

Student Name: Suryateja B.V.

Roll Number: 160729

Date: March 14, 2019

To solve part 1, we first need to compute $p(\mathbf{t}|\mathbf{X}, \mathbf{f}, \mathbf{Z})$.

$$p(\mathbf{t}|\mathbf{X}, \mathbf{f}, \mathbf{Z}) \propto p(\mathbf{f}|\mathbf{X}, \mathbf{t}, \mathbf{Z}) p(\mathbf{t}|\mathbf{Z}) = \mathcal{N}(\mathbf{f}|\bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} \mathbf{t}, \mathbf{P}) \mathcal{N}(\mathbf{t}|\mathbf{0}, \bar{\mathbf{K}})$$

where $\bar{\mathbf{k}} = [\bar{\mathbf{k}}_1 \bar{\mathbf{k}}_2 \dots \bar{\mathbf{k}}_N]$ is an $M \times N$ matrix and \mathbf{P} is a diagonal matrix with each diagonal entry given by $p_{ii} = \kappa(\mathbf{x}_i, \mathbf{x}_i) - \bar{\mathbf{k}}_i^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_i$. Consider the exponent in the RHS side. We have,

$$\left((\mathbf{f} - \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} \mathbf{t})^T \mathbf{P}^{-1} (\mathbf{f} - \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} \mathbf{t}) + \mathbf{t}^T \bar{\mathbf{K}}^{-1} \mathbf{t} \right) = \begin{bmatrix} \mathbf{t} \\ \mathbf{f} \end{bmatrix}^T \begin{bmatrix} \bar{\mathbf{K}}^{-1} + \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} \mathbf{P}^{-1} \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} & -\bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} \mathbf{P}^{-1} \\ -\mathbf{P}^{-1} \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} & \mathbf{P}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{t} \\ \mathbf{f} \end{bmatrix}$$

We know $p(\mathbf{t}|\mathbf{X}, \mathbf{f}, \mathbf{Z}) = \mathcal{N}(\mathbf{t}|\mu_{\mathbf{t}|\mathbf{f}}, \Sigma_{\mathbf{t}|\mathbf{f}})$. Using Gaussian conditional formulae, we get

$$\begin{aligned} \Sigma_{\mathbf{t}|\mathbf{f}} &= \Lambda_{\mathbf{tt}}^{-1} = \left(\bar{\mathbf{K}}^{-1} + \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} \mathbf{P}^{-1} \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} \right)^{-1} \\ &= \bar{\mathbf{K}} - \bar{\mathbf{k}} \left(\mathbf{P} + \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} \right)^{-1} \bar{\mathbf{k}}^T = \bar{\mathbf{K}} - \bar{\mathbf{k}} \mathbf{Q}^{-1} \bar{\mathbf{k}}^T \quad (\text{using Woodbury Inverse}) \end{aligned}$$

Note that the variance in \mathbf{t} has reduced after conditioning on \mathbf{f} . The matrix \mathbf{Q} is of the form,

$$\mathbf{Q} = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \bar{\mathbf{k}}_1^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_2 & \dots & \bar{\mathbf{k}}_1^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_N \\ \bar{\mathbf{k}}_2^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_1 & \kappa(\mathbf{x}_2, \mathbf{x}_2) & \dots & \bar{\mathbf{k}}_2^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_N \\ \vdots & \vdots & \ddots & \vdots \\ \bar{\mathbf{k}}_N^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_1 & \bar{\mathbf{k}}_N^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_2 & \dots & \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

We'd prefer the initial form for $\Sigma_{\mathbf{t}|\mathbf{f}}$ which involves an inverse of $M \times M$ matrix instead of the form obtained after using the Woodbury identity, since it involves inverse of \mathbf{Q} which is $N \times N$. The time to compute the initial form is of $\mathcal{O}(M^2 N)$ owing to the computation of $\bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} \mathbf{P}^{-1} \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1}$. Note that computing \mathbf{P}^{-1} is $\mathcal{O}(N)$ since \mathbf{P} is a diagonal matrix. Computing $\bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}$ and its transpose take $\mathcal{O}(M^2 N)$, and then multiplying both along with diagonal matrix \mathbf{P}^{-1} also takes $\mathcal{O}(M^2 N)$.

We can write $\mathbf{f} = \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} \mathbf{t} + \epsilon$, where $\epsilon \sim \mathcal{N}(\epsilon|\mathbf{0}, \mathbf{P})$. Then, $p(\mathbf{f}|\mathbf{X}, \mathbf{Z}) = \mathcal{N}(\mathbf{f}|\mu_{\mathbf{f}}, \Sigma_{\mathbf{f}})$ where $\mu_{\mathbf{f}} = \mathbb{E}[\mathbf{f}] = \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} (\mathbf{0}) + \mathbf{0} = \mathbf{0}$ and $\Sigma_{\mathbf{f}} = \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} + \mathbf{P} = \mathbf{Q}$. Using $\mu_{\mathbf{f}}$, we can now compute $\mu_{\mathbf{t}|\mathbf{f}}$ as follows :

$$\mu_{\mathbf{t}|\mathbf{f}} = \mu_{\mathbf{t}} - \Lambda_{\mathbf{tt}}^{-1} \Lambda_{\mathbf{tf}} (\mathbf{f} - \mu_{\mathbf{f}}) = \Sigma_{\mathbf{t}|\mathbf{f}} \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} \mathbf{P}^{-1} \mathbf{f}$$

We know $p(y_*|\mathbf{x}_*, \mathbf{f}, \mathbf{X}, \mathbf{Z}) = \int p(y_*|\mathbf{x}_*, \mathbf{f}, \mathbf{X}, \mathbf{Z}, \mathbf{t}) p(\mathbf{t}|\mathbf{f}, \mathbf{X}, \mathbf{Z}) d\mathbf{t}$.

Now, we can write $y_* = \mathbf{f}_* = \bar{\mathbf{k}}_*^T \bar{\mathbf{K}}^{-1} \mathbf{t} + \epsilon$, with $\mathbf{t} \sim \mathcal{N}(\mu_{\mathbf{t}|\mathbf{f}}, \Sigma_{\mathbf{t}|\mathbf{f}})$ and $\epsilon \sim \mathcal{N}(\epsilon|\mathbf{0}, \kappa(\mathbf{x}_*, \mathbf{x}_*) - \bar{\mathbf{k}}_*^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_*)$. Then, $p(y_*|\mathbf{x}_*, \mathbf{f}, \mathbf{X}, \mathbf{Z}) = \mathcal{N}(y_*|\mu_*, \Sigma_*)$ where

$$\begin{aligned} \mu_* &= \bar{\mathbf{k}}_*^T \bar{\mathbf{K}}^{-1} \mu_{\mathbf{t}|\mathbf{f}} = \bar{\mathbf{k}}_*^T \bar{\mathbf{K}}^{-1} \Sigma_{\mathbf{t}|\mathbf{f}} \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} \mathbf{P}^{-1} \mathbf{f} = \bar{\mathbf{k}}_*^T \bar{\mathbf{K}}^{-1} \left(\bar{\mathbf{K}}^{-1} + \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} \mathbf{P}^{-1} \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} \right)^{-1} \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} \mathbf{P}^{-1} \mathbf{f} \\ \Sigma_* &= \bar{\mathbf{k}}_*^T \bar{\mathbf{K}}^{-1} \Sigma_{\mathbf{t}|\mathbf{f}} \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_* + \kappa(\mathbf{x}_*, \mathbf{x}_*) - \bar{\mathbf{k}}_*^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}}_* \end{aligned}$$

Before, the posterior predictive involved the inversion of \mathbf{K} , which takes $\mathcal{O}(N^3)$ time complexity. With pseudo training data, we see that computation of $\Sigma_{\mathbf{t}|\mathbf{f}}$ takes $\mathcal{O}(M^2N)$, a significant reduction in time taken.

For part (2), we derive the MLE-II solution, that is, maximise the marginal likelihood $p(\mathbf{f}|\mathbf{X}, \mathbf{Z})$. Earlier, we have derived the same - $p(\mathbf{f}|\mathbf{X}, \mathbf{Z}) = \mathcal{N}(\mathbf{f}|\mu_{\mathbf{f}}, \Sigma_{\mathbf{f}})$ where $\mu_{\mathbf{f}} = \mathbf{0}$ and $\Sigma_{\mathbf{f}} = \bar{\mathbf{k}}^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{k}} + \mathbf{P} = \mathbf{Q}$. The MLE-II objective is given by,

$$\arg \max_{\mathbf{Z}} \log(p(\mathbf{f}|\mathbf{X}, \mathbf{Z})) = \arg \min_{\mathbf{Z}} (\log |\Sigma_{\mathbf{f}}| + \mathbf{f}^T \Sigma_{\mathbf{f}}^{-1} \mathbf{f})$$

There is dependence on \mathbf{Z} through \mathbf{Q} , as in, $\bar{\mathbf{k}}_i = [\kappa(\mathbf{x}_i, \mathbf{z}_1) \kappa(\mathbf{x}_i, \mathbf{z}_2) \cdots \kappa(\mathbf{x}_i, \mathbf{z}_M)]$

Student Name: Suryateja B.V.

Roll Number: 160729

Date: March 14, 2019

0.1 EM - I : not involving \mathbf{z}_n

We first marginalize over \mathbf{z}_n . Let us compute $p(\mathbf{x}_n | c_n = m, \Theta)$. Given $c_n = m$, we can write \mathbf{x}_n as $\mathbf{x}_n = \boldsymbol{\mu}_m + \mathbf{W}_m \mathbf{z}_n + \boldsymbol{\epsilon}_n$ where $\mathbf{z}_n \sim \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I}_K)$ and $\boldsymbol{\epsilon}_n \sim \mathcal{N}(\boldsymbol{\epsilon}_n | \mathbf{0}, \sigma_m^2 \mathbf{I}_D)$

$$\mathbb{E}[\mathbf{x}_n] = \boldsymbol{\mu}_m \quad \text{and} \quad V(\mathbf{x}_n) = \mathbf{W}_m \mathbf{W}_m^T + \sigma_m^2 \mathbf{I}_D = \mathbf{V}_m$$

So, $p(\mathbf{x}_n | c_n = m, \Theta) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_m, \mathbf{W}_m \mathbf{W}_m^T + \sigma_m^2 \mathbf{I}_D)$

Now let us compute the posterior probabilities of latent variables.

$$\begin{aligned} p(c_n = m | \mathbf{x}_n, \Theta) &\propto p(\mathbf{x}_n | c_n = m, \Theta) p(c_n = m) = \pi_m \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_m, \mathbf{V}_m) \\ \implies p(c_n = m | \mathbf{x}_n, \Theta) &= \frac{\pi_m \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_m, \mathbf{V}_m)}{\sum_{l=1}^M \pi_l \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_l, \mathbf{V}_l)} = r_{nm} \end{aligned}$$

The CLL is given as follows :

$$\begin{aligned} p(\mathbf{x}, \mathbf{c} | \Theta) &= \prod_{n=1}^N \prod_{m=1}^M (p(\mathbf{x}_n | c_n = m, \Theta) p(c_n = m | \Theta))^{\mathbb{I}[c_n=m]} \\ \implies \log p(\mathbf{x}, \mathbf{c} | \Theta) &= \sum_{n=1}^N \sum_{m=1}^M \mathbb{I}[c_n = m] \left(\log \pi_m - \frac{1}{2} \log |\mathbf{V}_m| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_m)^T \mathbf{V}_m^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_m) \right) \end{aligned}$$

The only expectation we need is $\mathbb{E}[\mathbb{I}[c_n = m]] = r_{nm}$. So the expected CLL is given by,

$$\sum_{n=1}^N \sum_{m=1}^M r_{nm} \left(\log \pi_m - \frac{1}{2} \log |\mathbf{V}_m| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_m)^T \mathbf{V}_m^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_m) \right) \quad (1)$$

The M-step update equations are as follows :

$$\begin{aligned} \hat{\pi}_m &= \frac{\sum_{n=1}^N r_{nm}}{N} = \frac{N_m}{N} \\ \hat{\boldsymbol{\mu}}_m &= \frac{\sum_{n=1}^N r_{nm} \mathbf{x}_n}{N_m} \\ \hat{\mathbf{V}}_m &= \frac{\sum_{n=1}^N r_{nm} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m) (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m)^T}{N_m} \end{aligned}$$

$$\mathbf{W}_m \mathbf{W}_m^T + \sigma_m^2 \mathbf{I}_D = \hat{\mathbf{V}}_m \implies \hat{\mathbf{W}}_m = \mathbf{U}_K (\mathbf{L}_K - \hat{\sigma}_m^2 \mathbf{I}_K)^{1/2} \mathbf{R} \quad \text{and} \quad \hat{\sigma}_m^2 = \frac{1}{D-K} \sum_{k=K+1}^D \lambda_k$$

where \mathbf{U}_K is a $D \times K$ matrix of top K eigen vectors of $\hat{\mathbf{V}}_m$, $\mathbf{L}_K : K \times K$ diagonal matrix of top K eigen values $\lambda_1, \lambda_2, \dots, \lambda_K$, \mathbf{R} is a $K \times K$ rotation matrix. While this method avoids z_n estimates, it is expensive due to eigen decomposition.

1 Initialize $\Theta = \{\pi_m, \boldsymbol{\mu}_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M = \Theta^{(0)}$. Set $t = 1$;

2 **E-Step**

$$r_{nm}^{(t)} = \frac{\pi_m^{(t-1)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_m^{(t-1)}, \mathbf{V}_m^{(t-1)})}{\sum_{l=1}^M \pi_m^{(t-1)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_l^{(t-1)}, \mathbf{V}_l^{(t-1)})} \quad \forall n, m$$

;

3 **M-Step** - Do for all m

$$\hat{\pi}_m^{(t)} = \frac{\sum_{n=1}^N r_{nm}^{(t)}}{N} = \frac{N_m^{(t)}}{N}$$

$$\hat{\boldsymbol{\mu}}_m^{(t)} = \frac{\sum_{n=1}^N r_{nm}^{(t)} \mathbf{x}_n}{N_m^{(t)}}$$

$$\hat{\mathbf{V}}_m^{(t)} = \frac{\sum_{n=1}^N r_{nm}^{(t)} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m^{(t)}) (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m^{(t)})^T}{N_m^{(t)}}$$

$$(\hat{\sigma}_m^2)^{(t)} = \frac{1}{D - K} \sum_{k=K+1}^D \lambda_k^{(t)}$$

$$\hat{\mathbf{W}}_m^{(t)} = \mathbf{U}_K^{(t)} \left(\mathbf{L}_K^{(t)} - (\hat{\sigma}_m^2)^{(t)} \mathbf{I}_K \right)^{1/2} \mathbf{R}^{(t)}$$

$$t = t + 1$$

;

4 Go to E-Step if not converged.

The stepwise online algorithm sketch is as follows :

1 Initialize $\Theta = \{\pi_m, \boldsymbol{\mu}_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M = \Theta^{(0)}$. Set $t = 1$;

2 Pick a random example \mathbf{x}_n ;

3 Compute $r_{nm}^{(t)}$ for all m ;

4 Compute learning rate ϵ_t ;

5 Compute $\hat{\Theta}$ using only example \mathbf{x}_n ;

6 $\Theta^{(t)} = (1 - \epsilon_t) \Theta^{(t-1)} + \epsilon_t \hat{\Theta}$;

7 Go to Step-2 if Θ not converged;

0.2 EM - II : Include Estimating \mathbf{z}_n

Conditional posterior of c_n is same as before (since it is independent of \mathbf{z}_n), and is given by:

$$p(c_n = m | \mathbf{x}_n, \mathbf{z}_n, \Theta) = \frac{\pi_m \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_m, \mathbf{V}_m)}{\sum_{l=1}^M \pi_m \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_l, \mathbf{V}_l)} = r_{nm}$$

We sort of have a mixture of M PPCA models. Each model coming from a particular class m . For a single PPCA model we had a latent variable \mathbf{z}_n for each \mathbf{x}_n . Now, in case of mixture of PPCA, we have M latent variables $\{\mathbf{z}_{nm}\}_{m=1}^M$ for each \mathbf{x}_n . Conditional posterior of \mathbf{z}_{nm} is as follows:

$$p(\mathbf{z}_{nm} | \mathbf{x}_n, c_n = m, \Theta) \propto p(\mathbf{x}_n | \mathbf{z}_{nm}, c_n = m, \Theta) p(\mathbf{z}_{nm} | \Theta) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_m + \mathbf{W}_m \mathbf{z}_{nm}, \sigma_m^2 \mathbf{I}_D) \mathcal{N}(\mathbf{z}_{nm} | 0, \mathbf{I}_K)$$

Using Gaussian conditional properties, we get $p(\mathbf{z}_{nm} | \mathbf{x}_n, c_n = m, \Theta) = \mathcal{N}(\mathbf{z}_{nm} | \boldsymbol{\mu}_{nm}, \boldsymbol{\Sigma}_{nm})$ where $\boldsymbol{\Sigma}_{nm} = \sigma_m^2 (\mathbf{W}_m^T \mathbf{W}_m + \sigma_m^2 \mathbf{I}_K)^{-1} = \sigma_m^2 \mathbf{M}_m^{-1}$ and $\boldsymbol{\mu}_{nm} = \mathbf{M}_m^{-1} \mathbf{W}_m^T (\mathbf{x}_n - \boldsymbol{\mu}_m)$. The joint posterior distribution is as follows :

$$p(\mathbf{z}_{nm}, c_n = m | \mathbf{x}_n, \Theta) = p(\mathbf{z}_{nm} | \mathbf{x}_n, c_n = m, \Theta) p(c_n = m | \mathbf{x}_n, \Theta) = r_{nm} \mathcal{N}(\mathbf{z}_{nm} | \boldsymbol{\mu}_{nm}, \boldsymbol{\Sigma}_{nm})$$

Note that \mathbf{z}_{nm} inherently holds the information that it belongs to m^{th} class, so we can say $p(\mathbf{z}_{nm} | \mathbf{x}_n, c_n = m, \Theta) = p(\mathbf{z}_{nm} | \mathbf{x}_n, \Theta)$. This means that the joint posterior contains two independent terms. Hence when we take expectations wrt joint posterior (as in M-step), we can decouple and take expectation wrt individual terms $p(\mathbf{z}_{nm} | \mathbf{x}_n, \Theta)$ and $p(c_n = m | \mathbf{x}_n, \Theta)$.

The CLL is given as follows :

$$\begin{aligned} p(\mathbf{x}, \mathbf{c}, \mathbf{z} | \Theta) &= \prod_{n=1}^N \prod_{m=1}^M (p(\mathbf{x}_n | \mathbf{z}_n, c_n = m, \Theta) p(\mathbf{z}_{nm} | c_n = m, \Theta) p(c_n = m | \Theta))^{\mathbb{I}[c_n=m]} \\ \log p(\mathbf{x}, \mathbf{c} | \Theta) &= \sum_{n=1}^N \sum_{m=1}^M \mathbb{I}[c_n = m] \left(\log \pi_m - \frac{D}{2} \log \sigma_m^2 - \frac{1}{2\sigma_m^2} (\mathbf{x}_n - \boldsymbol{\mu}_m - \mathbf{W}_m \mathbf{z}_{nm})^T (\mathbf{x}_n - \boldsymbol{\mu}_m - \mathbf{W}_m \mathbf{z}_{nm}) - \frac{1}{2} \mathbf{z}_{nm}^T \mathbf{z}_{nm} \right) \\ &= \sum_{n=1}^N \sum_{m=1}^M \mathbb{I}[c_n = m] \left(\log \pi_m - \frac{D}{2} \log \sigma_m^2 - \frac{1}{2\sigma_m^2} \|\mathbf{x}_n - \boldsymbol{\mu}_m\|^2 - \frac{1}{2\sigma_m^2} \text{tr}(\mathbf{z}_{nm} \mathbf{z}_{nm}^T \mathbf{W}_m^T \mathbf{W}_m) \right. \\ &\quad \left. + \frac{2}{2\sigma_m^2} (\mathbf{x}_n - \boldsymbol{\mu}_m)^T \mathbf{W}_m \mathbf{z}_{nm} - \frac{1}{2} \text{tr}(\mathbf{z}_{nm} \mathbf{z}_{nm}^T) \right) \\ \text{ECLL} &= \sum_{n=1}^N \sum_{m=1}^M \mathbb{E}[\mathbb{I}[c_n = m]] \left(\log \pi_m - \frac{D}{2} \log \sigma_m^2 - \frac{1}{2\sigma_m^2} \|\mathbf{x}_n - \boldsymbol{\mu}_m\|^2 - \frac{1}{2\sigma_m^2} \text{tr}(\mathbb{E}[\mathbf{z}_{nm} \mathbf{z}_{nm}^T] \mathbf{W}_m^T \mathbf{W}_m) \right. \\ &\quad \left. + \frac{2}{2\sigma_m^2} (\mathbf{x}_n - \boldsymbol{\mu}_m)^T \mathbf{W}_m \mathbb{E}[\mathbf{z}_{nm}] - \frac{1}{2} \text{tr}(\mathbb{E}[\mathbf{z}_{nm} \mathbf{z}_{nm}^T]) \right) \\ \text{where } \mathbb{E}[\mathbf{z}_{nm}] &= \boldsymbol{\mu}_{nm} \quad \text{and} \quad \mathbb{E}[\mathbf{z}_{nm} \mathbf{z}_{nm}^T] = \boldsymbol{\mu}_{nm} \boldsymbol{\mu}_{nm}^T + \boldsymbol{\Sigma}_{nm} \quad \text{and} \quad \mathbb{E}[\mathbb{I}[c_n = m]] = r_{nm} \end{aligned}$$

The M-Step updates are given as follows :

$$\begin{aligned}\hat{\pi}_m &= \frac{\sum_{n=1}^N r_{nm}}{N} = \frac{N_m}{N} \\ \hat{\boldsymbol{\mu}}_m &= \frac{\sum_{n=1}^N r_{nm} (\mathbf{x}_n - \mathbf{W}_m \mathbb{E}[\mathbf{z}_{nm}])}{N_m} \\ \hat{\mathbf{W}}_m &= \left(\sum_{n=1}^N r_{nm} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m) \mathbb{E}[\mathbf{z}_{nm}]^T \right) \left(\sum_{n=1}^N r_{nm} \mathbb{E}[\mathbf{z}_{nm} \mathbf{z}_{nm}^T] \right)^{-1} \\ \hat{\sigma}_m^2 &= \frac{1}{DN_m} \left(\sum_{n=1}^N \left(\|\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m\|^2 - 2 \mathbb{E}[\mathbf{z}_{nm}]^T \hat{\mathbf{W}}_m^T (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m) + \text{tr} \left(\mathbb{E}[\mathbf{z}_{nm} \mathbf{z}_{nm}^T] \hat{\mathbf{W}}_m^T \hat{\mathbf{W}}_m \right) \right) \right)\end{aligned}$$

1 Initialize $\Theta = \{\pi_m, \boldsymbol{\mu}_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M = \Theta^{(0)}$. Set $t = 1$;

2 E-Step $\forall n, m$

$$\begin{aligned}r_{nm}^{(t)} &= \frac{\pi_m^{(t-1)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_m^{(t-1)}, \mathbf{V}_m^{(t-1)})}{\sum_{l=1}^M \pi_m^{(t-1)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_l^{(t-1)}, \mathbf{V}_l^{(t-1)})} \\ \mathbf{M}_m^{(t)} &= (\mathbf{W}_m^T)^{(t-1)} \mathbf{W}_m^{(t-1)} + (\sigma_m^2)^{(t-1)} \mathbf{I}_K \\ \mathbb{E}[\mathbf{z}_{nm}]^{(t)} &= (\mathbf{M}_m^{-1})^{(t)} (\mathbf{W}_m^T)^{(t-1)} (\mathbf{x}_n - \boldsymbol{\mu}_m) \\ \mathbb{E}[\mathbf{z}_{nm} \mathbf{z}_{nm}^T]^{(t)} &= \mathbb{E}[\mathbf{z}_{nm}]^{(t-1)} (\mathbb{E}[\mathbf{z}_{nm}]^T)^{(t)} + (\sigma_m^2)^{(t-1)} (\mathbf{M}_m^{-1})^{(t)}\end{aligned}$$

;

3 M-Step - Do for all m

$$\begin{aligned}\hat{\pi}_m^{(t)} &= \frac{\sum_{n=1}^N r_{nm}^{(t)}}{N} = \frac{N_m^{(t)}}{N} \\ \hat{\boldsymbol{\mu}}_m^{(t)} &= \frac{\sum_{n=1}^N r_{nm}^{(t)} (\mathbf{x}_n - \mathbf{W}_m^{(t-1)} \mathbb{E}[\mathbf{z}_{nm}]^{(t)})}{N_m^{(t)}} \\ \hat{\mathbf{W}}_m^{(t)} &= \left(\sum_{n=1}^N r_{nm}^{(t)} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m^{(t)}) (\mathbb{E}[\mathbf{z}_{nm}]^T)^{(t)} \right) \left(\sum_{n=1}^N r_{nm}^{(t)} (\mathbb{E}[\mathbf{z}_{nm} \mathbf{z}_{nm}^T])^{(t)} \right)^{-1} \\ (\hat{\sigma}_m^2)^{(t)} &= \frac{1}{DN_m^{(t)}} \sum_{n=1}^N r_{nm}^{(t)} \left(\|\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m^{(t)}\|^2 - 2 (\mathbb{E}[\mathbf{z}_{nm}]^T)^{(t)} (\hat{\mathbf{W}}_m^T)^{(t)} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_m^{(t)}) \right. \\ &\quad \left. + \text{tr} \left((\mathbb{E}[\mathbf{z}_{nm} \mathbf{z}_{nm}^T])^{(t)} (\hat{\mathbf{W}}_m^T \hat{\mathbf{W}}_m)^{(t)} \right) \right)\end{aligned}$$

$$t = t + 1$$

;

4 Go to E-Step if not converged.

The stepwise online algorithm sketch is as follows :

- 1 Initialize $\Theta = \{\pi_m, \boldsymbol{\mu}_m, \mathbf{W}_m, \sigma_m^2\}_{m=1}^M = \Theta^{(0)}$. Set $t = 1$;
- 2 Pick a random example \mathbf{x}_n ;
- 3 Compute $r_{nm}^{(t)}, \mathbf{M}_m^{(t)}, (\mathbb{E}[\mathbf{z}_{nm}])^{(t)}, (\mathbb{E}[\mathbf{z}_{nm}\mathbf{z}_{nm}^T])^{(t)}$ for all m ;
- 4 Compute learning rate ϵ_t ;
- 5 Compute $\hat{\Theta}$ using only example \mathbf{x}_n ;
- 6 $\Theta^{(t)} = (1 - \epsilon_t) \Theta^{(t-1)} + \epsilon_t \hat{\Theta}$;
- 7 Go to Step-2 if Θ not converged;

Student Name: Suryateja B.V.

Roll Number: 160729

Date: March 14, 2019

Let the mean field approximation of posterior be given by $q(\mathbf{w}, \beta, \alpha) = q(\mathbf{w}|\phi_1) q(\beta|\phi_2) \prod_{d=1}^D q(\alpha_d|\phi_d)$. The complete data log likelihood is given by,

$$\begin{aligned} \log p(\mathbf{y}, \mathbf{w}, \beta, \alpha | \mathbf{X}, \theta) &= \log p(\mathbf{y} | \mathbf{w}, \beta, \alpha, \mathbf{X}, \theta) + \log p(\mathbf{w} | \alpha, \theta) + \log p(\beta | \theta) + \sum_{d=1}^D \log p(\alpha_d | \theta) \\ &= \log \beta \left(\frac{N}{2} + a_0 - 1 \right) - \frac{\beta}{2} \left((\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + 2b_0 \right) \\ &\quad - \frac{1}{2} \sum_{d=1}^D w_d^2 \alpha_d + \left(e_0 - \frac{1}{2} \right) \sum_{d=1}^D \log \alpha_d - f_0 \sum_{d=1}^D \alpha_d + \text{const.} \end{aligned}$$

where $\theta = \{a_0, b_0, e_0, f_0\}$, $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_D]$, $\mathbf{y} = [y_1, y_2, \dots, y_N]$ and \mathbf{X} is the matrix of $\{\mathbf{x}_n\}_{n=1}^N$. Now, using the mean-field VI algorithm, the optimal distributions are given as follows.

$$\begin{aligned} \log \hat{q}(\mathbf{w}) &= \mathbb{E}_{q(\alpha)q(\beta)} [\log p(\mathbf{y}, \mathbf{w}, \beta, \alpha | \mathbf{X}, \theta)] \\ &= -\frac{\mathbb{E}[\beta]}{2} \left(\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{w}^T \left(\mathbf{X}^T \mathbf{X} + \frac{1}{\mathbb{E}[\beta]} \mathbb{E}_\alpha [\text{diag}(\alpha_1, \dots, \alpha_D)] \right) \mathbf{w} + 2b_0 \right) \\ &\quad + \text{constant terms} \end{aligned}$$

Comparing the above with the log of normal distribution $\mathcal{N}(\mathbf{w} | \mu_N, \Sigma_N)$, we get,

$$\begin{aligned} \Sigma_N &= (\mathbb{E}[\beta] \mathbf{X}^T \mathbf{X} + \mathbb{E}_\alpha [\text{diag}(\alpha_1, \dots, \alpha_D)])^{-1} \\ \mu_N &= \Sigma_N \mathbb{E}[\beta] \mathbf{X}^T \mathbf{y} \end{aligned}$$

Now, let us solve for β .

$$\begin{aligned} \log \hat{q}(\beta) &= \mathbb{E}_{q(\alpha)q(\mathbf{w})} [\log p(\mathbf{y}, \mathbf{w}, \beta, \alpha | \mathbf{X}, \theta)] \\ &= \log \beta \left(\frac{N}{2} + a_0 - 1 \right) - \frac{\beta}{2} (\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbb{E}[\mathbf{w}] + \text{tr}(\mathbf{X}^T \mathbf{X} \mathbb{E}[\mathbf{w}\mathbf{w}^T]) + 2b_0) + \text{const.} \end{aligned}$$

Comparing the above with the log of Gamma($\beta | a_N, b_N$), we get,

$$\begin{aligned} a_N &= a_0 + \frac{N}{2} \\ b_N &= b_0 + \frac{1}{2} (\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbb{E}[\mathbf{w}] + \text{tr}(\mathbf{X}^T \mathbf{X} \mathbb{E}[\mathbf{w}\mathbf{w}^T])) \end{aligned}$$

Next, we solve for α .

$$\begin{aligned} \log \hat{q}(\alpha) &= \mathbb{E}_{q(\beta)q(\mathbf{w})} [\log p(\mathbf{y}, \mathbf{w}, \beta, \alpha | \mathbf{X}, \theta)] \\ &= -\frac{1}{2} \sum_{d=1}^D \mathbb{E}[w_d^2] \alpha_d + \left(e_0 - \frac{1}{2} \right) \sum_{d=1}^D \log \alpha_d - f_0 \sum_{d=1}^D \alpha_d + \text{const.} \end{aligned}$$

Comparing the above with $\prod_{d=1}^D \text{Gamma}(\alpha_d | e_{Nd}, f_{Nd})$, we get

$$\begin{aligned} e_{Nd} &= e_0 + \frac{1}{2} \quad \forall d \\ f_{Nd} &= f_0 + \frac{1}{2} \mathbb{E}[w_d^2] \end{aligned}$$

Now, let us write down the expectations.

$$\begin{aligned} \mathbb{E}[\mathbf{w}] &= \mu_N \\ \mathbb{E}[\mathbf{w}\mathbf{w}^T] &= \Sigma_N + \mu_N \mu_N^T \\ \mathbb{E}[\beta] &= \frac{a_N}{b_N} \\ \mathbb{E}[\alpha_d] &= \frac{e_{Nd}}{f_{Nd}} \\ \mathbb{E}[w_d^2] &= (\Sigma_N)_{dd} + \mu_{Nd}^2 \end{aligned}$$

Note that the updates of each of the optimum distributions depend on other distributions. So, we are required to perform cyclic updates. The algorithm is as follows: \mathbf{P} and \mathbf{K} are introduced

- 1 Given : $\mathbf{X}, \mathbf{y}, a_0, b_0, e_0, f_0$;
- 2 Set $e_{Nd} = e_0 + \frac{1}{2} \quad \forall d$. Set $a_N = a_0 + \frac{N}{2}$;
- 3 Set $\mathbf{K} = \mathbf{X}^T \mathbf{X}$. Set $\mathbf{P} = \mathbf{X}^T \mathbf{y}$. Set $\mathbf{Q} = \mathbf{y}^T \mathbf{y}$;
- 4 Initialize $b_N^{(0)} = b_0 + \frac{1}{2} (\mathbf{y}^T \mathbf{y} + \text{tr}(\mathbf{K}))$;
- 5 Initialize $f_{Nd}^{(0)} = f_0 + \frac{1}{2} \quad \forall d$;
- 6 Initialize $\mathbb{E}[\beta]^{(0)} = \frac{a_N^{(0)}}{b_N^{(0)}}$;
- 7 Initialize $\mathbb{E}[\alpha_d]^{(0)} = \frac{e_{Nd}^{(0)}}{f_{Nd}^{(0)}} \quad \forall d$;
- 8 Set $t = 0$. While ELBO not converged:

$$\begin{aligned} \Sigma_N^{(t+1)} &= \left(\mathbb{E}[\beta]^{(t)} \mathbf{K} + \mathbb{E}_\alpha [\text{diag}(\alpha_1, \dots, \alpha_D)]^{(t)} \right)^{-1} \\ \mu_N^{(t+1)} &= \Sigma_N^{(t+1)} \mathbb{E}[\beta]^{(t)} \mathbf{P} \\ b_N^{(t+1)} &= b_0 + \frac{1}{2} \left(\mathbf{Q} - 2\mathbf{P}^T \mathbb{E}[\mathbf{w}]^{(t+1)} + \text{tr} \left(\mathbf{K} \mathbb{E}[\mathbf{w}\mathbf{w}^T]^{(t+1)} \right) \right) \\ \mathbb{E}[\beta]^{(t+1)} &= \frac{a_N}{b_N^{(t+1)}} \\ f_{Nd}^{(t+1)} &= f_0 + \frac{1}{2} \mathbb{E}[w_d^2]^{(t+1)} \quad \forall d \\ \mathbb{E}[\alpha_d]^{(t+1)} &= \frac{e_{Nd}}{f_{Nd}^{(t+1)}} \quad \forall d \\ t &= t + 1 \end{aligned}$$

to avoid recomputation. Instead of making random initializations, we have assumed arbitrarily that \mathbf{w} is initially a standard normal variable.

Student Name: Suryateja B.V.

Roll Number: 160729

Date: March 14, 2019

1 BBVI

$$\begin{aligned}\nabla_{\phi} L(\phi) &= \mathbb{E}_q \left[\nabla_{\phi} \log q(\mathbf{w}|\phi) \left(\sum_{n=1}^N \log \sigma(y_n \mathbf{w}^T \mathbf{x}_n) + \log p(\mathbf{w}) - \log q(\mathbf{w}|\phi) \right) \right] \\ \nabla_{\mu} \log q(\mathbf{w}|\phi) &= \frac{d}{d\mu} \left(-\frac{1}{2} (\mathbf{w} - \mu)^T \Sigma^{-1} (\mathbf{w} - \mu) \right) = \Sigma^{-1} (\mathbf{w} - \mu) = (\mathbf{L}\mathbf{L}^T)^{-1} (\mathbf{w} - \mu) \\ \nabla_{\mathbf{L}} \log q(\mathbf{w}|\phi) &= \frac{d}{d\mathbf{L}} \left(-\frac{1}{2} \log |\mathbf{L}\mathbf{L}^T| \right) + \frac{d}{d\mathbf{L}} \left(-\frac{1}{2} (\mathbf{w} - \mu)^T \Sigma^{-1} (\mathbf{w} - \mu) \right) \frac{d\Sigma}{d\mathbf{L}} \\ &= -\mathbf{L}^{-T} + \left((\mathbf{L}\mathbf{L}^T)^{-1} (\mathbf{w} - \mu) (\mathbf{w} - \mu)^T (\mathbf{L}\mathbf{L}^T)^{-1} \right) \mathbf{L}\end{aligned}$$

$$\begin{aligned}\log p(\mathbf{w}) &= \frac{D}{2} \log \frac{\lambda}{2\pi} - \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\ \log q(\mathbf{w}|\phi) &= -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\mathbf{L}\mathbf{L}^T| - \frac{1}{2} (\mathbf{w} - \mu)^T (\mathbf{L}\mathbf{L}^T)^{-1} (\mathbf{w} - \mu)\end{aligned}$$

Instead of taking an expectation wrt $q(\mathbf{w}|\phi)$, we draw S samples from the distribution and do a Monte carlo approximation. The algorithm for BBVI with minibatch size of 1 is given in the next page. We note that the gradients derived does not depend on the type of model. Hence this algorithm is like a black-box (agnostic).

- 1 Initialize $\phi = \phi^{(0)} = \{\boldsymbol{\mu}^{(0)}, \mathbf{L}^{(0)}\}$. Set $t = 1$;
- 2 Pick a random example \mathbf{n} ;
- 3 Draw S samples $\{\mathbf{w}_s^{(t)}\}_{s=1}^S$ from $q(\mathbf{w}|\phi^{(t-1)})$;
- 4 Compute learning rates $\eta_{\mu}^{(t)}, \eta_L^{(t)}$;
- 5

$$\begin{aligned}\nabla_{\mu} L(\phi^{(t)}) &= \frac{1}{S} \sum_{s=1}^S \left(\nabla_{\mu} \log q(\mathbf{w}_s^{(t)}|\phi^{(t-1)}) \left(\log \sigma(y_n (\mathbf{w}_s^T)^{(t)} \mathbf{x}_n) + \log p(\mathbf{w}_s^{(t)}) \right. \right. \\ &\quad \left. \left. - \log q(\mathbf{w}_s^{(t)}|\phi^{(t-1)}) \right) \right) \\ \boldsymbol{\mu}^{(t)} &= \boldsymbol{\mu}^{(t-1)} + \eta_{\mu}^{(t)} \nabla_{\mu} L(\phi^{(t)}) \\ \nabla_L L(\phi^{(t)}) &= \frac{1}{S} \sum_{s=1}^S \left(\nabla_L \log q(\mathbf{w}_s^{(t)}|\boldsymbol{\mu}^{(t)}, \mathbf{L}^{(t-1)}) \left(\log \sigma(y_n (\mathbf{w}_s^T)^{(t)} \mathbf{x}_n) + \log p(\mathbf{w}_s^{(t)}) \right. \right. \\ &\quad \left. \left. - \log q(\mathbf{w}_s^{(t)}|\boldsymbol{\mu}^{(t)}, \mathbf{L}^{(t-1)}) \right) \right) \\ \mathbf{L}^{(t)} &= \mathbf{L}^{(t-1)} + \eta_L^{(t)} \nabla_L L(\phi^{(t)}) \\ t &= t + 1\end{aligned}$$

- 6 Go to Step-2 if ELBO not converged.

2 Reparametrization trick

Here we write $\mathbf{w} = \boldsymbol{\mu} + \mathbf{L}\mathbf{v}$, where $\mathbf{v} \sim \mathcal{N}(\mathbf{v}|0, 1)$. So, the elbo gradient is as follows -

$$\begin{aligned}\nabla_{\phi} L(\phi) &= \mathbb{E}_q \left[\nabla_{\phi} \left(\sum_{n=1}^N \log \sigma(y_n (\boldsymbol{\mu} + \mathbf{L}\mathbf{v})^T \mathbf{x}_n) + \log p(\boldsymbol{\mu} + \mathbf{L}\mathbf{v}) - \log q(\boldsymbol{\mu} + \mathbf{L}\mathbf{v}|\phi) \right) \right] \\ \nabla_{\phi} L(\phi) &= \mathbb{E}_{p(\mathbf{v})} [\nabla_{\phi} f(\boldsymbol{\mu} + \mathbf{L}\mathbf{v})] \\ \log p(\boldsymbol{\mu} + \mathbf{L}\mathbf{v}) &= \frac{D}{2} \log \frac{\lambda}{2\pi} - \frac{\lambda}{2} (\boldsymbol{\mu} + \mathbf{L}\mathbf{v})^T (\boldsymbol{\mu} + \mathbf{L}\mathbf{v}) \\ \log q(\boldsymbol{\mu} + \mathbf{L}\mathbf{v}|\phi) &= -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\mathbf{L}\mathbf{L}^T| - \frac{1}{2} \mathbf{v}^T \mathbf{v} \\ \nabla_{\mu} f(\boldsymbol{\mu} + \mathbf{L}\mathbf{v}) &= \sum_{n=1}^N \left(1 - \sigma(y_n (\boldsymbol{\mu} + \mathbf{L}\mathbf{v})^T \mathbf{x}_n) \right) y_n \mathbf{x}_n - \lambda (\boldsymbol{\mu} + \mathbf{L}\mathbf{v}) - 0 \\ \nabla_L f(\boldsymbol{\mu} + \mathbf{L}\mathbf{v}) &= \sum_{n=1}^N \left(1 - \sigma(y_n (\boldsymbol{\mu} + \mathbf{L}\mathbf{v})^T \mathbf{x}_n) \right) y_n \mathbf{x}_n \mathbf{v}^T - \lambda (\boldsymbol{\mu} + \mathbf{L}\mathbf{v}) \mathbf{v}^T - \mathbf{L}^{-T}\end{aligned}$$

We use monte carlo gradien again, by drawing S samples from the distribution $p(\mathbf{v})$. Algorithm is given in the next page

- 1 Initialize $\phi = \phi^{(0)} = \{\boldsymbol{\mu}^{(0)}, \mathbf{L}^{(0)}\}$. Set $t = 1$;
- 2 Pick a random example \mathbf{n} ;
- 3 Draw S samples $\left\{\mathbf{v}_s^{(t)}\right\}_{s=1}^S$ from $\mathcal{N}(\mathbf{v}|\mathbf{0}, \mathbf{I})$;
- 4 Compute learning rates $\eta_{\mu}^{(t)}, \eta_L^{(t)}$;
- 5

$$\nabla_{\mu} L\left(\phi^{(t)}\right)=\frac{1}{S} \sum_{s=1}^S \nabla_{\mu} f\left(\boldsymbol{\mu}^{(t-1)}+\mathbf{L}^{(t-1)} \mathbf{v}_s^{(t)}\right)$$

$$\boldsymbol{\mu}^{(t)}=\boldsymbol{\mu}^{(t-1)}+\eta_{\mu}^{(t)} \nabla_{\mu} L\left(\phi^{(t)}\right)$$

$$\nabla_L L\left(\phi^{(t)}\right)=\frac{1}{S} \sum_{s=1}^S \nabla_L f\left(\boldsymbol{\mu}^{(t)}+\mathbf{L}^{(t-1)} \mathbf{v}_s^{(t)}\right)$$

$$\mathbf{L}^{(t)}=\mathbf{L}^{(t-1)}+\eta_L^{(t)} \nabla_L L\left(\phi^{(t)}\right)$$

$$t=t+1$$

- 6 Go to Step-2 if ELBO not converged.

Student Name: Suryateja B.V.

Roll Number: 160729

Date: March 14, 2019

A few notes -

- logsumexp trick has been used to calculate since some probabilities are very small.
- The various mean images obtained resemble original numbers.
- Online EM is as expected very fast compared to full batch EM.
- Sometimes mean images, repeat, as in, resemble almost same numbers. This is due to either poor initialization, or the algorithm has captured two different types of the same digit. For example, 0 tilted towards right and 0 tilted towards left might be two different mean images.
- All images are available in submission file. You can also view images by varying K in the Jupyter notebook file.