = \text{Log}\{L(\theta | \{X_m\})\} = \text{Log} \{ p(\{X_m\} | \theta) \} = \sum_{m=1}^M \text{Log}\{p(X_m | \theta)\}$$

Si $p(X_m | \theta)$ est une simple loi normale, il suffit de trouver

$$\hat{\theta} = \arg\text{-max}_{\theta} \left\{ \sum_{m=1}^M p(X_m | \theta) \right\}$$

Ceci nous donne les formes habituelles des moments $\hat{\theta} = (\mu, \sigma)$

Maximum de vraisemblance pour le cas univariate

Soit un modèle normale $\mathcal{N}(\mu, \sigma)$. Pour estimer μ, σ :

$$\ell(\theta) = \text{Log}\{p(X_m | \theta)\} = -\frac{1}{2} \text{Log}\{2\sigma^2\} - \frac{1}{2\sigma^2} (X_m - \mu)^2$$

$$\frac{\partial \ell(\theta)}{\partial \mu} = \frac{1}{\sigma^2} (X_m - \mu)$$

$$\frac{\partial \ell(\theta)}{\partial \sigma^2} = -\frac{1}{2\sigma^2} + \frac{(X_m - \mu)^2}{2\sigma^4}$$

$$\mu, \frac{\partial \ell(\theta)}{\partial \sigma^2} = \frac{1}{\sigma^2} (X_m - \mu) - \frac{1}{2\sigma^2} + \frac{(X_m - \mu)^2}{2\sigma^4}$$

La maximum se trouve où le dérivé est nul.

$$\frac{\partial \ell(\theta)}{\partial \mu} = \sum_{m=1}^M \frac{1}{\sigma^2} (X_m - \hat{\mu}) = 0.$$

Avec un peu d'algèbre on a :

$$\sum_{m=1}^M \frac{1}{2} (X_m - \hat{\mu}) = 0.$$

$$\frac{1}{2} \sum_{m=1}^M X_m = \frac{1}{2} \sum_{m=1}^M \hat{\mu}$$

$$\sum_{m=1}^M X_m = M \hat{\mu}$$

$$\hat{\mu} = \frac{1}{M} \sum_{m=1}^M X_m$$

et de la même façon pour

$$\frac{\ell(\cdot)}{2} = -\frac{1}{2 \cdot 2} + \frac{(X_m - \mu)^2}{2 \cdot 4} = 0$$

$$\sum_{m=1}^M -\frac{1}{2 \cdot 2} + \frac{(X_m - \mu)^2}{2 \cdot 4} = 0$$

$$\sum_{m=1}^M \frac{1}{2 \cdot 2} = \sum_{m=1}^M \frac{(X_m - \mu)^2}{2 \cdot 4}$$

$$\frac{1}{2 \cdot 2} \sum_{m=1}^M 1 = \frac{1}{2 \cdot 4} \sum_{m=1}^M (X_m - \mu)^2$$

$$M = \frac{1}{2} \sum_{m=1}^M (X_m - \mu)^2 \quad \hat{\mu}^2 = \frac{1}{M} \sum_{m=1}^M (X_m - \hat{\mu})^2$$

Le cas multi-variate.

Pour \mathbf{x} composé de D caractéristiques, avec M exemples d'une classe $\mathbf{T} = \{\mathbf{X}_m\}$
 Issues un modèle normale $\mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$.

Le problème est d'estimer les

$$\hat{\boldsymbol{\mu}}, \hat{\mathbf{C}} = \max_{\boldsymbol{\mu}, \mathbf{C}} \{L(\boldsymbol{\mu}, \mathbf{C} | \mathbf{T})\} = \max_{\boldsymbol{\mu}, \mathbf{C}} \left\{ \sum_{m=1}^M \log p(\mathbf{X}_m | \boldsymbol{\mu}, \mathbf{C}) \right\}$$

Le maximum de $L(\boldsymbol{\mu}, \mathbf{C})$ est trouvé quand le gradient est null.

$$\frac{\partial L(\boldsymbol{\mu}, \mathbf{C})}{\partial \boldsymbol{\mu}} = \sum_{m=1}^M \frac{\partial \log p(\mathbf{X}_m | \boldsymbol{\mu}, \mathbf{C})}{\partial \boldsymbol{\mu}} = 0$$

$$\text{où le gradient est } \frac{\partial \log p(\mathbf{X}_m | \boldsymbol{\mu}, \mathbf{C})}{\partial \boldsymbol{\mu}} = \frac{1}{\mathbf{C}} (\mathbf{X}_m - \boldsymbol{\mu})$$

Le gradient nous permet de calculer une solution analytique.

Les estimations de $\hat{\boldsymbol{\mu}}, \hat{\mathbf{C}}$ sont obtenus par

$$\log p(\mathbf{X}_n | \boldsymbol{\mu}, \mathbf{C}) = -\frac{1}{2} \log \{2^{-n} \det(\mathbf{C})\} - \frac{1}{2} (\mathbf{X}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{X}_n - \boldsymbol{\mu})$$

$$\sum_{m=1}^M \frac{\partial \log p(\mathbf{X}_m | \boldsymbol{\mu}, \mathbf{C})}{\partial \boldsymbol{\mu}} = 0$$

Donne la formule classique :

$$\hat{\boldsymbol{\mu}} = \frac{1}{M} \sum_{m=1}^M \mathbf{X}_m \quad \text{et} \quad \hat{\mathbf{C}} = \frac{1}{M} \sum_{m=1}^M (\mathbf{X}_m - \hat{\boldsymbol{\mu}}) (\mathbf{X}_m - \hat{\boldsymbol{\mu}})^T$$

L'Algorithme EM (Expectation-Maximization)

L'algorithme EM s'applique à l'estimation de données cachés.

Il est utilisé notamment pour l'estimation des Modèles de Markov Cachées (HMM's) et pour l'estimation des Mélanges de Gaussiens.

Pour un mélange de Gaussiens, les variables cachées sont les sources des événements. On suppose que chaque événement est produit par un des N sources.

Pour chaque événement E_m On définit la variable "cachée" h_m

$h_m = N$ si l'événement E_m (avec caractéristique X_m) est issu de la source N .

Nous cherchons à estimer

$$p(x) = \prod_{n=1}^N \mathcal{N}(x; \mu_n, \sigma_n) \quad \text{Tels que} \quad \sum_{n=1}^N \pi_n = 1.$$

Il faut estimer N vecteurs : $\theta_n = (\pi_n, \mu_n, \sigma_n)$

En total il y a $3 \cdot N$ paramètres à estimer.

$$\theta = (\pi_1, \mu_1, \sigma_1, \pi_2, \mu_2, \sigma_2, \dots, \pi_N, \mu_N, \sigma_N)$$

$$\ell(\theta | \{X_m\}) = \text{Log} (L(\theta | X_1, X_2, \dots, X_M)) = \text{Log} \left\{ \prod_{m=1}^M p(X_m | \theta) \right\}$$

$$= \sum_{m=1}^M \text{Log} \{p(X_m | \theta)\}$$

$$= \sum_{m=1}^M \text{Log} \left(\sum_{n=1}^N \pi_n p(X_m | \theta_n) \right)$$

Il faut $\hat{\theta}_S = \max_{\theta} \{ \ell(\theta | \{X_m\}) \}$

Il n'y pas de solution analytique, parce qu'il n'y est pas une dérivé pour la logarithme de la somme.

Soit $\mathbf{T} = \{X_m\}$ $h(m,n)$ le donnée Complète.

On vas maximiser une mesure de qualité $Q(\theta, \theta^{(i)})$ définit par une esperence conditionnel :

$$Q(\theta, \theta^{(i)}) = E \{ \ell(\theta | \mathbf{T}) | \{X_m\}, \theta^{(i)} \}$$

L'étape E

Pour chaque événement E_m avec caractéristique X_m ,

On suppose qu'il manque l'information: h_m

$h_m = n$, la source de l'événement E_m .

On ne connaît pas h_m , mais on peut estimer les probabilités $P(h_m = n)$.
par une table $h(m, n)$.

Pour chaque X_m , et son source caché h_m

$$p(h_m, X_m | \theta) = p(h_m | X_m, \theta) p(X_m | \theta)$$

donc

$$p(h_m | X_m, \theta) = \frac{p(h_m, X_m | \theta)}{p(X_m | \theta)}$$

où

$$p(h_m = n, X_m | \theta) = \mathcal{N}(X_m; \mu_n, \sigma_n)$$

$$p(X_m | \theta) = \sum_{n=1}^N \mathcal{N}(X_m; \mu_n, \sigma_n)$$

Donc

$$p(h_m = n | X_1, X_2, \dots, X_M, \theta) = \frac{\mathcal{N}(X_m; \mu_n, \sigma_n)}{\sum_{j=1}^N \mathcal{N}(X_m; \mu_j, \sigma_j)}$$

Donc pour chaque itération (i) le premier étape E est :

$$h(m, n)^{(i)} = \frac{\mathcal{N}(X_m; \mu_n^{(i)}, \sigma_n^{(i)})}{\sum_{j=1}^N \mathcal{N}(X_m; \mu_j^{(i)}, \sigma_j^{(i)})}$$

L'étape M

Pour la deuxième étape, M, nous allons maximiser le "likelihood" $\text{Log}\{p(\mathbf{X}_m | \cdot)\}$
 On définit la "Expected Complete Data Log Likelihood" pour $\mathbf{T} = \{\mathbf{X}_m\} \cup \mathbf{h}(m,n)$

$$Q(\cdot, \cdot^{(i)}) = E\{\ell(\cdot | \mathbf{T}) | \{\mathbf{X}_m\}, \cdot^{(i)}\}$$

$$Q(\cdot, \cdot^{(i)}) = E\{\ell(\cdot | \{\mathbf{X}_m\}, \mathbf{h}(m,n)) | \{\mathbf{X}_m\}, \cdot^{(i)}\}$$

$$Q(\cdot, \cdot^{(i)}) = \sum_{m=1}^M \sum_{n=1}^N \text{Log}\{p(\mathbf{h}_{m=n}, \mathbf{X}_m | \cdot)\} p(\mathbf{h}_{m=n} | \mathbf{X}_m, \cdot^{(i)})$$

Pour chaque cycle dans l'itération nous allons chercher :

$$\cdot^{(i+1)} = \text{argmax}\{Q(\cdot | \cdot^{(i)})\}$$

Ce maximum est donné par :

$$S_n^{(i+1)} = \sum_{m=1}^M p(\mathbf{h}_{m=n} | \mathbf{X}_m, \cdot^{(i)}) = \sum_{m=1}^M h(m,n)$$

$$p_n^{(i+1)} = \frac{1}{M} \sum_{m=1}^M P(\mathbf{h}_{m=n} | \mathbf{X}_m, \cdot^{(i)}) = \frac{1}{M} S_n^{(i+1)}$$

$$\mu_n^{(i+1)} = \frac{1}{S_n^{(i)}} \sum_{m=1}^M p(\mathbf{h}_{m=n} | \mathbf{X}_m, \cdot^{(i)}) \mathbf{X}_m = \frac{1}{S_n^{(i+1)}} \sum_{m=1}^M h(m,n) \mathbf{X}_m$$

$$\begin{aligned} \sigma_n^{(i+1)} &= \frac{1}{S_n^{(i)}} \sum_{m=1}^M p(\mathbf{h}_{m=n} | \mathbf{X}_m, \cdot^{(i)}) (\mathbf{X}_m - \mu_n^{(i+1)})^2 \\ &= \frac{1}{S_n^{(i+1)}} \sum_{m=1}^M h(m,n) (\mathbf{X}_m - \mu_n^{(i+1)})^2 \end{aligned}$$

Avec notre table de probabilités $h(m,n)$

$$h(m,n) = P(\mathbf{h}_{m=n} | \mathbf{X}_m, \cdot^{(i)})$$

E (Expectation) :

$$h(m, n)^{(i)} := \frac{n^{(i)} \mathcal{N}(X_m; \mu_n^{(i)}, \sigma_n^{(i)})}{\sum_{j=1}^M n^{(i)} \mathcal{N}(X_m; \mu_j^{(i)}, \sigma_j^{(i)})}$$

M: (Maximisation)

$$S_n^{(i+1)} := \sum_{m=1}^M h(m, n)^{(i)}$$

$$n^{(i+1)} := \frac{1}{M} S_n^{(i+1)}$$

$$\mu_n^{(i+1)} := \frac{1}{S_n^{(i+1)}} \sum_{m=1}^M h(m, n)^{(i)} X_m$$

$$\sigma_n^{(i+1)} := \frac{1}{S_n^{(i+1)}} \sum_{m=1}^M h(m, n)^{(i)} (X_m - \mu_n^{(i+1)})^2$$

Dans le cas Multivariate ($D > 1$) la covariance C est composée de σ_{jk}^2 :

$$\sigma_{jkn}^{(i+1)} := \frac{1}{S_n^{(i+1)}} \sum_{m=1}^M h(m, n)^{(i)} (X_{jm} - \mu_{jn}^{(i+1)})(X_{km} - \mu_{kn}^{(i+1)})$$