

Systèmes Intelligents : Raisonnement et Reconnaissance

James L. Crowley

Deuxième Année ENSIMAG

Deuxième Sémestre 2008/2009

Séance 10

4 mai 2009

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

Mélange de Gaussiens.	2
Ebauche de l'Algorithme EM.....	3
Maximisation de la Vraisemblance.....	4
Maximum de vraisemblance pour le cas univariate.....	5
Le cas multi-variate.....	7
L'Algorithme EM (Expectation-Maximization).....	8
L'étape E.....	9
L'étape M.....	10

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.

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.

$$p(x) = \sum_{n=1}^N \pi_n \mathcal{N}(x; \mu_n, \sigma_n^2)$$

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

$$\pi_n = (\pi_n, \mu_n, \sigma_n^2)$$

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

$$\pi = (\pi_1, \mu_1, \sigma_1^2, \pi_2, \mu_2, \sigma_2^2, \dots, \pi_N, \mu_N, \sigma_N^2)$$

Si toutes les μ et σ^2 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^2 .

L'algorithme d'estimation s'appelle "Expectation-Maximization" ou "EM".

Ebauche de l'Algorithme EM

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

Fait une premier 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 $p(h_m | X_1, X_2, \dots, X_M, \quad ^{(i)})$

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

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

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

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

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 indépendants.

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 $\hat{\theta}$ 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 $L(\theta | S)$.

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

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

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

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

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

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 :

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

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

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

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

La maximum se trouve où le dérivé est nul.

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

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

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

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

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

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

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

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

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

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

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

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

Le gradient nous permet de calculer une solution analytique.
 Les estimations de $\hat{\boldsymbol{\mu}}, \hat{\mathbf{C}}$ sont obtenus par

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

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

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.

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

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

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

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

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

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 $Q^{(i)}$ définit par une esperence conditionnel

$$Q^{(i)} = E \{ l^{(i)} | \mathbf{T} \} = \text{Log} \left\{ \prod_{m=1}^M p(X_m | \theta) \right\}$$

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

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

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'etape M

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

$$Q = Q(\theta^{(i)}) - Q(\theta^{(i-1)})$$

Pour chaque cycle dans l'itération nous allons chercher :

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

Ce maximum est donné par :

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

$$\mu_n^{(i+1)} = \frac{1}{M} \sum_{m=1}^M P(h_{m=n} | \mathbf{X}_m, \theta^{(i)}) = \frac{1}{M} S_n^{(i+1)}$$

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

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

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

$$h(m, n) = P(h_m=n \mid X_m, \quad (i))$$

E (Expectation) :

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

M: (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 :

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