

Gaussian Mixture Model

何浩勋

16212799

EM algorithm for the ML fitting of the parameter metric mixture model

The mixture model is expressed as,

$$f(y_j; \Psi) = \sum_{i=1}^g \pi_i f_i(y_j; \theta_i)$$

where $\theta_i = (\mu_i, \sigma_i)$,

Let $\pi = (\pi_1, \dots, \pi_g)^T$, $\mu = (\mu_1, \dots, \mu_g)^T$, $\sigma = (\sigma_1, \dots, \sigma_g)^T$, $\xi = (\mu^T, \sigma^T)^T$, $\Psi = (\pi^T, \xi^T)^T$, $y = (y_1^T, \dots, y_n^T)$.

The log likelihood for Ψ is given by,

$$\log L(\Psi) = \sum_{j=1}^n \log \left[\sum_{i=1}^g \pi_i f_i(y_j; \theta_i) \right]$$

Solving the likelihood equation,

$$\partial \log L(\Psi) / \partial \Psi = 0$$

We have $\hat{\Psi}$, satisfies

$$\hat{\pi}_i = \sum_{j=1}^n \tau_i(y_j; \hat{\Psi})/n \quad (i = 1, \dots, g)$$

and

$$\sum_{i=1}^g \sum_{j=1}^n \tau_i(y_j; \hat{\Psi}) \partial \log f_i(y_j; \hat{\theta}_i) / \partial \xi$$

where

$$\tau_i(y_j; \hat{\Psi}) = \pi_i f_i(y_j; \theta_i) / \sum_{h=1}^g \pi_h f_h(y_j; \theta_h)$$

If y is viewed as being incomplete, as the associated component vectors, $z = (z_1, \dots, z_n)$ are not available, z_j is a g -dimensional vector with $z_{ij} = (z_j)_i = 1$ or 0 , according to whether y_j is did or did not arise from the i th component of the mixture ($i = 1, \dots, g, j = 1, \dots, n$). The complete vector is therefore declared to be,

$$y_c = (y^T, z^T)$$

The component-label vector z_1, \dots, z_n are taken to be the realized values of the random vectors Z_1, \dots, Z_n , where, for independent feature data, it is appropriate to assume that they are distributed unconditionally as,

$$Z_1, \dots, Z_n \sim \text{Mult}_g(1, \pi)$$

This assumption means that the distribution of the complete-data vector Y_C implies the appropriate distribution of the incomplete-data vector Y . The complete-data log likelihood for $L_c(\Psi)$, is given by

$$\log L_c(\Psi) = \sum_{i=1}^g \sum_{j=1}^n z_{ij} \{\log \pi_i + \log f_i(y_j; \theta_i)\}$$

E-step

$$w_{ij} = E_{\Psi^{(k)}}(Z_{ij}|y) = \pi_i^{(k)} f_i(y_j; \theta_i^{(k)}) / \sum_{h=1}^g \pi_h f_h(y_j; \theta_h^{(k)})$$

M-step

$$\pi_i^{(k+1)} = \sum_{j=1}^n w_{ij} / n \quad (i = 1, \dots, g)$$

$$\mu_i^{(k+1)} = \sum_{j=1}^n w_{ij} y_j / \sum_{j=1}^n w_{ij} \quad (i = 1, \dots, g)$$

$$\sigma_i^{(k+1)} = \sqrt{\sum_{j=1}^n w_{ij} (y_j - \mu_i^{(k+1)})^2 / \sum_{j=1}^n w_{ij}} \quad (i = 1, \dots, g)$$

Random starting values

Specifying a random start,

$$\mu_1^{(0)}, \dots, \mu_g^{(0)} \sim N(\bar{y}, V)$$

$$\Sigma_0^{(0)} = V$$

$$\pi_1 = \dots = \pi_g = 1/g$$

where $V = \sum_{j=1}^n (y_j - \bar{y})(y_j - \bar{y})^T / n$

Stopping criterion

$$l^{(k)} = \log L_c(\Psi^{(k)}) = \sum_{i=1}^g \sum_{j=1}^n z_{ij} \{\log \pi_i + \log f_i(y_j; \theta_i)\}$$

$$a^{(k)} = (l^{(k+1)} - l^{(k)}) / (l^{(k)} - l^{(k-1)}) \quad , k > 1$$

$$l_A^{(k+1)} = l^{(k)} + \frac{1}{1-a^{(k)}} (l^{(k+1)} - l^{(k)})$$

The EM algorithm can be stopped if

$$|l_A^{(k+1)} - l_A^{(k)}| < tol$$

Simulation

Example 1

I randomly generated Y_1, Y_2, Y_3 ,

$$Y_1 \sim N(-2, 1)$$

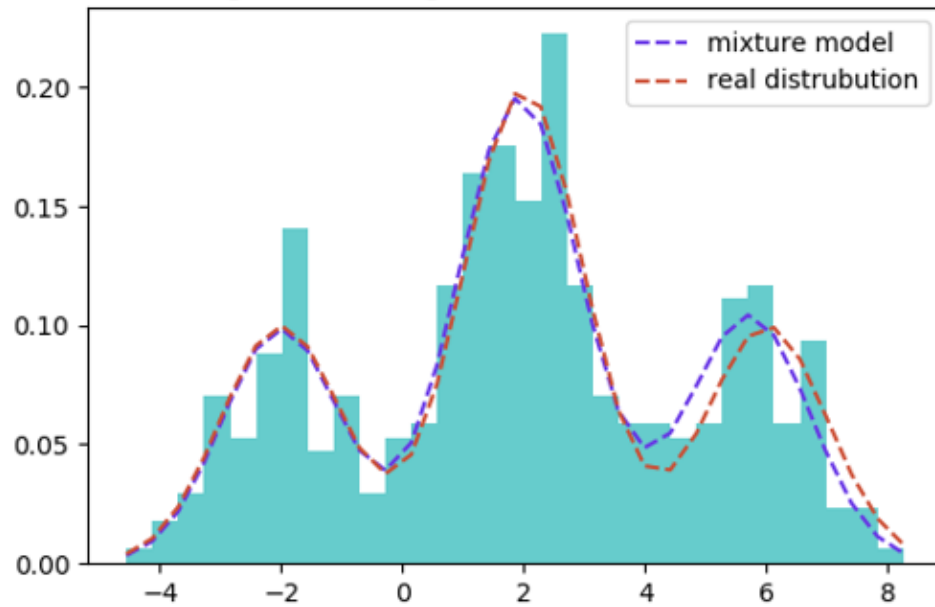
$$Y_2 \sim N(2, 1)$$

$$Y_3 \sim N(6, 1)$$

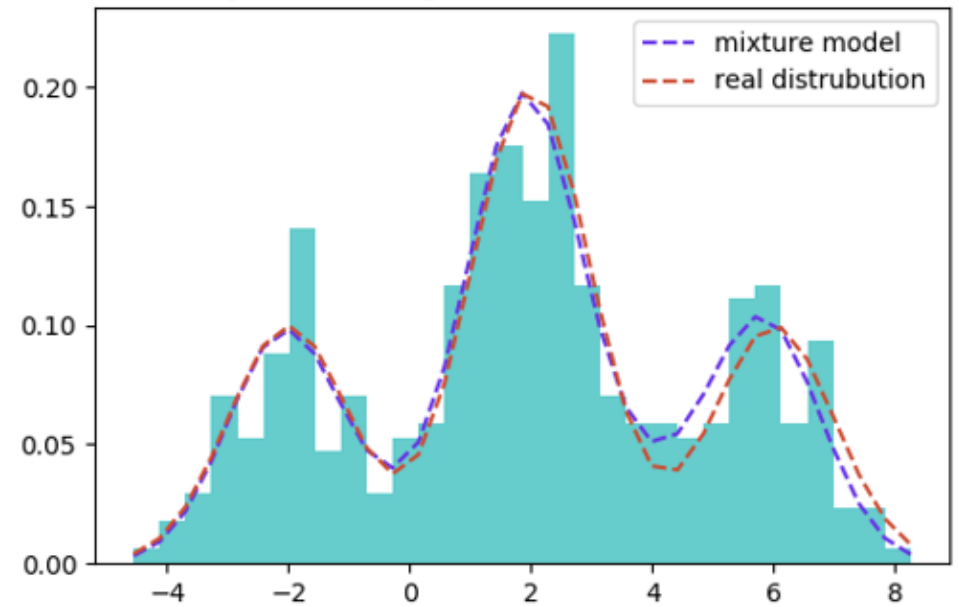
Y_1, Y_3 each have 100 observations, Y_2 have 200 observations, that is, 400 observations total. $Y = (Y_1^T, Y_2^T, Y_3^T)$

The mixture model fits the real distribution well when $g \leq 4$, we can see that there are signs of overfit when $g > 4$.

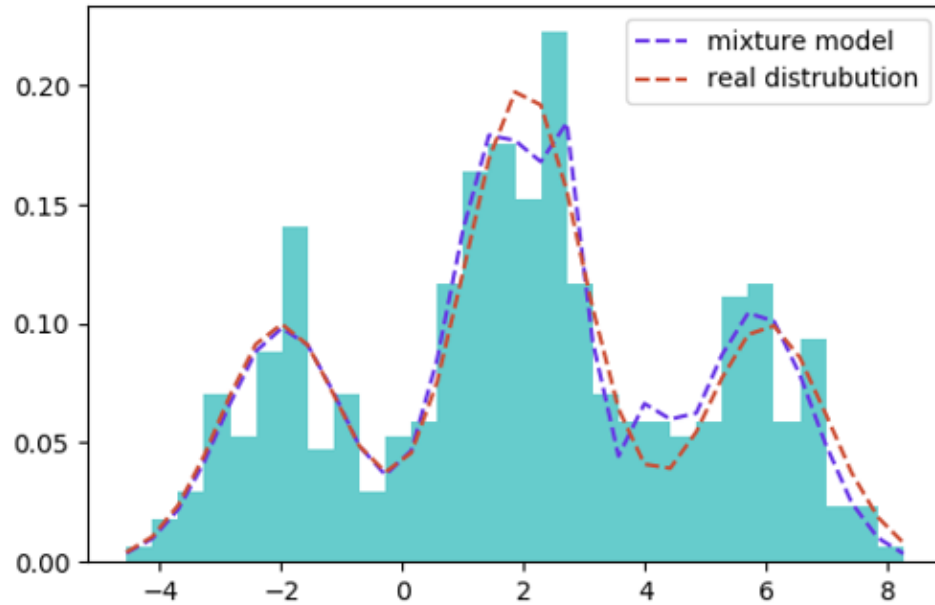
g=3 , convergence after 85 iterations



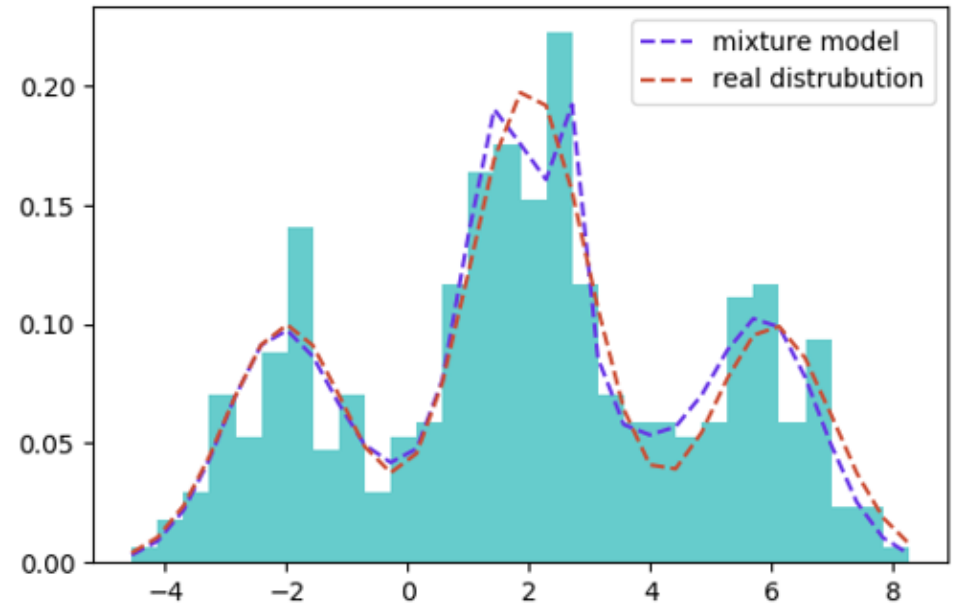
g=4 , convergence after 108 iterations



g=5 , convergence after 1749 iterations



g=6 , convergence after 1454 iterations



Example 2

I randomly generated Y_1, Y_2 ,

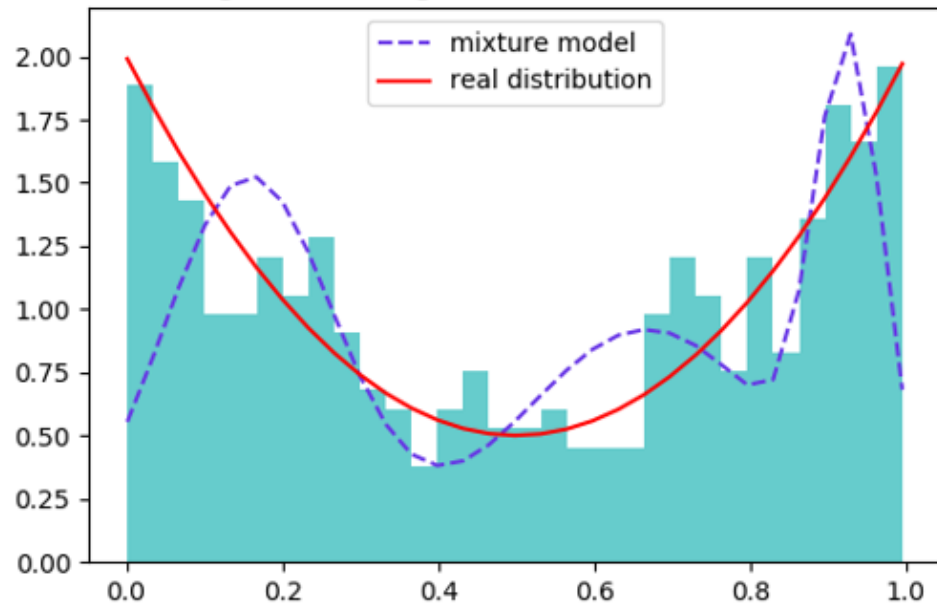
$$Y_1 \sim \text{Beta}(1, 4)$$

$$Y_2 \sim \text{Beta}(4, 1)$$

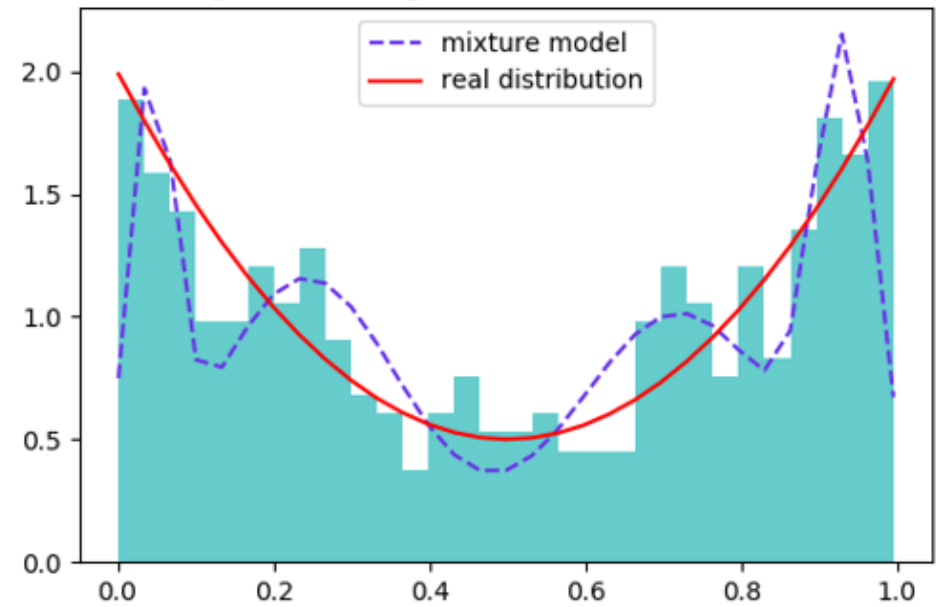
Y_1, Y_2 each have 200 observations, that is, 400 observations total. $Y = (Y_1^T, Y_2^T)$

When $g > 4$, I encountered a precision problem which causes 'Log(0) error', I had to stop EM algorithm and use the current result. However, the parameters seems converged.

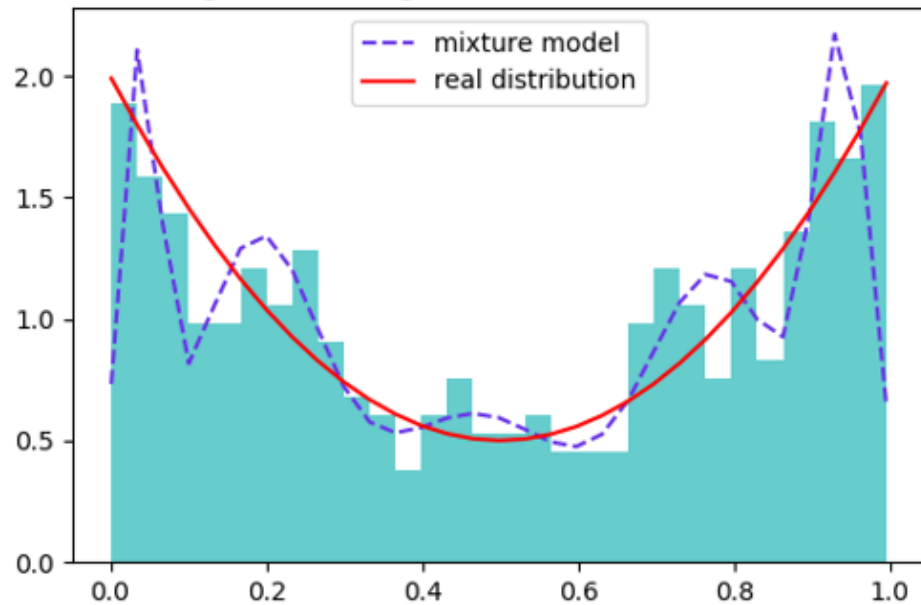
g=3 , convergence after 135 iterations



g=4 , convergence after 236 iterations



g=5 , convergence after 198 iterations



g=6 , convergence after 61 iterations

