

# EM Theory: (Statistician's View)

Data:  $\{X_i\}_{i=1}^n$  is the sample data points:  $\{Y_i\}_{i=1}^n$  is the corresponding labels of which Gaussian  $X_i$  is from;  
 $(X_i \in \mathbb{R}^m)$  (observed)  $(Y_i \in \mathbb{R}^K)$  (missing)

Model:  $Y_i \stackrel{iid}{\sim} \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K)$ ; (i.e.  $Y_i = (Y_i^1, \dots, Y_i^K)$ , exact one  $Y_i^k = 1$ , all others are 0,  
 $p(Y_i = (0, 0, \dots, 1, \dots, 0)) = p(Y_i^k = 1) = \pi_k$ )  
 $X_i | Y_i^k \sim N(\mu_k, \Sigma_k)$ ; (i.e. pdf.  $f(x_i | Y_i^k = 1) = N(x_i, \mu_k, \Sigma_k) = (2\pi)^{-\frac{m}{2}} |\Sigma_k|^{-\frac{1}{2}} e^{-\frac{1}{2}(x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)}$ )

Inference: parameters  $\theta = (\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K)$ ;  $X = \{X_i\}_{i=1}^n$ ;  $Y = \{Y_i\}_{i=1}^n$ ;

Bayes Rule:  $p(X_i, Y_i | \theta) = p(X_i | \theta) p(Y_i | X_i, \theta)$ , so  $p(X_i | \theta) = \frac{p(X_i, Y_i | \theta)}{p(Y_i | X_i, \theta)}$ ,  $\ln p(X_i | \theta) = \ln p(X_i, Y_i | \theta) - \ln p(Y_i | X_i, \theta)$ ;

$\Rightarrow$  log-likelihood:  $\ell(\theta | X) = \sum_{i=1}^n \ln p(X_i | \theta) = \sum_{i=1}^n [\ln p(X_i, Y_i | \theta) - \ln p(Y_i | X_i, \theta)]$ ;

Calculate expectation w.r.t.  $Y_i | X_i, \theta_e$ , i.e.  $\int \dots p(Y_i | X_i, \theta_e) dY_i$ , we have:

$$E[\ln p(X_i | \theta) | X_i, \theta_e] = \int \ln p(X_i | \theta) p(Y_i | X_i, \theta_e) dY_i = \ln p(X_i | \theta)$$

$$E[\ln p(X_i, Y_i | \theta) | X_i, \theta_e] = Q_i(\theta | \theta_e, X_i), \quad E[\ln p(Y_i | X_i, \theta) | X_i, \theta_e] = H_i(\theta | \theta_e, X_i);$$

$$\begin{aligned} \Rightarrow \ell(\theta | X) &= \sum_{i=1}^n \ln p(X_i | \theta) = \sum_{i=1}^n E[\ln p(X_i | \theta) | X_i, \theta_e] = \sum_{i=1}^n (E[\ln p(X_i, Y_i | \theta) | X_i, \theta_e] - E[\ln p(Y_i | X_i, \theta) | X_i, \theta_e]) \\ &= \sum_{i=1}^n Q_i(\theta | \theta_e, X_i) - \sum_{i=1}^n H_i(\theta | \theta_e, X_i); \end{aligned}$$

For  $H_i$ :  $H_i(\theta | \theta_e, X_i) - H_i(\theta_e | \theta_e, X_i) = E[\ln \frac{p(Y_i | X_i, \theta)}{p(Y_i | X_i, \theta_e)} | X_i, \theta_e]$ , by Jensen's inequality:  $E(\ln X) \leq \ln EX$

$$\begin{aligned} &\leq \ln E[\frac{p(Y_i | X_i, \theta)}{p(Y_i | X_i, \theta_e)} | X_i, \theta_e] \\ &= \ln \int \frac{p(Y_i | X_i, \theta)}{p(Y_i | X_i, \theta_e)} \cdot p(Y_i | X_i, \theta_e) dY_i = \ln \int p(Y_i | X_i, \theta_e) dY_i = 0; \end{aligned}$$

$$\Rightarrow \forall \theta, H_i(\theta | \theta_e, X_i) \leq H_i(\theta_e | \theta_e, X_i), \text{ so } \sum_{i=1}^n H_i(\theta | \theta_e, X_i) \leq \sum_{i=1}^n H_i(\theta_e | \theta_e, X_i);$$

If we set  $\theta_{e+1} = \arg \max_{\theta} \sum_{i=1}^n Q_i(\theta | \theta_e, X_i)$ , then:  $\begin{cases} \sum_{i=1}^n Q_i(\theta_{e+1} | \theta_e, X_i) \geq \sum_{i=1}^n Q_i(\theta_e | \theta_e, X_i) \\ \sum_{i=1}^n H_i(\theta_{e+1} | \theta_e, X_i) \leq \sum_{i=1}^n H_i(\theta_e | \theta_e, X_i) \end{cases}$

$$\Rightarrow \ell(\theta_{e+1} | X) \geq \ell(\theta_e | X);$$

Therefore if we iteratively update  $\theta_e$  as above, then  $\ell(\theta_e | X)$  will be nondecreasing,

so in this way we can obtain the (local) maximum log-likelihood;

Now to update

$$\theta_{t+1} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n Q_i(\theta | \theta_t, X_i) \quad , \quad Q_i(\theta | \theta_t, X_i) = \int \ln p(X_i, Y_i | \theta) p(Y_i | X_i, \theta_t) dY_i ;$$

denote  $Z_i = k$  if  $y_i^k = 1$ ,

which is the distribution  $X_i$  is from;

denote  $\alpha_{ik} = p(Z_i = k | X_i, \theta_t)$ ;

$$\begin{aligned} &= \int [\ln p(X_i | Y_i, \theta) + \ln p(Y_i | \theta)] p(Y_i | X_i, \theta_t) dY_i ; \\ &= \sum_{k=1}^K p(Z_i = k | X_i, \theta_t) [\ln N(X_i, \mu_k, \Sigma_k) + \ln \pi_k] ; \quad (Y_i \text{ discrete}) \\ &= \sum_{k=1}^K \alpha_{ik} \left[ -\frac{n}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (X_i - \mu_k)^T \Sigma_k^{-1} (X_i - \mu_k) + \ln \pi_k \right] ; \end{aligned}$$

$$\begin{aligned} \text{Here: } \alpha_{ik} &= p(Z_i = k | X_i, \theta_t) = \frac{p(X_i | Z_i = k, \theta_t) p(Z_i = k | \theta_t)}{\sum_{r=1}^K p(X_i | Z_i = r, \theta_t) p(Z_i = r | \theta_t)} \\ &= \frac{\pi_{kt} \cdot N(X_i, \mu_{kt}, \Sigma_{kt})}{\sum_{r=1}^K \pi_{rt} N(X_i, \mu_{rt}, \Sigma_{rt})} \end{aligned}$$

$$(\text{Bayes'ian Rule: } p(Z_i = k | X_i, \theta_t) = \frac{p(X_i | Z_i = k, \theta_t) p(Z_i = k | \theta_t)}{p(X_i | \theta_t)})$$

$\theta_t = (\pi_{1t} \dots \pi_{Kt}, \mu_{1t} \dots \mu_{Kt}, \Sigma_{1t} \dots \Sigma_{Kt})$ ,  
the previous  $\theta$  values, known;

So  $\alpha_{ik}$  is known;

This calculation is the E-step for all  $i = 1, 2, \dots, n$ ;

$$\Rightarrow \theta_{t+1} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \sum_{k=1}^K \alpha_{ik} \left[ \ln \pi_k - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (X_i - \mu_k)^T \Sigma_k^{-1} (X_i - \mu_k) \right] \quad , \quad \alpha_{ik} \cdot -\frac{n}{2} \ln 2\pi \text{ is constant, removed ;}$$

$$\left( \sum_{k=1}^K \alpha_{ik} = \sum_{k=1}^K p(Z_i = k | X_i, \theta_t) = 1 \right)$$

Usual way: let derivative to be 0;

$$\textcircled{1} \pi_{kt+1}: \sum_{k=1}^K \pi_k = 1, \text{ so } Q(\theta) = Q(\theta) - \lambda \left( \sum_{k=1}^K \pi_k - 1 \right), \quad \frac{\partial Q}{\partial \pi_k} = \frac{1}{\pi_k} \sum_{i=1}^n \alpha_{ik} - \lambda, \quad \frac{1}{\pi_k} = \frac{1}{\lambda} \sum_{i=1}^n \alpha_{ik}; \quad \sum_{k=1}^K \frac{1}{\pi_k} = 1 \text{ gives } \lambda = \frac{1}{\sum_{i=1}^n \sum_{k=1}^K \alpha_{ik}} = \frac{1}{n};$$

$$\Rightarrow \pi_{kt+1} = \hat{\pi}_k = \frac{1}{n} \sum_{i=1}^n \alpha_{ik};$$

$$\textcircled{2} \mu_{kt+1}: \frac{\partial Q}{\partial \mu_k} = \sum_{i=1}^n \alpha_{ik} \Sigma_k^{-1} (X_i - \mu_k) = \Sigma_k^{-1} \sum_{i=1}^n \alpha_{ik} (X_i - \mu_k), \Rightarrow \mu_{kt+1} = \frac{\sum_{i=1}^n \alpha_{ik} X_i}{\sum_{i=1}^n \alpha_{ik}};$$

$$\begin{aligned} \textcircled{3} \Sigma_{kt+1}: \frac{\partial Q}{\partial \Sigma_k} &= \sum_{i=1}^n \alpha_{ik} \left[ -\frac{1}{2} \Sigma_k^{-1} + \frac{1}{2} \Sigma_k^{-1} (X_i - \mu_k)(X_i - \mu_k)^T \Sigma_k^{-1} \right] \Rightarrow \Sigma_{kt+1} = \frac{\sum_{i=1}^n \alpha_{ik} (X_i - \mu_k)(X_i - \mu_k)^T}{\sum_{i=1}^n \alpha_{ik}} ; \\ &= \frac{1}{2} \Sigma_k^{-1} \sum_{i=1}^n \alpha_{ik} [(X_i - \mu_k)(X_i - \mu_k)^T \Sigma_k^{-1} - I] \end{aligned}$$

(Matrix differentiation can be searched online, which is not required)

These update-step is the M-step for all  $k = 1, 2, \dots, K$ ;