

GMMs by Mean Field Variational Inference

Jiajun He

In Bayesian framework, we want to find a set of parameters

$$\boldsymbol{\theta} = \{\boldsymbol{\pi}, z_1, z_2, \dots, z_N, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K\}$$

, which maximizes the posterior probability $P(\boldsymbol{\theta}|\mathbf{X})$. But the form of the posterior distribution is relatively complicated in GMMs, so we use $q(\boldsymbol{\theta})$ to approximate it.

According to Bayes Rule, the probability of observing $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ is

$$P(\mathbf{X}) = \frac{P(\boldsymbol{\theta}, \mathbf{X})}{P(\boldsymbol{\theta}|\mathbf{X})} = \frac{P(\boldsymbol{\theta}, \mathbf{X})}{q(\boldsymbol{\theta})} \frac{q(\boldsymbol{\theta})}{P(\boldsymbol{\theta}|\mathbf{X})}$$

Take logarithm and expectation on both sides, we get

$$\int q(\boldsymbol{\theta}) \log P(\mathbf{X}) d\boldsymbol{\theta} = \int q(\boldsymbol{\theta}) \log \frac{P(\boldsymbol{\theta}, \mathbf{X})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta} + \int q(\boldsymbol{\theta}) \log \frac{q(\boldsymbol{\theta})}{P(\boldsymbol{\theta}|\mathbf{X})} d\boldsymbol{\theta}$$

i.e.,

$$\log P(\mathbf{X}) = \int q(\boldsymbol{\theta}) \log \frac{P(\boldsymbol{\theta}, \mathbf{X})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta} + D_{KL}[q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathbf{X})]$$

The term on the left-hand side is a constant, so minimizing KL-divergence between $q(\boldsymbol{\theta})$ and real posterior distribution is equivalent to maximizing the first term on the right-hand side, which is named as *ELBO*.

Under mean field assumption, i.e., $q(\boldsymbol{\theta}) = \prod_{\theta_i \in \boldsymbol{\theta}} q(\theta_i)$, each variable can be optimized successively:

Specifically, for each $\theta_i \in \boldsymbol{\theta}$,

$$\begin{aligned} ELBO &= \int q(\boldsymbol{\theta}) \log \left(\frac{P(\boldsymbol{\theta}, \mathbf{X})}{q(\boldsymbol{\theta})} \right) d\boldsymbol{\theta} \\ &= \int q(\boldsymbol{\theta}) \log P(\boldsymbol{\theta}, \mathbf{X}) d\boldsymbol{\theta} - \int q(\boldsymbol{\theta}) \log q(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \mathbb{E}_{\theta_i} [\mathbb{E}_{\boldsymbol{\theta}-\{\theta_i\}} [\log P(\boldsymbol{\theta}, \mathbf{X})]] - \mathbb{E}_{\boldsymbol{\theta}} [\log q(\theta_i)] - \sum_{j \neq i} \mathbb{E}_{\boldsymbol{\theta}} [\log q(\theta_j)] \\ &= \int q(\theta_i) \mathbb{E}_{\boldsymbol{\theta}-\{\theta_i\}} [\log P(\boldsymbol{\theta}, \mathbf{X})] d\theta_i - \int q(\theta_i) \log q(\theta_i) d\theta_i - \sum_{j \neq i} \int q(\theta_j) \log q(\theta_j) d\theta_j \end{aligned}$$

The last term is a constant for θ_i , i.e.,

$$ELBO = \int q(\theta_i) \mathbb{E}_{\boldsymbol{\theta}-\{\theta_i\}} [\log P(\boldsymbol{\theta}, \mathbf{X})] d\theta_i - \int q(\theta_i) \log q(\theta_i) d\theta_i + const$$

As the expectation term is positive, we can view it as an log-probability $\log \tilde{P}$, then

$$ELBO = -D_{KL}[q(\theta_i)||\tilde{P}] + const$$

To maximum *ELBO* is equivalent to maximum the KL-divergence between each $q(\theta_i)$ and corresponding \tilde{P} .

Here, our task is just to find an MAP assignment, so to furthermore simplify the calculation, I set all q to be one-point distribution, in which situation, *ELBO* is maximized when

$$q(\theta_i) = \mathbb{I}\{\theta_i = \text{mode}\{\tilde{P}\}\}$$

Hereinafter, to simplify the notation, I will use $\theta_i = a$ to represent $q(\theta_i) = \mathbb{I}\{\theta_i = a\}$.

In our Gaussian Mixture Model, in each iteration step, I successively optimize $\boldsymbol{\pi}, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K, z_1, z_2, \dots, z_N$ as follows:

Optimize $\boldsymbol{\pi}$:

$$\begin{aligned} \log q(\boldsymbol{\pi}) &= \int q(z_1)q(z_2)\dots q(z_N)q(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)q(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)\dots q(\boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K) \\ &\quad \log \left\{ P(\boldsymbol{\pi}) \prod_{i=1}^k P(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \prod_{j=1}^N [P(z_j|\boldsymbol{\pi})P(\mathbf{x}_j|z_j, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K)] \right\} \\ &\quad d z_1 d z_2 \dots d z_N d \boldsymbol{\mu}_1 d \boldsymbol{\Sigma}_1 d \boldsymbol{\mu}_2 d \boldsymbol{\Sigma}_2 \dots d \boldsymbol{\mu}_K d \boldsymbol{\Sigma}_K \end{aligned}$$

All terms without $\boldsymbol{\pi}$ in it can be viewed as constants, and integrate out variables that do not occur in log term, we can get

$$\begin{aligned} \log q(\boldsymbol{\pi}) &= \int q(z_1)q(z_2)\dots q(z_N) \log P(\boldsymbol{\pi})P(z_1|\boldsymbol{\pi})P(z_2|\boldsymbol{\pi})\dots P(z_N|\boldsymbol{\pi}) d z_1 d z_2 \dots d z_N + \text{const} \\ &= \log P(\boldsymbol{\pi}) + \log P(z_1, z_2, \dots, z_N|\boldsymbol{\pi}) + \text{const} \end{aligned}$$

The last step is correct as $q(z_i)$ is one-point distribution according to my assumption.

Considering that $\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha})$ and $P(z_1, z_2, \dots, z_N|\boldsymbol{\pi}) \sim \text{Multi}(\boldsymbol{\pi})$, the right-hand side is just the posterior Dirichlet distribution. Therefore, $\boldsymbol{\pi} = \text{mode}\{\text{Dir}(\alpha_1 + N_1, \alpha_2 + N_2, \dots, \alpha_K + N_K)\}$ where N_k is the number of observations in cluster k .

Optimize $\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i$:

When optimizing $\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i$, given that the entanglement of them does not bring us any inconvenience because we have a well-formed distribution over $\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i$, I do not entangle them into two q . Similar to $\boldsymbol{\pi}$, we can cancel out irrelevant variables and only get

$$\log q(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) = \log P(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) + \sum_{z_j=i} \log P(\mathbf{x}_j|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) + \text{const}$$

The right-hand side is just the posterior Normal-Inverse-Wisart distribution. Thus, $\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j = \text{mode}\left\{\mathcal{NIW}\left(\frac{\lambda\boldsymbol{\mu}_0 + N_j\bar{\mathbf{x}}}{\lambda + N_j}, \lambda + N_j, \nu + N_j, \boldsymbol{\Psi} + \mathbf{C} + \frac{\lambda N_j}{\lambda + N_j}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T\right)\right\}$, where $\mathbf{C} = \sum_i (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$ and $\bar{\mathbf{x}}$ is the mean of samples in class j ;

Optimize z_1, z_2, \dots, z_N :

Via similar calculation, it is easy to find that each z_i can be optimized as $z_i = \max_{j=1}^k \{\pi_j P(\mathbf{x}_i|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)\}$.

Predict z for new \mathbf{x} :

After finding the MAP assignments for all parameters (i.e., one-point posterior distribution), we can predict class labels for new input \mathbf{x} .

$$\begin{aligned}
& P(z|\mathbf{x}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) \\
&= \int P(z|\mathbf{x}, \boldsymbol{\pi}, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K, z_1, z_2, \dots, z_N) \\
&\quad P(\boldsymbol{\pi}, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K, z_1, z_2, \dots, z_N | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) \\
&\quad d\boldsymbol{\pi} d\boldsymbol{\mu}_1 d\boldsymbol{\Sigma}_1 d\boldsymbol{\mu}_2 d\boldsymbol{\Sigma}_2 \dots d\boldsymbol{\mu}_K d\boldsymbol{\Sigma}_K dz_1 dz_2 \dots dz_N \\
&= \int P(z|\mathbf{x}, \boldsymbol{\pi}, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K) \\
&\quad P(\boldsymbol{\pi}, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K, z_1, z_2, \dots, z_N | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) \\
&\quad d\boldsymbol{\pi} d\boldsymbol{\mu}_1 d\boldsymbol{\Sigma}_1 d\boldsymbol{\mu}_2 d\boldsymbol{\Sigma}_2 \dots d\boldsymbol{\mu}_K d\boldsymbol{\Sigma}_K dz_1 dz_2 \dots dz_N \\
&= P(z|\mathbf{x}, \boldsymbol{\pi}^{MAP}, \boldsymbol{\mu}_1^{MAP}, \boldsymbol{\Sigma}_1^{MAP}, \boldsymbol{\mu}_2^{MAP}, \boldsymbol{\Sigma}_2^{MAP}, \dots, \boldsymbol{\mu}_K^{MAP}, \boldsymbol{\Sigma}_K^{MAP}) \\
&\propto P(\mathbf{x}|z, \boldsymbol{\mu}_1^{MAP}, \boldsymbol{\Sigma}_1^{MAP}, \boldsymbol{\mu}_2^{MAP}, \boldsymbol{\Sigma}_2^{MAP}, \dots, \boldsymbol{\mu}_K^{MAP}, \boldsymbol{\Sigma}_K^{MAP}) P(z|\boldsymbol{\pi}^{MAP})
\end{aligned}$$