

EM アルゴリズムについて

前川眞一

東京工業大学大学院 社会理工学研究科 人間行動システム専攻

090708,14

1 EM アルゴリズム

パラメタ ξ と θ の事後分布の密度関数 $h(\xi, \theta | \mathbf{y})$ が与えられている場合に、そこから θ を積分消去した ξ の周辺事後分布

$$g(\xi | \mathbf{y}) = \int_{\Theta} h(\xi, \theta | \mathbf{y}) d\theta \quad (1)$$

の最頻値を求める問題を考える。

この問題は

$$\frac{\partial \log g(\xi | \mathbf{y})}{\partial \xi} = \frac{\partial \log \int_{\Theta} h(\xi, \theta | \mathbf{y}) d\theta}{\partial \xi} = 0 \quad (2)$$

の解を求めることと同値であるが、定積分が難しい場合には、その解を求めることは難しい。Dempster, et al(19xx) が提案した EM アルゴリズムは、上記の定積分をより簡単な形の定積分に置き換えることにより数値的に周辺最頻値を求める方法である。

いま、 $k(\theta | \xi, \mathbf{y})$ を ξ を与えたときの θ の事後分布の密度関数とすれば、 ξ の周辺事後分布の密度関数は、それぞれの密度関数の定義から

$$g(\xi | \mathbf{y}) = \frac{h(\xi, \theta | \mathbf{y})}{k(\theta | \xi, \mathbf{y})} \quad (3)$$

と書くことが出来る。したがって、その対数は

$$\log g(\xi | \mathbf{y}) = \log h(\xi, \theta | \mathbf{y}) - \log k(\theta | \xi, \mathbf{y}) \quad (4)$$

となる。

次に、パラメタ ξ の仮の値を $\xi^{(\ell)}$ として、それが与えられた時の θ の事後分布の密度関数を $k(\theta | \xi^{(\ell)}, \mathbf{y})$ とし、それに関して、(4) 式の両辺の期待値をとったものを考える。しかし、実は、(4) 式は θ には依存しない量であるため、左辺はそのままであることに注意すると、

$$\log g(\xi | \mathbf{y}) = Q(\xi | \xi^{(\ell)}) - H(\xi | \xi^{(\ell)}) \quad (5)$$

ただし、

$$Q(\xi | \xi^{(\ell)}) = \int_{\Theta} \log h(\xi, \theta | \mathbf{y}) \times k(\theta | \xi^{(\ell)}, \mathbf{y}) d\theta \quad (6)$$

$$H(\xi | \xi^{(\ell)}) = \int_{\Theta} \log k(\theta | \xi, \mathbf{y}) \times k(\theta | \xi^{(\ell)}, \mathbf{y}) d\theta \quad (7)$$

と書くことが出来る。

EM アルゴリズムとは、(1) 式を直接に ξ に関して最大化するのではなく、 $\xi^{(\ell)}$ が与えられた場合に、

$$Q(\xi^{(\ell+1)} | \xi^{(\ell)}) > Q(\xi^{(\ell)} | \xi^{(\ell)}) \quad (8)$$

を満たすような $\xi^{(\ell+1)}$ を逐次求めていきながら、最終的に (4) 式、したがって (1) 式 を最大にする ξ を求めていく方法である。 $\xi^{(\ell+1)}$ としては、例えば

$$\frac{\partial Q(\xi|\xi^{(\ell)})}{\partial \xi} = 0 \quad (9)$$

の解を用いることが考えられる。ただし、 $\xi^{(\ell)}$ は定数として取り扱う。

なお、 (8) 式 を満たすような $\xi^{(\ell+1)}$ は (5) 式 をも増加させる、すなわち

$$d^{(\ell)} = \log g(\xi^{(\ell+1)}|\mathbf{y}) - \log g(\xi^{(\ell)}|\mathbf{y}) > 0 \quad (10)$$

であることは、以下のようにして示すことが出来る。上式は

$$d^{(\ell)} = Q(\xi^{(\ell+1)}|\xi^{(\ell)}) - H(\xi^{(\ell+1)}|\xi^{(\ell)}) - \left(Q(\xi^{(\ell)}|\xi^{(\ell)}) - H(\xi^{(\ell)}|\xi^{(\ell)}) \right) \quad (11)$$

$$= Q(\xi^{(\ell+1)}|\xi^{(\ell)}) - Q(\xi^{(\ell)}|\xi^{(\ell)}) - \left(H(\xi^{(\ell+1)}|\xi^{(\ell)}) - H(\xi^{(\ell)}|\xi^{(\ell)}) \right) \quad (12)$$

とかけるが、その右辺の Q を含む部分は (8) 式 より正である。また、 H を含む部分は、対数関数が凹関数であることを利用してジェンセンの不等式を用いれば、

$$\begin{aligned} H(\xi^{(\ell+1)}|\xi^{(\ell)}) - H(\xi^{(\ell)}|\xi^{(\ell)}) &= \int_{\Theta} \log k(\theta|\xi^{(\ell+1)}, \mathbf{y}) \times k(\theta|\xi^{(\ell)}, \mathbf{y}) d\theta \\ &\quad - \int_{\Theta} \log k(\theta|\xi^{(\ell)}, \mathbf{y}) \times k(\theta|\xi^{(\ell)}, \mathbf{y}) d\theta \end{aligned} \quad (13)$$

$$= \int_{\Theta} \log \frac{k(\theta|\xi^{(\ell+1)}, \mathbf{y})}{k(\theta|\xi^{(\ell)}, \mathbf{y})} \times k(\theta|\xi^{(\ell)}, \mathbf{y}) d\theta \quad (14)$$

$$\leq \log \int_{\Theta} \frac{k(\theta|\xi^{(\ell+1)}, \mathbf{y})}{k(\theta|\xi^{(\ell)}, \mathbf{y})} \times k(\theta|\xi^{(\ell)}, \mathbf{y}) d\theta \quad (15)$$

$$= \log \int_{\Theta} k(\theta|\xi^{(\ell+1)}, \mathbf{y}) d\theta \quad (16)$$

$$= \log(1) = 0 \quad (17)$$

となり、負の値となる。したがって、 (10) 式 が示せた。なお、上式に於いて、最終行は、 k が密度関数であるため、その積分が 1 となることを用いた。

なお、上記の議論に於いて、同時事後分布を、パラメタが与えられたときのデータ \mathbf{y} の同時分布の密度関数と $f(\mathbf{y}|\xi, \theta)$ とパラメタの事前分布 $h(\theta, \xi)$ とで表現すれば、

$$h(\xi, \theta|\mathbf{y}) \propto f(\mathbf{y}|\xi, \theta)h(\theta, \xi) = f(\mathbf{y}|\xi, \theta)h(\theta|\xi)h(\xi) \quad (18)$$

となるが、これより

$$g(\xi|\mathbf{y}) \propto \int_{\Theta} f(\mathbf{y}|\xi, \theta)h(\theta|\xi)d\theta h(\xi) = L(\xi) h(\xi) \quad (19)$$

となる。したがって、パラメタ ξ に関する事前分布が一様分布の場合には、EM アルゴリズムは、所謂周辺尤度 $L(\xi)$ をを最大にするアルゴリズムと見なすことが出来る。実は、この形が本来の EM アルゴリズムであり、その意味で、 $Q(\xi|\xi^{(\ell)})$ は期待対数完全データ尤度と呼ばれることがある。

EM アルゴリズムを用いて周辺事後分布の最頻値が簡単に求まるか否かは、関数 Q 、すなわち、同時事後分布の密度関数の期待値が簡単に求められ、且つ、 Q の値を増加させる、もしくは最大にする

る $\xi^{(\ell+1)}$ が簡単に求められるか否かにかかっている。なお、この期待値を求める操作のことを E ステップと呼び、それを最大にする (もしくはその値を前回から増加させる) $\xi^{(\ell+1)}$ の値を求める操作のことを M ステップと呼んでいる。

EM アルゴリズムを用いて得られるパラメタ $\xi^{(\ell)}$ の系列は、以上に説明により、決して周辺事後分布の密度関数の値を減少させることがないため、最終的に得られる ξ の値は、周辺事後分布の最頻値であると考えられる。しかし、最頻値が複数存在する事後分布の場合、最も密度関数の値が大きい最頻値に到達するという保証はなく、最大値ではなく極大値に到達している可能性も否めない。また、EM アルゴリズムは、自動的に、周辺事後分布の最頻値付近における集中の度合い、例えば、事後標準偏差等を与えるものではない。しかし、周辺事後分布の要約として点推定値のみが必要な場合には、有用な方法である。

以下に、階層的線形モデルと潜在クラスモデルに関して、EM アルゴリズムを用いて周辺事後分布の最頻値を求める方法に関して記す。

註

シェンセンの不等式 Jensen's inequality

$f(x)$ を x の密度関数、 $r(x)$ を x の任意の実数値関数、 $q(x)$ を x の任意の凹関数とすると、

$$q\left(\int_X r(x)f(x)dx\right) \geq \int_X q(r(x))f(x)dx \quad (20)$$

すなわち

$$q(\mathcal{E}(r(\mathbf{x}))) \geq \mathcal{E}(q(r(\mathbf{x}))) \quad (21)$$

である。

凹関数 concave function

任意の変数の値 x と y に於いて、以下の関係が成り立つ関数を凸関数 convex function と呼ぶ。ただし a は $0 < a < 1$ を満たす定数である。

$$q(ax + (1-a)y) \leq aq(x) + (1-a)q(y) \quad (22)$$

$-q(x)$ が凸関数の場合 q は凹関数と呼ばれる。

2 階層的線形モデル

Q=3 層の階層線形モデル

$$\mathbf{y}|\mathbf{v}_1, \mathbf{C}_1 \sim N(\mathbf{G}_0\mathbf{v}_1, \mathbf{C}_1) \quad (23)$$

$$\mathbf{v}_1|\mathbf{v}_2, \mathbf{C}_2 \sim N(\mathbf{G}_1\mathbf{v}_2, \mathbf{C}_2) \quad (24)$$

$$\mathbf{v}_2 \sim N(\mathbf{G}_2\mathbf{v}_3, \mathbf{C}_3) \quad (25)$$

を考え、パラメタを

$$\xi = \{\mathbf{v}_2, \mathbf{C}_1, \mathbf{C}_2\}, \quad \theta = \mathbf{v}_1 \quad (26)$$

と分割する。また、 (v_3, C_3) 並びに C_1, C_2 の事前分布のハイパーパラメタ D_1, D_2 は定数とする。

いま、このモデルの ξ の周辺事後分布から最頻値を求めたいものとする。 ξ の周辺事後分布は、上記のモデルから θ を積分消去した周辺尤度に ξ の事前分布の密度関数を掛け合わせることで得られるが、周辺尤度が

$$y|v_2, C_1, C_2 \sim N(G_0 G_1 v_2, G_0 C_2 G_0' + C_1) \quad (27)$$

となり、 $|G_0 C_2 G_0' + C_1|$ や $(G_0 C_2 G_0' + C_1)^{-1}$ 等が出現するため、簡単に周辺事後分布の最頻値を求めることができない。

そこで EM アルゴリズムを用いて、周辺事後分布の最頻値を求めることを考える。まず、 $\xi^{(\ell)}$ を与えたときの θ の条件付き分布であるが、これは第 xx 章の (??) 式より、

$$v_1|y, C_1, C_2, v_2 \sim N(\theta_1, \Theta_1) \quad (28)$$

ただし、

$$\theta_1 = \Theta_1(C_2^{-1} G_1 v_2 + G_0' C_1^{-1} G_0 \widehat{v}_1) \quad (29)$$

$$\Theta_1 = (G_0' C_1^{-1} G_0 + C_2^{-1})^{-1} \quad (30)$$

$$\widehat{v}_1 = (G_0' C_1^{-1} G_0)^{-1} G_0' C_1^{-1} y \quad (31)$$

である。なお、上記の条件付き分布のパラメタ $\theta_1, \Theta_1, \widehat{v}_1$ の中に現れる C_1, C_2, v_2 は $\xi^{(\ell)}$ の値である。次に、同時事後分布の密度関数の対数は、 $h(\xi) = h_v(v_2)h_{C_1}(C_1)h_{C_2}(C_2)$ をパラメタの事前分布として

$$\begin{aligned} \log h(\xi, \theta|y) = & -\frac{1}{2} \left[\log |C_1| + (y - G_0 v_1)' C_1^{-1} (y - G_0 v_1) \right. \\ & \left. + \log |C_2| + (v_1 - G_1 v_2)' C_2^{-1} (v_1 - G_1 v_2) \right] + \log(h(\xi)) + const \end{aligned} \quad (32)$$

であるから、その v_1 に関する期待値は、付録の公式から、

$$\begin{aligned} Q(\xi, \xi^{(\ell)}) = & -\frac{1}{2} \left[\log |C_1| + (y - G_0 \theta_1)' C_1^{-1} (y - G_0 \theta_1) + \text{tr} \left(\Theta_1 G_0' C_1^{-1} G_0 \right) \right. \\ & \left. + \log |C_2| + (\theta_1 - G_1 v_2)' C_2^{-1} (\theta_1 - G_1 v_2) \right] + \text{tr} \left(\Theta_1 C_2^{-1} \right) + \log(h(\xi)) + const \end{aligned} \quad (33)$$

となるが、これが E ステップに対応する。したがって、上式を ξ に関して最大化する手続きが M ステップに当たるが、たとえば ξ のうち v_2 の事前分布が無情報分布である場合には、 v_2 に関しては、付録の残差の 2 乗和の分解公式から

$$v_2 = (G_1' C_2^{-1} G_1)^{-1} G_1' C_2^{-1} \theta_1 \quad (34)$$

がその解となることが分かる。また、 v_2 の事前分布が正規分布である場合には、一般ガウス・マルコフモデルの場合と同様に平方を完成を行えば、その事後平均に相当するものが Q を最大とする解であることが分かる。

また、(33) 式を以下のように変形すると

$$\begin{aligned} Q(\xi, \xi^{(\ell)}) = & -\frac{1}{2} \left[\log |C_1| + \text{tr} \left(C_1^{-1} ((y - G_0 \theta_1)(y - G_0 \theta_1)' + G_0 \Theta_1 G_0') \right) + \log(h_{C_1}(C_1)) \right] \\ & -\frac{1}{2} \left[\log |C_2| + \text{tr} \left(C_2^{-1} ((\theta_1 - G_1 v_2)(\theta_1 - G_1 v_2)' + \Theta_1) \right) + \log(h_{C_2}(C_2)) \right] \end{aligned} \quad (35)$$

となり、周辺尤度から導かれる周辺事後分布の密度関数よりも簡単な形をしていることが分かる。すなわち、 C_1 と C_2 は独立に推定することが可能である。特に、 Q を最大とする C_1 の構成要素である $\sigma_k, k = 1, 2, \dots, m$ は、準自然共役事前分布である逆カイ分布を用いれば、それぞれ独立に解析的に求めることが出来る。また、 C_2 の構成要素である Φ も同様である。

3 潜在クラスモデル

潜在クラスモデルにおいては、パラメタ ξ と ρ を ξ とし、各固体のクラスへの所属を示す潜在変数 $d_i, i = 1, 2, \dots, N$ を θ として取り扱う。したがって、前節で示した期待値をとるための積分を和で置き換える必要が生じる。

$\xi^{(\ell)}$ を与えたときの θ の事後分布は第 xx 章の (??) 式で示した多項分布となる。

また、(??) 式に示した、事後分布の密度関数の対数は

$$\log h(\xi, \theta | \mathbf{y}) = \sum_{i=1}^N \sum_{q=1}^Q d_{iq} \log (f(y_i | \xi_q) \rho_q) + \sum_{q=1}^Q (\alpha_q - 1) \log \rho_q + \log h(\xi) \quad (36)$$

$$= \sum_{q=1}^Q \sum_{i=1}^N d_{iq} (\log f(y_i | \xi_q) + \log \rho_q) + \sum_{q=1}^Q (\alpha_q - 1) \log \rho_q + \log h(\xi) \quad (37)$$

$$= \sum_{q=1}^Q \sum_{i=1}^N d_{iq} \log f(y_i | \xi_q) + \sum_{q=1}^Q \left(\sum_{i=1}^N d_{iq} \right) \log \rho_q + \sum_{q=1}^Q (\alpha_q - 1) \log \rho_q + \log h(\xi) \quad (38)$$

$$= \sum_{q=1}^Q \sum_{i=1}^N d_{iq} \log f(y_i | \xi_q) + \sum_{q=1}^Q \left(\sum_{i=1}^N d_{iq} + \alpha_q - 1 \right) \log \rho_q \quad (39)$$

となるが、 d の事後分布より、

$$\mathcal{E}(d_{iq} | y_i, \rho, \xi) = h_{iq} \quad (40)$$

であるから、 Q は

$$Q(\xi | \xi^{(\ell)}) = \sum_{q=1}^Q \sum_{i=1}^N h_{iq} \log f(y_i | \xi_q) + \log h(\xi) + \sum_{q=1}^Q \left(\sum_{i=1}^N h_{iq} + \alpha_q - 1 \right) \log \rho_q \quad (41)$$

となる。

したがって、M step では、 $\xi^{(\ell)}$ の関数である $h_{iq}, i = 1, 2, \dots, N, q = 1, 2, \dots, Q$ を定数と見なし、上式を最大化する事になるが、(41) 式の最初の項は、 y_i という値の観測値が h_{iq} 個観測された場合の尤度関数

$$f_h(\mathbf{y} | \xi_q) = \prod_{i=1}^N f(y_i | \xi_q)^{h_{iq}} \quad (42)$$

の積の対数であると考えられる。また、事前分布として、各クラスのパラメタの事前分布は独立である

$$h(\xi) = \prod_{q=1}^Q h(\xi_q) \quad (43)$$

と仮定すると、結局、(41)式の最初の2項は

$$h(\boldsymbol{\xi}_q | \mathbf{y}, \mathbf{h}) \propto \prod_{i=1}^N f(y_i | \boldsymbol{\xi}_q)^{h_{iq}} \times h(\boldsymbol{\xi}_q) \quad (44)$$

、すなわち、 y_i という値の観測値が h_{iq} 個観測された場合の事後分布の密度関数の対数となっていることが分かる。したがって、この事後分布の最頻値を与える $\boldsymbol{\xi}_q$ を求めれば、それが M step における Q を最大とするパラメタの値となる。

例えば、混合正規分布や switching regression の例では、 $\boldsymbol{\xi}_q$ の事前分布として事前共役事前分布をもちれば、この Q を最大とする $\boldsymbol{\xi}$ の値は解析的に求められるため、M step における計算は簡単である。

また、生起確率のパラメタ $\boldsymbol{\rho}$ に関しては、その和が 1 であるという制約を考えれば、

$$\rho_q = \frac{\sum_{i=1}^N h_{iq} + \alpha_q - 1}{\sum_{q=1}^Q \left(\sum_{i=1}^N h_{iq} + \alpha_q - 1 \right)} \quad (45)$$

が Q を最大とするパラメタの値となる。