

# Mean-field Variational Inference for Bayesian Mixture of Gaussian

December 27, 2022

## Abstract

Notes for the Bayesian Mixture of Gaussians section in Bishop's PRML. Compared to the book's treatment, this derivation more detailed. When choosing priors, we rely on Theorem 2.2 in Beal's PhD thesis, i.e., the priors should be conjugate to the complete-data log likelihood.

## Table of contents

1 Complete-data log likelihood	1
2 Prior for $\pi$	1
3 Joint prior for $\mu_k$ and $\Sigma_k$ :	1
4 Computing $q(Z)$ : variational E-step	1
5 Computing $q(\pi, \mu_{1:C}, \Lambda_{1:C})$ : variational M-step	1
5.1 $q^*(\pi)$	1
5.2 $q(\mu_k, \Lambda_k)$	1
6 Quantities required for computing $q(Z)$	1
7 Quantities required for computing $q^*(\pi)$ and $q(\mu_k, \Lambda_k)$	?

## 1 Complete-data log likelihood

Just write out the mixture of Gaussians assuming that the latent variables are observed:

$$\prod_{i=1}^n p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta}) = \prod_{i=1}^n \prod_{k=1}^c \pi_k^{z_{ik}} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})^{z_{ik}}$$

where  $\boldsymbol{\theta} = (\pi, \boldsymbol{\mu}_{1:k}, \boldsymbol{\Lambda}_{1:k})$ .

## 2 Prior for $\pi$

Dirichlet:

$$\text{Dir}(\pi | \boldsymbol{\alpha}_0) \propto \prod_{k=1}^c \pi_k^{(\alpha_0)_k - 1}$$

Relevant part of the complete-data log likelihood:

$$\prod_{i=1}^n \prod_{k=1}^c \pi_k^{z_{ik}}$$

Proof that they are conjugate:

$$\left( \prod_{k=1}^c \pi_k^{(\alpha_0)_k - 1} \right) \left( \prod_{i=1}^n \prod_{k=1}^c \pi_k^{z_{ik}} \right) = \left( \prod_{k=1}^c \pi_k^{(\alpha_0)_k - 1} \right) \left( \prod_{k=1}^c \pi_k^{\sum_{i=1}^n z_{ik}} \right) = \prod_{k=1}^c \pi_k^{(\alpha_0)_k + \sum_{i=1}^n z_{ik} - 1}$$

### 3 Joint prior for $\mu_k$ and $\Sigma_k$ :

Gaussian-Wishart:

$$p(\mu_k, \Lambda_k | \mathbf{m}_0, \beta_0, v_0) = \mathcal{N}(\mu_k | \mathbf{m}_0, (\beta_0 \Lambda_k)^{-1}) \text{Wi}(\Lambda_k | \mathbf{L}_0, v_0)$$

Relevant part of the complete-data log likelihood:

$$\prod_{i=1}^n \mathcal{N}(\mathbf{x}_i | \mu_k, \Lambda_k^{-1})^{z_{ik}} = \prod_{i \text{ s.t. } z_{i,k}=1} \mathcal{N}(\mathbf{x}_i | \mu_k, \Lambda_k^{-1})$$

Proof that they are conjugate to each other: omitted, see Section 4.6.3.3 of Murphy.

### 4 Computing $q(\mathbf{Z})$ : variational E-step

For notational clarity, we let  $\mu = \{\mu_k\}$  and  $\Lambda = \{\Lambda_k\}$ .

Starting from the CAVI update rule:

$$\begin{aligned} & \log q(\mathbf{Z}) \\ &= \mathbb{E}_{q(\pi, \mu, \Lambda)} [\log p(\mathbf{X}, \mathbf{Z}, \pi, \mu, \Lambda)] + \text{const} \\ &= \mathbb{E}_{q(\pi, \mu, \Lambda)} [\log p(\mathbf{Z} | \pi) + \log p(\mathbf{X} | \mathbf{Z}, \mu, \Lambda) + \log p(\pi, \mu, \Lambda)] + \text{const} \\ &= \mathbb{E}_{q(\pi, \mu, \Lambda)} [\log p(\mathbf{Z} | \pi) + \log p(\mathbf{X} | \mathbf{Z}, \mu, \Lambda)] + \text{const} \\ &= \mathbb{E}_{q(\pi, \mu, \Lambda)} [\log p(\mathbf{Z} | \pi)] + \mathbb{E}_{q(\pi, \mu, \Lambda)} [\log p(\mathbf{X} | \mathbf{Z}, \mu, \Lambda)] + \text{const} \\ &= \mathbb{E}_{q(\pi)} [\log p(\mathbf{Z} | \pi)] + \mathbb{E}_{q(\mu, \Lambda)} [\log p(\mathbf{X} | \mathbf{Z}, \mu, \Lambda)] + \text{const} \\ &= \mathbb{E}_{q(\pi)} \left[ \sum_{i=1}^n \sum_{k=1}^c z_{ik} \log \pi_k \right] + \mathbb{E}_{q(\mu, \Lambda)} \left[ \sum_{i=1}^n \sum_{k=1}^c z_{ik} \log \mathcal{N}(\mathbf{x}_i | \mu_k, \Lambda_k^{-1}) \right] + \text{const} \\ &= \sum_{i=1}^n \sum_{k=1}^c z_{ik} \mathbb{E}_{q(\pi)} [\log \pi_k] + \sum_{i=1}^n \sum_{k=1}^c z_{ik} \mathbb{E}_{q(\mu_k, \Lambda_k)} \left[ \frac{1}{2} \ln |\Lambda_k| - \frac{D}{2} \ln 2\pi - \frac{1}{2} (\mathbf{x}_i - \mu_k)^T \Lambda_k (\mathbf{x}_i - \mu_k) \right] + \text{const} \\ &= \sum_{i=1}^n \sum_{k=1}^c z_{ik} \mathbb{E}_{q(\pi)} [\log \pi_k] + \sum_{i=1}^n \sum_{k=1}^c z_{ik} \left[ \frac{1}{2} \mathbb{E}_{q(\Lambda_k)} [\ln |\Lambda_k|] - \frac{D}{2} \ln 2\pi - \frac{1}{2} \mathbb{E}_{q(\mu_k, \Lambda_k)} [(\mathbf{x}_i - \mu_k)^T \Lambda_k (\mathbf{x}_i - \mu_k)] \right] + \text{const} \\ &= \sum_{i=1}^n \sum_{k=1}^c z_{ik} \left( \mathbb{E}_{q(\pi)} [\log \pi_k] + \frac{1}{2} \mathbb{E}_{q(\Lambda_k)} [\ln |\Lambda_k|] - \frac{D}{2} \ln 2\pi - \frac{1}{2} \mathbb{E}_{q(\mu_k, \Lambda_k)} [(\mathbf{x}_i - \mu_k)^T \Lambda_k (\mathbf{x}_i - \mu_k)] \right) + \text{const} \\ &\triangleq \sum_{i=1}^n \sum_{k=1}^c z_{ik} \log \rho_{ik} \end{aligned}$$

Therefore

$$q(\mathbf{Z}) \propto \prod_{i=1}^n \prod_{k=1}^c \rho_{ik}^{z_{ik}}$$

Requiring that the distributuon is normalized on each row:

$$q(\mathbf{Z}) = \prod_{i=1}^n \prod_{k=1}^c r_{ik}^{z_{ik}}$$

where  $r_{ik}$  is the  $\rho_{ik}$  divided by row sum.

## 5 Computing $q(\boldsymbol{\pi}, \boldsymbol{\mu}_{1:C}, \boldsymbol{\Lambda}_{1:C})$ : variational M-step

Starting from the CAVI update rule:

$$\begin{aligned} & \ln q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \\ &= \mathbb{E}_{q(\mathbf{Z})}[\log p(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \mathbf{Z}, \mathbf{X})] + \text{const} \\ &= \mathbb{E}_{q(\mathbf{Z})}[\log p(\boldsymbol{\pi}) + \log p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) + \log p(\mathbf{Z} | \boldsymbol{\pi}) + \log p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] + \text{const} \\ &= \log p(\boldsymbol{\pi}) + \log p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) + \mathbb{E}_{q(\mathbf{Z})}[\log p(\mathbf{Z} | \boldsymbol{\pi})] + \mathbb{E}_{q(\mathbf{Z})}[\log p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] + \text{const} \\ &= (\log p(\boldsymbol{\pi}) + \mathbb{E}_{q(\mathbf{Z})}[\log p(\mathbf{Z} | \boldsymbol{\pi})]) + \sum_{k=1}^C \left( p(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) + \sum_{i=1}^N r_{ik} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) \right) + \text{const}, \quad (1) \end{aligned}$$

which contains separate terms for  $\boldsymbol{\pi}$  and each  $(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$ . This implies that the approximate posterior further factors into

$$\begin{aligned} \log q(\boldsymbol{\pi}, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_C, \boldsymbol{\Lambda}_1, \dots, \boldsymbol{\Lambda}_C) &= \log q(\boldsymbol{\pi}) + \sum_{k=1}^C \ln q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) \\ q(\boldsymbol{\pi}, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_C, \boldsymbol{\Lambda}_1, \dots, \boldsymbol{\Lambda}_C) &= q(\boldsymbol{\pi}) \prod_{k=1}^C q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) \end{aligned}$$

so that we can optimize over each component distribution separately.

### 5.1 $q(\boldsymbol{\pi})$

Isolating the terms containing  $\boldsymbol{\pi}$  from Equation 1, we see that

$$\begin{aligned} \log q^*(\boldsymbol{\pi}) &= \log p(\boldsymbol{\pi}) + \mathbb{E}_{q(\mathbf{Z})}[\log p(\mathbf{Z} | \boldsymbol{\pi})] + \text{const} \\ &= \sum_{k=1}^C (\alpha_0 - 1) \ln \pi_k + \sum_{k=1}^C \left( \sum_{i=1}^N r_{ik} \right) \ln \pi_k + \text{const} \\ &= \sum_{k=1}^C \left( \alpha_0 + \sum_{i=1}^N r_{ik} - 1 \right) \ln \pi_k + \sum_{k=1}^C \ln \pi_k + \text{const} \\ &\Rightarrow \\ q^*(\boldsymbol{\pi}) &= \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}) \quad \text{with} \quad \alpha_k = \alpha_0 + N_k \end{aligned}$$

### 5.2 $q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$

Isolating the terms containing  $(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$  from Equation 1, we see that

$$\log q^*(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) = \log p(\boldsymbol{\mu}_k | \boldsymbol{\Lambda}_k) + \log p(\boldsymbol{\Lambda}_k) + \sum_{i=1}^N r_{ik} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) + \text{const}$$

By the Gaussian-Wishart prior we specified earlier, we have

$$\begin{aligned}
p(\boldsymbol{\mu}_k | \boldsymbol{\Lambda}_k) &= \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, (\beta_0 \boldsymbol{\Lambda}_k)^{-1}) \\
&\propto |\beta_0 \boldsymbol{\Lambda}_k|^{1/2} \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k - \mathbf{m}_0)^T (\beta_0 \boldsymbol{\Lambda}_k) (\boldsymbol{\mu}_k - \mathbf{m}_0)\right) \\
&\propto |\boldsymbol{\Lambda}_k|^{1/2} \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k - \mathbf{m}_0)^T (\beta_0 \boldsymbol{\Lambda}_k) (\boldsymbol{\mu}_k - \mathbf{m}_0)\right) \\
&\Rightarrow \\
\log p(\boldsymbol{\mu}_k | \boldsymbol{\Lambda}_k) &= \frac{1}{2} \log |\boldsymbol{\Lambda}_k| - \frac{1}{2}(\boldsymbol{\mu}_k - \mathbf{m}_0)^T (\beta_0 \boldsymbol{\Lambda}_k) (\boldsymbol{\mu}_k - \mathbf{m}_0) + \text{const}
\end{aligned}$$

and

$$\begin{aligned}
p(\boldsymbol{\Lambda}_k) &= \mathcal{W}(\boldsymbol{\Lambda}_k | \mathbf{W}_0, v_0) \\
&\propto |\boldsymbol{\Lambda}_k|^{(v_0 - C - 1)/2} \exp\left(-\frac{1}{2} \text{Tr}(\mathbf{W}_0^{-1} \boldsymbol{\Lambda}_k)\right) \\
&\Rightarrow \\
\log p(\boldsymbol{\Lambda}_k) &= \frac{(v_0 - C - 1)}{2} \log |\boldsymbol{\Lambda}_k| - \frac{1}{2} \text{Tr}(\mathbf{W}_0^{-1} \boldsymbol{\Lambda}_k) + \text{const}.
\end{aligned}$$

The normal likelihoods for all the observations can be written as

$$\sum_{i=1}^N r_{ik} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) = \sum_{i=1}^N r_{ik} \left( \frac{1}{2} \log |\boldsymbol{\Lambda}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Lambda}_k (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) + \text{const}.$$

Now, we can assemble all terms together:

$$\begin{aligned}
\log q^*(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) &= \frac{1}{2} \log |\boldsymbol{\Lambda}_k| - \frac{1}{2}(\boldsymbol{\mu}_k - \mathbf{m}_0)^T (\beta_0 \boldsymbol{\Lambda}_k) (\boldsymbol{\mu}_k - \mathbf{m}_0) \\
&\quad + \frac{(v_0 - C - 1)}{2} \log |\boldsymbol{\Lambda}_k| - \frac{1}{2} \text{Tr}(\mathbf{W}_0^{-1} \boldsymbol{\Lambda}_k) \\
&\quad + \sum_{i=1}^N r_{ik} \left( \frac{1}{2} \log |\boldsymbol{\Lambda}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Lambda}_k (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) + \text{const}.
\end{aligned}$$

It may not be immediate clear what we should do here, so it's good to remind ourselves that by Theorem 2 of Beal tells us that the approximate posterior over  $(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$  would also be Gaussian-Wishart.

Firstly, we see that

$$\frac{(v_0 - C - 1)}{2} \log |\boldsymbol{\Lambda}_k| + \frac{(\sum_{i=1}^N r_{ik})}{2} \log |\boldsymbol{\Lambda}_k| = \frac{(v_0 + N_k - C - 1)}{2} \log |\boldsymbol{\Lambda}_k|,$$

which means that the Wishart posterior distribution on  $\boldsymbol{\Lambda}_k$  have DOF  $v_k = v_0 + N - C$ .

Secondly, we need manipulate these terms

$$-\frac{1}{2}(\boldsymbol{\mu}_k - \mathbf{m}_0)^T (\beta_0 \boldsymbol{\Lambda}_k) (\boldsymbol{\mu}_k - \mathbf{m}_0) \quad -\frac{1}{2} \text{Tr}(\mathbf{W}_0^{-1} \boldsymbol{\Lambda}_k) \quad \sum_{i=1}^N r_{ik} \left( -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Lambda}_k (\mathbf{x}_i - \boldsymbol{\mu}_k) \right)$$

into two terms of the form

$$-\frac{1}{2}(\boldsymbol{\mu}_k - \mathbf{m}_0)^T (\beta_0 \boldsymbol{\Lambda}_k) (\boldsymbol{\mu}_k - \mathbf{m}_0) \quad -\frac{1}{2} \text{Tr}(\mathbf{W}_0^{-1} \boldsymbol{\Lambda}_k)$$

where the  $\mu_k$  in the second term shouldn't contain  $\mu_k$  because it would be a parameter of the Wishart posterior. For simplicity, we will drop the shared  $-1/2$  term shortly.

To begin, we do an expansion with a trick:

$$\begin{aligned}
& \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Lambda}_k (\mathbf{x}_i - \boldsymbol{\mu}_k) \\
&= \sum_{i=1}^N r_{ik} ((\mathbf{x}_i - \bar{\mathbf{x}}_k) - (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k))^T \boldsymbol{\Lambda}_k ((\mathbf{x}_i - \bar{\mathbf{x}}_k) - (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k)) \\
&= \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \bar{\mathbf{x}}_k) \boldsymbol{\Lambda}_k (\mathbf{x}_i - \bar{\mathbf{x}}_k)^T - 2 (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Lambda}_k \underbrace{\sum_{i=1}^N r_{ik} (\mathbf{x}_i - \bar{\mathbf{x}}_k)}_0 + N_k (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Lambda}_k (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k) \\
&= N_k (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Lambda}_k (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k) + \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \bar{\mathbf{x}}_k)^T \boldsymbol{\Lambda}_k (\mathbf{x}_i - \bar{\mathbf{x}}_k) \\
&= N_k (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Lambda}_k (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k) + \text{Tr} \left[ \left( \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \bar{\mathbf{x}}_k) (\mathbf{x}_i - \bar{\mathbf{x}}_k)^T \right) \boldsymbol{\Lambda}_k \right],
\end{aligned}$$

which is great because the first term here contains  $\boldsymbol{\mu}_k$  while the second term here doesn't and can be directly merged with  $\text{Tr}(\mathbf{W}_0^{-1} \boldsymbol{\Lambda}_k)$ .

We can merge the first term here with

$$\begin{aligned}
& \beta_0 (\boldsymbol{\mu}_k - \mathbf{m}_0)^T \boldsymbol{\Lambda}_k (\boldsymbol{\mu}_k - \mathbf{m}_0) + N_k (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Lambda}_k (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k) \\
&= \text{Tr} [(\beta_0 (\boldsymbol{\mu}_k - \mathbf{m}_0) (\boldsymbol{\mu}_k - \mathbf{m}_0)^T + N_k (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k) (\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k)^T) \boldsymbol{\Lambda}_k] \\
&= \text{Tr} [(\beta_0 (\boldsymbol{\mu}_k \boldsymbol{\mu}_k^T - \boldsymbol{\mu}_k \mathbf{m}_0^T - \mathbf{m}_0 \boldsymbol{\mu}_k^T + \mathbf{m}_0 \mathbf{m}_0^T) + N_k (\boldsymbol{\mu}_k \boldsymbol{\mu}_k^T - \boldsymbol{\mu}_k \bar{\mathbf{x}}_k^T - \bar{\mathbf{x}}_k \boldsymbol{\mu}_k^T + \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T)) \boldsymbol{\Lambda}_k] \\
&= \text{Tr} [((\beta_0 + N_k) \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T - \boldsymbol{\mu}_k (\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)^T - (\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k) \boldsymbol{\mu}_k^T + \beta_0 \mathbf{m}_0 \mathbf{m}_0^T + N_k \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T) \boldsymbol{\Lambda}_k] \\
&= \text{Tr} \left[ \left( (\beta_0 + N_k) \left( \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T - \boldsymbol{\mu}_k \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)^T}{\beta_0 + N_k} - \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)^T}{\beta_0 + N_k} \boldsymbol{\mu}_k^T \right) + \beta_0 \mathbf{m}_0 \mathbf{m}_0^T + N_k \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T \right) \boldsymbol{\Lambda}_k \right] \\
&= \text{Tr} \left[ \left( (\beta_0 + N_k) \left( \boldsymbol{\mu}_k - \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)}{\beta_0 + N_k} \right) \left( \boldsymbol{\mu}_k - \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)}{\beta_0 + N_k} \right)^T - \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k) (\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)^T}{\beta_0 + N_k} + \beta_0 \mathbf{m}_0 \mathbf{m}_0^T + N_k \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T \right) \boldsymbol{\Lambda}_k \right] \text{ (complete the square)} \\
&= \text{Tr} \left[ \left( \dots - \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k) (\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)^T}{\beta_0 + N_k} + \beta_0 \mathbf{m}_0 \mathbf{m}_0^T + N_k \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T \right) \boldsymbol{\Lambda}_k \right] \\
&= \text{Tr} \left[ \left( \dots - \frac{\beta_0^2 \mathbf{m}_0 \mathbf{m}_0^T}{\beta_0 + N_k} - \frac{\beta_0 N_k \mathbf{m}_0 \bar{\mathbf{x}}_k^T}{\beta_0 + N_k} - \frac{\beta_0 N_k \bar{\mathbf{x}}_k \mathbf{m}_0^T}{\beta_0 + N_k} - \frac{N_k^2 \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T}{\beta_0 + N_k} + \beta_0 \mathbf{m}_0 \mathbf{m}_0^T + N_k \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T \right) \boldsymbol{\Lambda}_k \right] \\
&= \text{Tr} \left[ \left( \dots - \frac{\beta_0^2 \mathbf{m}_0 \mathbf{m}_0^T}{\beta_0 + N_k} - \frac{\beta_0 N_k \mathbf{m}_0 \bar{\mathbf{x}}_k^T}{\beta_0 + N_k} - \frac{\beta_0 N_k \bar{\mathbf{x}}_k \mathbf{m}_0^T}{\beta_0 + N_k} - \frac{N_k^2 \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T}{\beta_0 + N_k} + \frac{\beta_0^2 + \beta_0 N_k}{\beta_0 + N_k} \mathbf{m}_0 \mathbf{m}_0^T + \frac{\beta_0 N_k + N_k^2}{\beta_0 + N_k} \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T \right) \boldsymbol{\Lambda}_k \right] \\
&= \text{Tr} \left[ \left( \dots - \frac{\beta_0 N_k \mathbf{m}_0 \bar{\mathbf{x}}_k^T}{\beta_0 + N_k} - \frac{\beta_0 N_k \bar{\mathbf{x}}_k \mathbf{m}_0^T}{\beta_0 + N_k} + \frac{\beta_0 N_k}{\beta_0 + N_k} \mathbf{m}_0 \mathbf{m}_0^T + \frac{\beta_0 N_k}{\beta_0 + N_k} \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T \right) \boldsymbol{\Lambda}_k \right] \\
&= \text{Tr} \left[ \left( \dots + \left( \frac{\beta_0 N_k}{\beta_0 + N_k} \right) (\mathbf{m}_0 \mathbf{m}_0^T - \mathbf{m}_0 \bar{\mathbf{x}}_k^T - \bar{\mathbf{x}}_k \mathbf{m}_0^T + \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T) \right) \boldsymbol{\Lambda}_k \right] \\
&= \text{Tr} \left[ \left( \dots + \left( \frac{\beta_0 N_k}{\beta_0 + N_k} \right) (\bar{\mathbf{x}}_k - \mathbf{m}_0) (\bar{\mathbf{x}}_k - \mathbf{m}_0)^T \right) \boldsymbol{\Lambda}_k \right]
\end{aligned}$$

$$\begin{aligned}
&= \text{Tr} \left[ \left( \cdots + \left( \frac{\beta_0 N_k}{\beta_0 + N_k} \right) (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T \right) \mathbf{\Lambda}_k \right] \\
&= \text{Tr} \left[ \left( (\beta_0 + N_k) \left( \boldsymbol{\mu}_k - \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)}{\beta_0 + N_k} \right) \left( \boldsymbol{\mu}_k - \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)}{\beta_0 + N_k} \right)^T + \left( \frac{\beta_0 N_k}{\beta_0 + N_k} \right) (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T \right) \mathbf{\Lambda}_k \right] \\
&= \left( \boldsymbol{\mu}_k - \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)}{\beta_0 + N_k} \right)^T ((\beta_0 + N_k) \mathbf{\Lambda}_k) \left( \boldsymbol{\mu}_k - \frac{(\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k)}{\beta_0 + N_k} \right) + \text{Tr} \left[ \left( \frac{\beta_0 N_k}{\beta_0 + N_k} \right) (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T \mathbf{\Lambda}_k \right],
\end{aligned}$$

where the first term here has the required form of  $(\boldsymbol{\mu}_k - \mathbf{?}_1)^T (\mathbf{?} \mathbf{\Lambda}_k) (\boldsymbol{\mu}_k - \mathbf{?}_1)$  and the second term can be directly merged with  $\text{Tr}(\mathbf{W}_0^{-1} \mathbf{\Lambda}_k)$ .

Finally, we have

$$\begin{aligned}
&\log q^*(\boldsymbol{\mu}_k, \mathbf{\Lambda}_k) \\
&= -\frac{1}{2} \left( \boldsymbol{\mu}_k - \frac{\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k}{\beta_0 + N_k} \right)^T ((\beta_0 + N_k) \mathbf{\Lambda}_k) \left( \boldsymbol{\mu}_k - \frac{\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k}{\beta_0 + N_k} \right) + \\
&\quad \frac{(v_0 + \sum_{i=1}^N r_{ik} - 1)}{2} \log |\mathbf{\Lambda}_k| - \frac{1}{2} \text{Tr} \left( \left( \mathbf{W}_0^{-1} + N_k (\mathbf{x}_i - \bar{\mathbf{x}}_k)(\mathbf{x}_i - \bar{\mathbf{x}}_k)^T + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T \right) \mathbf{\Lambda}_k \right) + \text{const} \\
&= \log \mathcal{N} \left( \boldsymbol{\mu}_k \left| \frac{\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k}{\beta_0 + N_k}, ((\beta_0 + N_k) \mathbf{\Lambda}_k)^{-1} \right. \right) + \log \mathcal{W} \left( \mathbf{\Lambda}_k \left| \left( \mathbf{W}_0^{-1} + N_k (\mathbf{x}_n - \bar{\mathbf{x}}_k)(\mathbf{x}_n - \bar{\mathbf{x}}_k)^T + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T \right)^{-1}, v_0 + N_k \right. \right).
\end{aligned}$$

Or we can define

$$\begin{aligned}
\mathbf{m}_k &= \frac{\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k}{\beta_0 + N_k} \\
\beta_k &= \beta_0 + N_k \\
\mathbf{W}_k^{-1} &= \mathbf{W}_0^{-1} + N_k (\mathbf{x}_n - \bar{\mathbf{x}}_k)(\mathbf{x}_n - \bar{\mathbf{x}}_k)^T + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T \\
v_k &= v_0 + N_k
\end{aligned}$$

This results turn out to be very similar to the standard Gaussian-Wishart posterior without soft assignment. The approach here should be directly applicable to that, too.

## 6 Quantities required for computing $q(\mathbf{Z})$

Standard results:

$$\rho_{ik} = \mathbb{E}_{q(\pi)}[\log \pi_k] + \frac{1}{2} \mathbb{E}_{q(\mathbf{\Lambda}_k)}[\ln |\mathbf{\Lambda}_k|] - \frac{D}{2} \ln 2\pi - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\mu}_k, \mathbf{\Lambda}_k)}[(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \mathbf{\Lambda}_k (\mathbf{x}_i - \boldsymbol{\mu}_k)]$$

Please refer to Bishop's PRML from computing these expectations.

## 7 Quantities required for computing $q(\pi)$ and $q(\boldsymbol{\mu}_k, \mathbf{\Lambda}_k)$

Please refer to Bishop's PRML from computing the associated expectations.