

Homework 4 STOR 767: Theory Part

Due November 12, 2020

Student Name: Leo Li (PID: 730031954)

Instructions

- Edit this L^AT_EX file with your solutions and generate a PDF file from it. Upload **both the tex and the pdf** file to Sakai.
- Use proper fonts for a clear presentation:
 x for an observed value; X for a random variable; \mathbf{X} for a vector; $\mathbf{\Sigma}$ for a matrix.
- You are allowed to work with other students but homework should be in your own words. Identical solutions will receive a **0** in grade and will be investigated.

1. Consider the Gaussian mixture model for clustering with latent $Z \sim \text{Multinomial}(1, (\eta_1, \dots, \eta_K))$ and $\mathbf{X}|Z = j \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j), \forall j = 1, \dots, K$. Develop the E-step and M-step for the parameters $\eta_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j$'s.

SOLUTION.

In this question, we would like to derive the EM algorithm to obtain the parameters in the Gaussian mixture model for clustering. We start by letting $Y = (\mathbf{X}_1, Z_1), \dots, (\mathbf{X}_n, Z_n)$ be the full data, and $Y_{obs} = \mathbf{X}_1, \dots, \mathbf{X}_n$ be the observed data. The full data likelihood function is derived as,

$$L_n = \prod_{i=1}^n \prod_{j=1}^K \eta_j^{I(Z_i=j)} P_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{X}_i)^{I(Z_i=j)}.$$

Then the log-likelihood function is derived as,

$$l_n = \sum_{i=1}^n \sum_{j=1}^K I(Z_i = j) (\log(\eta_j) + \log(P_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{X}_i))).$$

We also initialize $\psi^{(0)} = (\eta_1^{(0)}, \dots, \eta_K^{(0)}, \boldsymbol{\mu}_1^{(0)}, \dots, \boldsymbol{\mu}_K^{(0)}, \boldsymbol{\Sigma}_1^{(0)}, \dots, \boldsymbol{\Sigma}_K^{(0)})$, and for $t = 0, 1, 2, \dots$, the E-step and M-step are developed as follow:

- E-step:

We calculate the conditional expectation of the full data log-likelihood function given the observed data and $\psi^{(t)}$:

$$\begin{aligned} E[l_n(Y; \psi) | Y_{obs}, \psi^{(t)}] &= E\left[\sum_{i=1}^n \sum_{j=1}^K I(Z_i = j)(\log(\eta_j) + \log(P_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{X}_i))) | \{\mathbf{X}_i\}_1^n, \psi^{(t)}\right] \\ &= \sum_{i=1}^n \sum_{j=1}^K \gamma_{ij}^{(t)} (\log(\eta_j) + \log(P_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{X}_i))), \end{aligned}$$

where,

$$\begin{aligned} \gamma_{ij}^{(t)} &= E[I(Z_i = j) | \{\mathbf{X}_i\}_1^n, \psi^{(t)}] \\ &= E[I(Z_i = j) | \mathbf{X}_i, \psi^{(t)}] \\ &= P(Z_i = j | \mathbf{X}_i, \psi^{(t)}) \\ &= \frac{P(\mathbf{X}_i | Z_i = j, \psi^{(t)}) P(Z_i = j | \psi^{(t)})}{P(\mathbf{X}_i | \psi^{(t)})} \\ &= \frac{\eta_j^{(t)} P_{\boldsymbol{\mu}_j^{(t)}, \boldsymbol{\Sigma}_j^{(t)}}(\mathbf{X}_i)}{\sum_{j=1}^K \eta_j^{(t)} P_{\boldsymbol{\mu}_j^{(t)}, \boldsymbol{\Sigma}_j^{(t)}}(\mathbf{X}_i)}. \end{aligned}$$

- M-step

At the M-step, we would like to maximize the conditional expectation that we obtained in the E-step to find new estimates for the parameters. Specifically, we would like to obtain the updated estimate $\psi^{(t+1)}$ such that,

$$\begin{aligned} \psi^{(t+1)} &= \underset{\psi}{\operatorname{argmax}} E\left[\sum_{i=1}^n \sum_{j=1}^K I(Z_i = j)(\log(\eta_j) + \log(P_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{X}_i))) | \{\mathbf{X}_i\}_1^n, \psi^{(t)}\right] \\ &= \underset{\psi}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \sum_{j=1}^K \gamma_{ij}^{(t)} (\log(\eta_j) + \log(P_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{X}_i))) \right\}. \end{aligned}$$

Then, we can find the estimate $\boldsymbol{\mu}_j^{(t+1)}$ by differentiating the objective function and setting the derivative to zero,

$$\frac{\partial}{\partial \boldsymbol{\mu}_j} \left\{ \sum_{i=1}^n \sum_{j=1}^K \gamma_{ij}^{(t)} (\log(\eta_j) + \log(P_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{X}_i))) \right\} = \sum_{i=1}^n \gamma_{ij}^{(t)} \boldsymbol{\Sigma}_j^{-1} (\mathbf{X}_i - \boldsymbol{\mu}_j) = \mathbf{0}.$$

By solving the above equation, we have that,

$$\boldsymbol{\mu}_j^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{ij}^{(t)} \mathbf{X}_i}{\sum_{i=1}^n \gamma_{ij}^{(t)}}.$$

Similarly, we can find the estimate $\boldsymbol{\Sigma}_j^{(t+1)}$ by using the same technique,

$$\begin{aligned} & \frac{\partial}{\partial \boldsymbol{\Sigma}_j} \left\{ \sum_{i=1}^n \sum_{j=1}^K \gamma_{ij}^{(t)} (\log(\eta_j) + \log(P_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{X}_i))) \right\} \\ &= \frac{1}{2} \sum_{i=1}^n \gamma_{ij}^{(t)} \boldsymbol{\Sigma}_j - \frac{1}{2} \sum_{i=1}^n \gamma_{ij}^{(t)} (\mathbf{X}_i - \boldsymbol{\mu}_j)(\mathbf{X}_i - \boldsymbol{\mu}_j)^T \\ &= \mathbf{0}. \end{aligned}$$

By solving the above equations and using the invariance property, we have that,

$$\boldsymbol{\Sigma}_j^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{ij}^{(t)} (\mathbf{X}_i - \boldsymbol{\mu}_j^{(t+1)})(\mathbf{X}_i - \boldsymbol{\mu}_j^{(t+1)})^T}{\sum_{i=1}^n \gamma_{ij}^{(t)}}.$$

Finally, to find the estimate $\eta_j^{(t+1)}$, we would need to maximize the objective function with respect to η_j subject to the constraint such that,

$$\sum_{j=1}^K \eta_j^{(t+1)} = 1$$

By using the Lagrange Multiplier, we have that,

$$\frac{\partial}{\partial \eta_j} \left\{ \sum_{i=1}^n \sum_{j=1}^K \gamma_{ij}^{(t)} (\log(\eta_j) + \log(P_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{X}_i))) \right\} = \sum_{i=1}^n \frac{\gamma_{ij}^{(t)}}{\eta_j} = \lambda,$$

so that,

$$\begin{aligned} \eta_j &= \sum_{i=1}^n \frac{\gamma_{ij}^{(t)}}{\lambda}, \\ \sum_{j=1}^K \eta_j^{(t+1)} &= \frac{\sum_{i=1}^n \sum_{j=1}^K \gamma_{ij}^{(t)}}{\lambda} = \frac{n}{\lambda} = 1, \\ \lambda &= n. \end{aligned}$$

Therefore, the estimate $\eta_j^{(t+1)}$ is obtained by solving

$$\sum_{i=1}^n \frac{\gamma_{ij}^{(t)}}{\eta_j} = n,$$

which leads to,

$$\eta_j^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{ij}^{(t)}.$$

To sum up, at the M-step, we have that,

$$\eta_j^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{ij}^{(t)},$$

$$\boldsymbol{\mu}_j^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{ij}^{(t)} \mathbf{X}_i}{\sum_{i=1}^n \gamma_{ij}^{(t)}},$$

$$\boldsymbol{\Sigma}_j^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{ij}^{(t)} (\mathbf{X}_i - \boldsymbol{\mu}_j^{(t+1)}) (\mathbf{X}_i - \boldsymbol{\mu}_j^{(t+1)})^T}{\sum_{i=1}^n \gamma_{ij}^{(t)}}.$$

□

2. Consider the Gaussian graphical model $\mathbf{X} \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with precision matrix $\Theta = \boldsymbol{\Sigma}^{-1}$. Show that $\Theta_{jk} = 0$ if and only if $X_j \perp\!\!\!\perp X_k | \{X_1, \dots, X_d\} \setminus \{X_j, X_k\}$.

PROOF.

For the Gaussian graphical model $\mathbf{X} \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, if we partition the random vector \mathbf{X} into two components, $\mathbf{X}_A \in \mathbb{R}^a$ and $\mathbf{X}_B \in \mathbb{R}^b$ such that $a + b = d$, and let $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ be partitioned accordingly, i.e.,

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{bmatrix},$$

and

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{A,A} & \boldsymbol{\Sigma}_{A,B} \\ \boldsymbol{\Sigma}_{B,A} & \boldsymbol{\Sigma}_{B,B} \end{bmatrix},$$

then,

- The marginal distribution of \mathbf{X}_A is $N(\boldsymbol{\mu}_A, \boldsymbol{\Sigma}_{A,A})$.
- The conditional distribution of $\mathbf{X}_A | \mathbf{X}_B = \mathbf{x}_B$ is $N(\boldsymbol{\mu}_{A|B}, \boldsymbol{\Sigma}_{A|B})$, where

$$\boldsymbol{\mu}_{A|B} = \boldsymbol{\mu}_A + \boldsymbol{\Sigma}_{A,B} \boldsymbol{\Sigma}_{B,B}^{-1} (\mathbf{x}_B - \boldsymbol{\mu}_B),$$

$$\boldsymbol{\Sigma}_{A|B} = \boldsymbol{\Sigma}_{A,A} - \boldsymbol{\Sigma}_{A,B} \boldsymbol{\Sigma}_{B,B}^{-1} \boldsymbol{\Sigma}_{B,A}.$$

Now, we can start by proving the following claim:

- **Claim:** Under the setting of $\mathbf{X} \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, for $j, k \in \{1, \dots, d\}$ with $j \neq k$, we have that, $X_j \perp\!\!\!\perp X_k$ if and only if $\sigma_{jk} = 0$, where X_j and X_k are j^{th} and k^{th} components of \mathbf{X} , and σ_{jk} is the $(j, k)^{th}$ element of $\boldsymbol{\Sigma}$.

- **Proof:** We know that the marginal distribution of $(X_j, X_k)^T$ is $N((\mu_j, \mu_k)^T, \boldsymbol{\Sigma}_{\{j,k\}})$, where,

$$\boldsymbol{\Sigma}_{\{j,k\}} = \begin{bmatrix} \sigma_{jj} & \sigma_{jk} \\ \sigma_{kj} & \sigma_{kk} \end{bmatrix}.$$

Then the marginal distribution of X_j is $N(\mu_j, \sigma_{jj})$, and the conditional distribution of $X_j|X_k = x_k$ is,

$$N(\mu_j + \sigma_{jk}\sigma_{kk}^{-1}(x_k - \mu_k), \sigma_{jj} - \sigma_{jk}\sigma_{kk}^{-1}\sigma_{kj}).$$

Then, $X_j \perp\!\!\!\perp X_k$ if and only if the marginal distribution of X_j , $f_{X_j}(x_j)$, equals to the conditional distribution of $X_j|X_k = x_k$, $f_{X_j|X_k=x_k}(x_j)$ (by the definition of independence), which happens if and only if $\sigma_{jk} = 0$ (both directions are very obvious), as desired.

Once we have proved the above claim, we move onto the original statement that we would like to prove. Then we have that,

- $X_j \perp\!\!\!\perp X_k | \{X_1, \dots, X_d\} \setminus \{X_j, X_k\}$ if and only if the conditional covariance matrix, $\Sigma_{\{j,k\}||[d] \setminus \{j,k\}}$, is diagonal (by the claim that we have proved earlier).
- $\Sigma_{\{j,k\}||[d] \setminus \{j,k\}}$ is diagonal if and only if its inverse, $\Sigma_{\{j,k\}||[d] \setminus \{j,k\}}^{-1}$, is diagonal (because the inverse of a non-singular diagonal matrix is also diagonal).
- $\Sigma_{\{j,k\}||[d] \setminus \{j,k\}}^{-1} = \Theta_{\{j,k\}}$, where

$$\Theta_{\{j,k\}} = \begin{bmatrix} \Theta_{jj} & \Theta_{jk} \\ \Theta_{kj} & \Theta_{kk} \end{bmatrix}$$

(this can be verified by deriving the inverse of block matrix).

- $\Theta_{\{j,k\}}$ is diagonal matrix is equivalent to $\Theta_{jk} = 0$.

Therefore, $\Theta_{jk} = 0$ if and only if $X_j \perp\!\!\!\perp X_k | \{X_1, \dots, X_d\} \setminus \{X_j, X_k\}$, as desired. \square