

# Final Exam

20213073 Donggyu Kim

June 17, 2021

1. (a) The marginal distribution  $p(\mathbf{x}; \boldsymbol{\theta})$  is

$$\begin{aligned}
 p(\mathbf{x}; \boldsymbol{\theta}) &= \int p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \\
 &= \int p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) p(\mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \\
 &= \int \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}_d) \mathcal{N}(\mathbf{z}|\mathbf{0}_h, \mathbf{I}_h) d\mathbf{z} \\
 &= \int \frac{1}{(2\pi)^{\frac{d+h}{2}} \sigma^d} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu})^\top (\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu}) - \frac{1}{2}\mathbf{z}^\top \mathbf{z}\right) d\mathbf{z}
 \end{aligned}$$

The exponent part can be rewritten as follows

$$\begin{aligned}
 &-\frac{1}{2\sigma^2}(\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu})^\top (\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu}) - \frac{1}{2}\mathbf{z}^\top \mathbf{z} \\
 &= -\frac{1}{2\sigma^2} \left( \mathbf{z}^\top \mathbf{V} \mathbf{z} - 2\mathbf{z}^\top \mathbf{W}^\top (\mathbf{x} - \boldsymbol{\mu}) + (\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{x} - \boldsymbol{\mu}) \right) \\
 &= -\frac{1}{2\sigma^2} \left( (\mathbf{z} - \boldsymbol{\alpha})^\top \mathbf{V} (\mathbf{z} - \boldsymbol{\alpha}) + (\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{x} - \boldsymbol{\mu}) - \boldsymbol{\alpha}^\top \mathbf{V} \boldsymbol{\alpha} \right)
 \end{aligned}$$

where  $\mathbf{V} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_h$  and  $\boldsymbol{\alpha} = \mathbf{V}^{-1} \mathbf{W}^\top (\mathbf{x} - \boldsymbol{\mu})$ .

Using the matrix inverse identity, we further get

$$\begin{aligned}
 &(\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{x} - \boldsymbol{\mu}) - \boldsymbol{\alpha}^\top \mathbf{V} \boldsymbol{\alpha} \\
 &= (\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{I}_d - \mathbf{W} \mathbf{V}^{-1} \mathbf{W}^\top) (\mathbf{x} - \boldsymbol{\mu}) \\
 &= (\mathbf{x} - \boldsymbol{\mu})^\top \left( \mathbf{I}_d - (\mathbf{W}/\sigma) \left( \mathbf{I}_h + \mathbf{W}^\top \mathbf{W} / \sigma^2 \right)^{-1} (\mathbf{W}^\top / \sigma) \right) (\mathbf{x} - \boldsymbol{\mu}) \\
 &= (\mathbf{x} - \boldsymbol{\mu})^\top \left( \mathbf{I}_d + \mathbf{W} \mathbf{W}^\top / \sigma^2 \right)^{-1} (\mathbf{x} - \boldsymbol{\mu})
 \end{aligned}$$

Using the fact  $\det(\mathbf{W} \mathbf{W}^\top + \sigma^2 \mathbf{I}_d) = \sigma^{2d} \det(\mathbf{I}_h + \mathbf{W}^\top \mathbf{W} / \sigma^2)$ , we have

$$\begin{aligned}
 &\int \exp\left(-\frac{1}{2} \left( (\mathbf{z} - \boldsymbol{\alpha})^\top (\mathbf{V}/\sigma^2) (\mathbf{z} - \boldsymbol{\alpha}) \right)\right) d\mathbf{z} \\
 &= \int (2\pi)^{\frac{h}{2}} \det((\mathbf{V}/\sigma^2)^{-1})^{\frac{1}{2}} \mathcal{N}(\mathbf{z}|\boldsymbol{\alpha}, (\mathbf{V}/\sigma^2)^{-1}) d\mathbf{z} \\
 &= (2\pi)^{\frac{h}{2}} \det\left(\mathbf{I}_h + \mathbf{W}^\top \mathbf{W} / \sigma^2\right)^{-\frac{1}{2}} \\
 &= (2\pi)^{\frac{h}{2}} \sigma^d \det\left(\mathbf{W} \mathbf{W}^\top + \sigma^2 \mathbf{I}_d\right)^{-\frac{1}{2}}
 \end{aligned}$$

Thus, we can finally get

$$\begin{aligned}
p(\mathbf{x}; \boldsymbol{\theta}) &= \int \frac{1}{(2\pi)^{\frac{d+h}{2}} \sigma^d} \exp \left( -\frac{1}{2} \left( (\mathbf{z} - \boldsymbol{\alpha})^\top (\mathbf{V}/\sigma^2) (\mathbf{z} - \boldsymbol{\alpha}) \right. \right. \\
&\quad \left. \left. + (\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right) \right) d\mathbf{z} \\
&= \frac{(2\pi)^{\frac{h}{2}} \sigma^d \det(\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)^{-\frac{1}{2}}}{(2\pi)^{\frac{d+h}{2}} \sigma^d} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right) \\
&= \frac{1}{(2\pi)^{\frac{d}{2}} \det(\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right) \\
&= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)
\end{aligned}$$

And the conditional distribution  $p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})$  is

$$\begin{aligned}
p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}) &= \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{p(\mathbf{x}; \boldsymbol{\theta})} \\
&= \frac{p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) p(\mathbf{z}; \boldsymbol{\theta})}{p(\mathbf{x}; \boldsymbol{\theta})} \\
&= \frac{\mathcal{N}(\mathbf{x} | \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}_d) \mathcal{N}(\mathbf{z} | \mathbf{0}_h, \mathbf{I}_h)}{\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)} \\
&= \frac{1}{(2\pi)^{\frac{h}{2}} \sigma^d \det(\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)^{-\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{z} - \boldsymbol{\alpha})^\top (\mathbf{V}/\sigma^2) (\mathbf{z} - \boldsymbol{\alpha}) \right) \\
&= \frac{1}{(2\pi)^{\frac{h}{2}} \det(\mathbf{I}_h + \mathbf{W}^\top \mathbf{W} / \sigma^2)^{-\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{z} - \boldsymbol{\alpha})^\top (\mathbf{V}/\sigma^2) (\mathbf{z} - \boldsymbol{\alpha}) \right) \\
&= \frac{1}{(2\pi)^{\frac{h}{2}} \det((\mathbf{V}/\sigma^2)^{-1})^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{z} - \boldsymbol{\alpha})^\top (\mathbf{V}/\sigma^2) (\mathbf{z} - \boldsymbol{\alpha}) \right) \\
&= \mathcal{N}(\mathbf{z} | \boldsymbol{\alpha}, (\mathbf{V}/\sigma^2)^{-1}) \\
&= \mathcal{N}(\mathbf{z} | \mathbf{V}^{-1} \mathbf{W}^\top (\mathbf{x} - \boldsymbol{\mu}), \sigma^2 \mathbf{V}^{-1})
\end{aligned}$$

(b) The log-likelihood  $\log p(\mathbf{X})$  is as follows.

$$\begin{aligned}
\log p(\mathbf{X}) &= \sum_{i=1}^n \log p(\mathbf{x}_i; \boldsymbol{\theta}) \\
&= \sum_{i=1}^n \left( -\frac{d}{2} \log 2\pi - \frac{1}{2} \log \det(\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d) \right. \\
&\quad \left. - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right)
\end{aligned}$$

Then we can derive its maximum-likelihood estimator of the parameter  $\boldsymbol{\mu}$ .

$$\begin{aligned}
\frac{\partial \log p(\mathbf{X})}{\partial \boldsymbol{\mu}} &= \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)^{-1} \\
&= \left( \sum_{i=1}^n \mathbf{x}_i - n\boldsymbol{\mu} \right)^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d)^{-1} = 0 \Rightarrow \boldsymbol{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i
\end{aligned}$$

(c) The expectation of log-likelihood computed in E-step is

$$\begin{aligned} \mathbb{E}_{p(\mathbf{z}_i|\mathbf{x}_i, \boldsymbol{\theta}_t)}[\log p(\mathbf{x}_i, \mathbf{z}_i; \boldsymbol{\theta})] &= -\frac{d+h}{2} \log 2\pi - d \log \sigma - \frac{1}{2\sigma^2} (\mathbf{x}_i - \boldsymbol{\mu})^\top (\mathbf{x}_i - \boldsymbol{\mu}) \\ &\quad - \frac{1}{2\sigma^2} \mathbb{E}_{p(\mathbf{z}_i|\mathbf{x}_i, \boldsymbol{\theta}_t)} \left[ \mathbf{z}_i^\top \mathbf{V} \mathbf{z}_i - 2\mathbf{z}_i^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \right] \end{aligned}$$

Let  $\boldsymbol{\alpha}_{it} = \mathbf{V}_t^{-1} \mathbf{W}_t^\top (\mathbf{x}_i - \boldsymbol{\mu})$ . Then,

$$\begin{aligned} &\mathbb{E}_{p(\mathbf{z}_i|\mathbf{x}_i, \boldsymbol{\theta}_t)} \left[ \mathbf{z}_i^\top \mathbf{V} \mathbf{z}_i - 2\mathbf{z}_i^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \right] \\ &= \mathbb{E}_{\mathcal{N}(\mathbf{z}_i|\boldsymbol{\alpha}_{it}, \sigma_t^2 \mathbf{V}_t^{-1})} \left[ (\mathbf{z}_i - \boldsymbol{\alpha}_{it})^\top \mathbf{V} (\mathbf{z}_i - \boldsymbol{\alpha}_{it}) + 2\mathbf{z}_i^\top \mathbf{V} \boldsymbol{\alpha}_{it} - \boldsymbol{\alpha}_{it}^\top \mathbf{V} \boldsymbol{\alpha}_{it} - 2\mathbf{z}_i^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \right] \\ &= \mathbb{E}_{\mathcal{N}(\mathbf{z}_i|\boldsymbol{\alpha}_{it}, \sigma_t^2 \mathbf{V}_t^{-1})} \left[ (\mathbf{z}_i - \boldsymbol{\alpha}_{it})^\top \mathbf{V} (\mathbf{z}_i - \boldsymbol{\alpha}_{it}) \right] + \boldsymbol{\alpha}_{it}^\top \mathbf{V} \boldsymbol{\alpha}_{it} - 2\boldsymbol{\alpha}_{it}^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= \mathbb{E}_{\mathcal{N}(\mathbf{z}_i|\boldsymbol{\alpha}_{it}, \sigma_t^2 \mathbf{V}_t^{-1})} \left[ \text{Tr} \left( \mathbf{V} (\mathbf{z}_i - \boldsymbol{\alpha}_{it}) (\mathbf{z}_i - \boldsymbol{\alpha}_{it})^\top \right) \right] + \boldsymbol{\alpha}_{it}^\top \mathbf{V} \boldsymbol{\alpha}_{it} - 2\boldsymbol{\alpha}_{it}^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= \text{Tr} \left( \mathbb{E}_{\mathcal{N}(\mathbf{z}_i|\boldsymbol{\alpha}_{it}, \sigma_t^2 \mathbf{V}_t^{-1})} \left[ \mathbf{V} (\mathbf{z}_i - \boldsymbol{\alpha}_{it}) (\mathbf{z}_i - \boldsymbol{\alpha}_{it})^\top \right] \right) + \boldsymbol{\alpha}_{it}^\top \mathbf{V} \boldsymbol{\alpha}_{it} - 2\boldsymbol{\alpha}_{it}^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= \text{Tr} \left( \mathbf{V} \mathbb{E}_{\mathcal{N}(\mathbf{z}_i|\boldsymbol{\alpha}_{it}, \sigma_t^2 \mathbf{V}_t^{-1})} \left[ (\mathbf{z}_i - \boldsymbol{\alpha}_{it}) (\mathbf{z}_i - \boldsymbol{\alpha}_{it})^\top \right] \right) + \boldsymbol{\alpha}_{it}^\top \mathbf{V} \boldsymbol{\alpha}_{it} - 2\boldsymbol{\alpha}_{it}^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= \sigma_t^2 \text{Tr} (\mathbf{V} \mathbf{V}_t^{-1}) + \boldsymbol{\alpha}_{it}^\top \mathbf{V} \boldsymbol{\alpha}_{it} - 2\boldsymbol{\alpha}_{it}^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \end{aligned}$$

Since  $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$  is fixed, we can also compute

$$\begin{aligned} &\sum_{i=1}^n \left( \boldsymbol{\alpha}_{it}^\top \mathbf{V} \boldsymbol{\alpha}_{it} - 2\boldsymbol{\alpha}_{it}^\top \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \right) \\ &= \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{V} \mathbf{V}_t^{-1} \mathbf{W}_t^\top (\mathbf{x}_i - \boldsymbol{\mu}) - 2 \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= n \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{V} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \right) - 2n \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{S} \right) \end{aligned}$$

Therefore, the complete-data log-likelihood  $\tilde{\mathcal{L}}(\boldsymbol{\theta})$  is computed as

$$\begin{aligned} \tilde{\mathcal{L}}(\boldsymbol{\theta}) &= \sum_{i=1}^n \mathbb{E}_{p(\mathbf{z}_i|\mathbf{x}_i, \boldsymbol{\theta}_t)} [\log p(\mathbf{x}_i, \mathbf{z}_i; \boldsymbol{\theta})] \\ &= -\frac{n(d+h)}{2} \log 2\pi - nd \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top (\mathbf{x}_i - \boldsymbol{\mu}) - \frac{n\sigma_t^2}{2\sigma^2} \text{Tr} (\mathbf{V} \mathbf{V}_t^{-1}) \\ &\quad - \frac{n}{2\sigma^2} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{V} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \right) + \frac{n}{\sigma^2} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{S} \right) \\ &= -\frac{n(d+h)}{2} \log 2\pi - nd \log \sigma - \frac{n}{2\sigma^2} \text{Tr} (\mathbf{S}) - \frac{n\sigma_t^2}{2\sigma^2} \text{Tr} (\mathbf{V} \mathbf{V}_t^{-1}) \\ &\quad - \frac{n}{2\sigma^2} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{V} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \right) + \frac{n}{\sigma^2} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{S} \right) \end{aligned}$$

Taking the gradient w.r.t.  $\mathbf{W}$ ,

$$\frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta})}{\partial \mathbf{W}} = -\frac{n\sigma_t^2}{2\sigma^2} \frac{\partial \text{Tr} (\mathbf{V} \mathbf{V}_t^{-1})}{\partial \mathbf{W}} - \frac{n}{2\sigma^2} \frac{\partial \text{Tr} (\mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{V} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S})}{\partial \mathbf{W}} + \frac{n}{\sigma^2} \frac{\partial \text{Tr} (\mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{S})}{\partial \mathbf{W}}$$

First we compute

$$\begin{aligned}\frac{\partial \text{Tr}(\mathbf{V}\mathbf{V}_t^{-1})}{\partial \mathbf{W}} &= \frac{\partial \text{Tr}((\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_h) \mathbf{V}_t^{-1})}{\partial \mathbf{W}} \\ &= \frac{\partial \text{Tr}(\mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1})}{\partial \mathbf{W}} \\ &= 2\mathbf{V}_t^{-1} \mathbf{W}^\top\end{aligned}$$

Next,

$$\begin{aligned}\frac{\partial \text{Tr}(\mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{V} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S})}{\partial \mathbf{W}} &= \frac{\partial \text{Tr}(\mathbf{W}_t \mathbf{V}_t^{-1} (\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_h) \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S})}{\partial \mathbf{W}} \\ &= \frac{\partial \text{Tr}(\mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S})}{\partial \mathbf{W}} \\ &= \left( \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} + \left( \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \right)^\top \right) \mathbf{W}^\top \\ &= 2\mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top\end{aligned}$$

Lastly,

$$\frac{\partial \text{Tr}(\mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{S})}{\partial \mathbf{W}} = (\mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1})^\top = \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S}$$

So the final form of the gradient w.r.t.  $\mathbf{W}$  is

$$\begin{aligned}\frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta})}{\partial \mathbf{W}} &= -\frac{n\sigma_t^2}{\sigma^2} \mathbf{V}_t^{-1} \mathbf{W}^\top - \frac{n}{\sigma^2} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top + \frac{n}{\sigma^2} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \\ &= -\frac{n}{\sigma^2} \mathbf{V}_t^{-1} \left( \left( \sigma_t^2 \mathbf{I}_h + \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \right) \mathbf{W}^\top - \mathbf{W}_t^\top \mathbf{S} \right)\end{aligned}$$

Thus, the M-step update equation for  $\mathbf{W}_{t+1}$  is

$$\begin{aligned}\frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta})}{\partial \mathbf{W}} = 0 &\Rightarrow \mathbf{W}_{t+1} = \left( \left( \sigma_t^2 \mathbf{I}_h + \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \right)^{-1} \mathbf{W}_t^\top \mathbf{S} \right)^\top \\ &= \mathbf{S} \mathbf{W}_t \left( \sigma_t^2 \mathbf{I}_h + \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \right)^{-1}\end{aligned}$$

Now let's take the gradient w.r.t.  $\sigma$ .

$$\begin{aligned}\frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta})}{\partial \sigma} &= -\frac{nd}{\sigma} + \frac{n}{\sigma^3} \text{Tr}(\mathbf{S}) - \frac{\partial}{\partial \sigma} \left( \frac{n\sigma_t^2}{2\sigma^2} \text{Tr}(\mathbf{V} \mathbf{V}_t^{-1}) \right) \\ &\quad - \frac{\partial}{\partial \sigma} \left( \frac{n}{2\sigma^2} \text{Tr}(\mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{V} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S}) \right) - \frac{2n}{\sigma^3} \text{Tr}(\mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{S})\end{aligned}$$

Firstly,

$$\begin{aligned}\frac{\partial}{\partial \sigma} \left( \frac{n\sigma_t^2}{2\sigma^2} \text{Tr}(\mathbf{V} \mathbf{V}_t^{-1}) \right) &= \frac{\partial}{\partial \sigma} \left( \frac{n\sigma_t^2}{2\sigma^2} \text{Tr}((\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_h) \mathbf{V}_t^{-1}) \right) \\ &= \frac{\partial}{\partial \sigma} \left( \frac{n\sigma_t^2}{2\sigma^2} \text{Tr}(\mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1}) \right) \\ &= -\frac{n\sigma_t^2}{\sigma^3} \text{Tr}(\mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1})\end{aligned}$$

And then we compute

$$\begin{aligned}
& \frac{\partial}{\partial \sigma} \left( \frac{n}{2\sigma^2} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{V} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \right) \right) \\
&= \frac{\partial}{\partial \sigma} \left( \frac{n}{2\sigma^2} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} (\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_h) \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \right) \right) \\
&= \frac{\partial}{\partial \sigma} \left( \frac{n}{2\sigma^2} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \right) \right) \\
&= -\frac{n}{\sigma^3} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \right)
\end{aligned}$$

Using the properties of trace, we have

$$\begin{aligned}
\frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta})}{\partial \sigma} &= -\frac{nd}{\sigma} + \frac{n}{\sigma^3} \text{Tr}(\mathbf{S}) + \frac{n\sigma_t^2}{\sigma^3} \text{Tr} \left( \mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1} \right) \\
&\quad + \frac{n}{\sigma^3} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \right) - \frac{2n}{\sigma^3} \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{S} \right) \\
&= -\frac{n}{\sigma^3} \left( \sigma^2 d - \text{Tr}(\mathbf{S}) - \sigma_t^2 \text{Tr} \left( \mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1} \right) \right. \\
&\quad \left. - \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{W} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \right) + 2 \text{Tr} \left( \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \mathbf{S} \right) \right) \\
&= -\frac{n}{\sigma^3} \left( \sigma^2 d - \text{Tr}(\mathbf{S}) - \sigma_t^2 \text{Tr} \left( \mathbf{W} \mathbf{V}_t^{-1} \mathbf{W}^\top \right) \right. \\
&\quad \left. - \text{Tr} \left( \mathbf{W} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \right) + 2 \text{Tr} \left( \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \right) \right) \\
&= -\frac{n}{\sigma^3} \left( \sigma^2 d - \text{Tr}(\mathbf{S}) - \text{Tr} \left( \sigma_t^2 \mathbf{W} \mathbf{V}_t^{-1} \mathbf{W}^\top + \mathbf{W} \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \right) \right. \\
&\quad \left. + 2 \text{Tr} \left( \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \right) \right) \\
&= -\frac{n}{\sigma^3} \left( \sigma^2 d - \text{Tr}(\mathbf{S}) - \text{Tr} \left( \mathbf{W} \left( \sigma_t^2 \mathbf{I}_h + \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \right) \mathbf{V}_t^{-1} \mathbf{W}^\top \right) \right. \\
&\quad \left. + 2 \text{Tr} \left( \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \right) \right) \\
&= -\frac{n}{\sigma^3} \left( \sigma^2 d - \text{Tr}(\mathbf{S}) - \text{Tr} \left( \mathbf{W} \left( \sigma_t^2 \mathbf{I}_h + \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t \right) \mathbf{V}_t^{-1} \mathbf{W}^\top \right) \right. \\
&\quad \left. + 2 \text{Tr} \left( \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}^\top \right) \right)
\end{aligned}$$

Substituting  $\mathbf{W} = \mathbf{W}_{t+1} = \mathbf{S} \mathbf{W}_t (\sigma_t^2 \mathbf{I}_h + \mathbf{V}_t^{-1} \mathbf{W}_t^\top \mathbf{S} \mathbf{W}_t)^{-1}$ , we finally get the M-step update equation for  $\sigma_{t+1}$ .

$$\begin{aligned}
\frac{\partial \tilde{\mathcal{L}}(\boldsymbol{\theta})}{\partial \sigma} &= -\frac{n}{\sigma^3} \left( \sigma^2 d - \text{Tr}(\mathbf{S}) - \text{Tr} \left( \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}_{t+1}^\top \right) + 2 \text{Tr} \left( \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}_{t+1}^\top \right) \right) \\
&= -\frac{n}{\sigma^3} \left( \sigma^2 d - \text{Tr} \left( \mathbf{S} - \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}_{t+1}^\top \right) \right) = 0 \\
&\Rightarrow \sigma_{t+1}^2 = \frac{1}{d} \text{Tr} \left( \mathbf{S} - \mathbf{S} \mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}_{t+1}^\top \right)
\end{aligned}$$

2. (a) According to Problem 1-(a), we can directly get

$$\begin{aligned} p(\mathbf{z}|c, \mathbf{x}; \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{z} | \mathbf{V}_c^{-1} \mathbf{W}_c^\top (\mathbf{x} - \boldsymbol{\mu}_c), \sigma_c^2 \mathbf{V}_c^{-1}) \\ p(\mathbf{x}|c; \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_c, \mathbf{W}_c \mathbf{W}_c^\top + \sigma_c^2 \mathbf{I}_d) \end{aligned}$$

where  $\mathbf{V}_c = \mathbf{W}_c^\top \mathbf{W}_c + \sigma_c^2 \mathbf{I}_h$ .

So we can compute

$$\begin{aligned} p(\mathbf{x}; \boldsymbol{\theta}) &= \sum_{j=1}^k p(c = j, \mathbf{x}; \boldsymbol{\theta}) \\ &= \sum_{j=1}^k p(c = j; \boldsymbol{\theta}) p(\mathbf{x}|c = j; \boldsymbol{\theta}) \\ &= \sum_{j=1}^k \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \mathbf{W}_j \mathbf{W}_j^\top + \sigma_j^2 \mathbf{I}_d) \end{aligned}$$

Then we can derive

$$\begin{aligned} p(c|\mathbf{x}; \boldsymbol{\theta}) &= \frac{p(c, \mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{x}; \boldsymbol{\theta})} \\ &= \frac{\pi_c \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_c, \mathbf{W}_c \mathbf{W}_c^\top + \sigma_c^2 \mathbf{I}_d)}{\sum_{j=1}^k \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \mathbf{W}_j \mathbf{W}_j^\top + \sigma_j^2 \mathbf{I}_d)} \end{aligned}$$

- (b) Since  $\boldsymbol{\pi}$  is under the constraint  $\sum_{j=1}^k \pi_j = 1$ , we can introduce the Lagrangian multiplier

$$\mathcal{L}(\boldsymbol{\theta}, \lambda) = \mathcal{L}(\boldsymbol{\theta}) + \lambda \left( 1 - \sum_{j=1}^k \pi_j \right)$$

By taking the gradient w.r.t.  $\pi_j$ , we have

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta}, \lambda)}{\partial \pi_j} = \sum_{i=1}^n \frac{r_{t,i,j}}{\pi_j} - \lambda = 0 \Rightarrow \pi_{t+1,j} = \frac{1}{\lambda} \sum_{i=1}^n r_{t,i,j}$$

Using the constraint  $\sum_{j=1}^k \pi_j = 1$  again, we can finalize the M-step update for  $\boldsymbol{\pi}$ .

$$\begin{aligned} \sum_{j=1}^k \pi_{t+1,j} &= \frac{1}{\lambda} \sum_{i=1}^n \sum_{j=1}^k r_{t,i,j} \\ &= \frac{1}{\lambda} \sum_{i=1}^n 1 = \frac{n}{\lambda} = 1 \Rightarrow \lambda = n \\ \therefore \pi_{t+1,j} &= \frac{1}{n} \sum_{i=1}^n r_{t,i,j} \end{aligned}$$

By taking the gradient w.r.t.  $\boldsymbol{\mu}$ , we can get the M-step update for  $\boldsymbol{\mu}$ .

$$\begin{aligned}\frac{\partial \mathcal{L}(\boldsymbol{\theta}, \lambda)}{\partial \boldsymbol{\mu}_j} &= \sum_{i=1}^n r_{t,i,j} \frac{\partial}{\partial \boldsymbol{\mu}_j} \left( -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top (\mathbf{W}_j \mathbf{W}_j^\top + \sigma_j^2 \mathbf{I}_d)^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) \right) \\ &= \sum_{i=1}^n r_{t,i,j} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top (\mathbf{W}_j \mathbf{W}_j^\top + \sigma_j^2 \mathbf{I}_d)^{-1} = \mathbf{0} \\ \Rightarrow \boldsymbol{\mu}_{t+1,j} &= \frac{\sum_{i=1}^n r_{t,i,j} \mathbf{x}_i}{\sum_{i=1}^n r_{t,i,j}}\end{aligned}$$

(c) We can rewrite  $q(\mathbf{z}_i, c_i; \boldsymbol{\theta}_t)$  as follows.

$$\begin{aligned}q(\mathbf{z}_i, c_i; \boldsymbol{\theta}_t) &= p(\mathbf{z}_i | c_i, \mathbf{x}_i; \boldsymbol{\theta}_t) p(c_i | \mathbf{x}_i; \boldsymbol{\theta}_t) \\ &= r_{t,i,c_i} \mathcal{N}(\mathbf{z}_i | \mathbf{V}_{t,c_i}^{-1} \mathbf{W}_{t,c_i}^\top (\mathbf{x}_i - \boldsymbol{\mu}_{t,c_i}), \sigma_{t,c_i}^2 \mathbf{V}_{t,c_i}^{-1})\end{aligned}$$

Then, by using the results of Problem 1-(c), we have

$$\begin{aligned}\tilde{\mathcal{L}}(\boldsymbol{\theta}) &= \sum_{i=1}^n \mathbb{E}_q[\log p(\mathbf{x}_i, \mathbf{z}_i, c_i; \boldsymbol{\theta})] \\ &= \sum_{i=1}^n \sum_{j=1}^k r_{t,i,j} \mathbb{E}_{p(\mathbf{z}_i | c_i=j, \mathbf{x}_i; \boldsymbol{\theta}_t)}[\log p(\mathbf{x}_i, \mathbf{z}_i, c_i=j; \boldsymbol{\theta})] \\ &= \sum_{i=1}^n \sum_{j=1}^k r_{t,i,j} \mathbb{E}_{p(\mathbf{z}_i | c_i=j, \mathbf{x}_i; \boldsymbol{\theta}_t)}[\log p(\mathbf{x}_i, \mathbf{z}_i | c_i=j; \boldsymbol{\theta}) + \log p(c_i=j; \boldsymbol{\theta})] \\ &= \sum_{i=1}^n \sum_{j=1}^k r_{t,i,j} (\log \pi_j + \mathbb{E}_{p(\mathbf{z}_i | c_i=j, \mathbf{x}_i; \boldsymbol{\theta}_t)}[\log p(\mathbf{x}_i, \mathbf{z}_i | c_i=j; \boldsymbol{\theta})])\end{aligned}$$

And we already know that

$$\begin{aligned}&\mathbb{E}_{p(\mathbf{z}_i | c_i=j, \mathbf{x}_i; \boldsymbol{\theta}_t)}[\log p(\mathbf{x}_i, \mathbf{z}_i | c_i=j; \boldsymbol{\theta})] \\ &= -\frac{d+h}{2} \log 2\pi - d \log \sigma_j - \frac{1}{2\sigma_j^2} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top (\mathbf{x}_i - \boldsymbol{\mu}_j) - \frac{\sigma_{t,j}^2}{2\sigma_j^2} \text{Tr}(\mathbf{V}_j \mathbf{V}_{t,j}^{-1}) \\ &\quad - \frac{1}{2\sigma_j^2} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \mathbf{W}_{t,j} \mathbf{V}_{t,j}^{-1} \mathbf{V}_j \mathbf{V}_{t,j}^{-1} \mathbf{W}_{t,j}^\top (\mathbf{x}_i - \boldsymbol{\mu}_j) \\ &\quad + \frac{1}{2\sigma_j^2} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \mathbf{W}_{t,j} \mathbf{V}_{t,j}^{-1} \mathbf{W}_j^\top (\mathbf{x}_i - \boldsymbol{\mu}_j)\end{aligned}$$

where  $\mathbf{W}_j, \mathbf{W}_{t,j}, \mathbf{V}_j, \mathbf{V}_{t,j}$  are defined for  $c_i=j$  in the same manner as Problem 1. Let's define  $T_j$  and  $n_j$  as follows to write  $\tilde{\mathcal{L}}(\boldsymbol{\theta})$  concisely.

$$\begin{aligned}n_j &= \sum_{i=1}^n r_{t,i,j} \\ T_j &= \frac{1}{n_j} \sum_{i=1}^n r_{t,i,j} (\mathbf{x}_i - \boldsymbol{\mu}_j) (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top\end{aligned}$$

Then we can rewrite  $\tilde{\mathcal{L}}(\boldsymbol{\theta})$

$$\begin{aligned}
\tilde{\mathcal{L}}(\boldsymbol{\theta}) &= \sum_{i=1}^n \sum_{j=1}^k r_{t,i,j} \left( \log \pi_j + \mathbb{E}_{p(\mathbf{z}_i|c_i=j, \mathbf{x}_i; \boldsymbol{\theta}_t)} [\log p(\mathbf{x}_i, \mathbf{z}_i|c_i=j; \boldsymbol{\theta})] \right) \\
&= \sum_{i=1}^n \sum_{j=1}^k r_{t,i,j} \left( \log \pi_j - \frac{d+h}{2} \log 2\pi - d \log \sigma_j \right) \\
&\quad - \frac{n_j}{2\sigma_j^2} \sum_{j=1}^k \left( \text{Tr}(\mathbf{T}_j) + \sigma_{t,j}^2 \text{Tr}(\mathbf{V}_j \mathbf{V}_{t,j}^{-1}) + \text{Tr}(\mathbf{W}_{t,j} \mathbf{V}_{t,j}^{-1} \mathbf{V}_j \mathbf{V}_{t,j}^{-1} \mathbf{W}_{t,j}^\top \mathbf{T}_j) \right. \\
&\quad \left. - 2 \text{Tr}(\mathbf{W}_{t,j} \mathbf{V}_{t,j}^{-1} \mathbf{W}_j^\top \mathbf{T}_j) \right) \\
&= - \sum_{j=1}^k n_j d \log \sigma_j \\
&\quad - \frac{n_j}{2\sigma_j^2} \sum_{j=1}^k \left( \text{Tr}(\mathbf{T}_j) + \sigma_t^2 \text{Tr}(\mathbf{V}_j \mathbf{V}_{t,j}^{-1}) + \text{Tr}(\mathbf{W}_{t,j} \mathbf{V}_{t,j}^{-1} \mathbf{V}_j \mathbf{V}_{t,j}^{-1} \mathbf{W}_{t,j}^\top \mathbf{T}_j) \right. \\
&\quad \left. - 2 \text{Tr}(\mathbf{W}_{t,j} \mathbf{V}_{t,j}^{-1} \mathbf{W}_j^\top \mathbf{T}_j) \right) + \text{const.}
\end{aligned}$$

which is almost same to the one in Problem 1.

Thus, by taking the gradient as in Problem 1, the M-step equations are

$$\begin{aligned}
\mathbf{W}_{t+1,j} &= \mathbf{T}_j \mathbf{W}_{t,j} \left( \sigma_{t,j}^2 \mathbf{I}_h + \mathbf{V}_{t,j}^{-1} \mathbf{W}_{t,j}^\top \mathbf{T}_j \mathbf{W}_{t,j} \right)^{-1} \\
\sigma_{t+1,j}^2 &= \frac{1}{d} \text{Tr} \left( \mathbf{T}_j - \mathbf{T}_j \mathbf{W}_{t,j} \mathbf{V}_{t,j}^{-1} \mathbf{W}_{t+1,j}^\top \right)
\end{aligned}$$