

Student Name: Gurpreet Singh
Roll Number: 150259
Date: February 2, 2018

Question 1

Here, I state the notations, terms given to us and some observations.

1. Likelihood

$$\mathbb{P}[\mathbf{y} | X, \mathbf{w}, \beta] = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^T \mathbf{x}_n, \beta^{-1})$$

2. Let $\boldsymbol{\lambda}$ represent $[\lambda_1, \lambda_2 \dots \lambda_d]^T$ and Λ represent the $D \times D$ diagonal matrix with element at d, d be equal to λ_d

3. Prior

$$\mathbb{P}[\mathbf{w} | \boldsymbol{\lambda}] = \prod_{d=1}^D \mathcal{N}(w_d | 0, \lambda_d^{-1})$$

Part A

I have shown that the marginal likelihood is in closed form for both the cases, however am showing the working for only the second case, and then simplify for the first case by putting $\forall d \in [D], \lambda_d = \lambda$.

In order to compute the marginal likelihood, we can compute the posterior of our parameter *i.e.* \mathbf{w} , and then use Bayes Rule to compute the marginal likelihood.

$$\mathbb{P}[\mathbf{y} | X, \beta, \boldsymbol{\lambda}] = \frac{\mathbb{P}[\mathbf{y} | X, \mathbf{w}, \beta] \mathbb{P}[\mathbf{w} | \boldsymbol{\lambda}]}{\mathbb{P}[\mathbf{w} | \mathbf{y}, X, \beta, \boldsymbol{\lambda}]} \quad (1)$$

We start by computing the posterior over \mathbf{w} .

$$\begin{aligned} \mathbb{P}[\mathbf{w} | \mathbf{y}, X, \beta, \boldsymbol{\lambda}] &\propto \mathbb{P}[\mathbf{y} | X, \mathbf{w}, \beta] \mathbb{P}[\mathbf{w} | \boldsymbol{\lambda}] \\ &\propto \prod_{n=1}^N \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \beta^{-1}) \prod_{d=1}^D \mathcal{N}(w_d, \lambda_d^{-1}) \\ &\propto \prod_{n=1}^N \exp\left(\frac{-\beta}{2} (y_n - \mathbf{w}^T \mathbf{x}_n)^2\right) \prod_{d=1}^D \exp\left(\frac{-\lambda_d}{2} (w_d)^2\right) \\ &\propto \exp\left(\frac{-1}{2} \left(\beta \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \sum_{d=1}^D \lambda_d w_d^2 \right)\right) \\ &\propto \exp\left(\frac{-1}{2} (\beta \mathbf{w}^T X^T X \mathbf{w} - 2 \mathbf{w}^T X^T \mathbf{y} + \mathbf{w}^T \Lambda \mathbf{w})\right) \end{aligned}$$

Note. I have ignored all terms independent of \mathbf{w} during the above computations

The last term looks like a Gaussian distribution, and hence, we can use the Completing the Squares trick to find out the exact distribution. Therefore, we get

$$\mathbb{P}[\mathbf{w} | \mathbf{y}, X, \beta, \boldsymbol{\lambda}] = \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_w, \Sigma_w) \quad (2)$$

where

$$\begin{aligned} \Sigma_w &= (\beta X^T X + \Lambda)^{-1} \\ \boldsymbol{\mu}_w &= \left(X^T X + \frac{1}{\beta} \Lambda \right)^{-1} X^T \mathbf{y} \end{aligned}$$

Therefore, we have the posterior. Now we can put this in Equation 1.

$$\begin{aligned} \mathbb{P}[\mathbf{y} | X, \beta, \boldsymbol{\lambda}] &= \frac{\mathbb{P}[\mathbf{y} | X, \mathbf{w}, \beta] \mathbb{P}[\mathbf{w} | \boldsymbol{\lambda}]}{\mathbb{P}[\mathbf{w} | \mathbf{y}, X, \beta, \boldsymbol{\lambda}]} \\ &\propto \exp \left(\frac{-1}{2} \left(\beta \sum_{n=1}^N y_n^2 + \beta \mathbf{w}^T X^T X \mathbf{w} - 2\beta \mathbf{w}^T X^T \mathbf{y} + \mathbf{w}^T \Lambda \mathbf{w} \right) + \frac{1}{2} \left((\mathbf{w} - \boldsymbol{\mu}_w)^T \Sigma_w^{-1} (\mathbf{w} - \boldsymbol{\mu}_w) \right) \right) \\ &\propto \exp \left(\frac{-1}{2} \left(\beta \sum_{n=1}^N y_n^2 - \boldsymbol{\mu}_w^T \Sigma_w^{-1} \boldsymbol{\mu}_w \right) \right) \end{aligned}$$

Putting $\boldsymbol{\mu}_w = \beta \cdot \Sigma_w X^T \mathbf{y}$,

$$\begin{aligned} \mathbb{P}[\mathbf{y} | X, \beta, \boldsymbol{\lambda}] &\propto \exp \left(\frac{-1}{2} \left(\beta \sum_{n=1}^N y_n^2 - (\beta \cdot \Sigma_w X^T \mathbf{y})^T \Sigma_w^{-1} \beta \cdot \Sigma_w X^T \mathbf{y} \right) \right) \\ &\propto \exp \left(\frac{-1}{2} \left(\mathbf{y}^T \left(\beta \mathbf{I} - \beta^2 X \Sigma_w^T X^T \right) \mathbf{y} \right) \right) \end{aligned}$$

This again is a gaussian distribution. Therefore, we can write the final form of this distribution as follows

$$\mathbb{P}[\mathbf{y} | X, \beta, \boldsymbol{\lambda}] = \mathcal{N}(\mathbf{y} | \mathbf{0}, \Sigma) \quad (3)$$

where $\Sigma = (\beta \mathbf{I} - \beta^2 X \Sigma_w X^T)^{-1}$

We can further simplify the term of the gaussian using a property of matrix inverses. Consider the inverse property given in Equation 166 in [1], which states

$$(\mathbf{I} + AB)^{-1} = \mathbf{I} - A(\mathbf{I} + BA)^{-1} B \quad (4)$$

First let us rewrite Σ

$$\begin{aligned} \Sigma &= \frac{1}{\beta} \left(\mathbf{I} - \beta X (\beta X^T X + \Lambda)^{-1} X^T \right) \\ \Rightarrow \beta \Sigma &= \left(\mathbf{I} + \left(-X \left(X^T X + \frac{1}{\beta} \Lambda \right)^{-1} \right) X^T \right) \end{aligned}$$

We want to represent this in the form of the term given in the LHS of Equation 4.

Therefore, let $A = \left(-X \left(X^T X + \frac{1}{\beta} \Lambda \right)^{-1} \right)$ and $B = X^T$. Then

$$\begin{aligned}
\beta \Sigma &= (\mathbf{I} + AB)^{-1} \\
&= \mathbf{I} - A(\mathbf{I} + BA)^{-1} B \\
&= \mathbf{I} - A \left(\mathbf{I} - X^T X \left(X^T X + \frac{\Lambda}{\beta} \right)^{-1} \right)^{-1} B \\
&= \mathbf{I} - A \left(\mathbf{I} - \left(X^T X + \frac{\Lambda}{\beta} - \frac{\Lambda}{\beta} \right) \left(X^T X + \frac{\Lambda}{\beta} \right)^{-1} \right)^{-1} B \\
&= \mathbf{I} - A \left(\frac{\Lambda}{\beta} \left(X^T X + \frac{\Lambda}{\beta} \right)^{-1} \right)^{-1} B \\
&= \mathbf{I} - \left(-X \left(X^T X + \frac{\Lambda}{\beta} \right)^{-1} \left(X^T X + \frac{\Lambda}{\beta} \right) \left(\frac{\Lambda}{\beta} \right)^{-1} \right) X^T \\
&= \mathbf{I} + \beta X \Lambda^{-1} X^T \\
\Rightarrow \Sigma &= \frac{1}{\beta} \mathbf{I} + X \Lambda^{-1} X^T
\end{aligned}$$

Hence, we can now rewrite our answer,

$$\mathbb{P}[\mathbf{y} | X, \beta, \boldsymbol{\lambda}] = \mathcal{N}(\mathbf{y} | \mathbf{0}, \Sigma) \quad (5)$$

where $\Sigma = \beta^{-1} \mathbf{I} + X \Lambda^{-1} X^T$

If for all $d = 1 \dots D$, $\lambda_d = \lambda$, we can simply put this in the result obtained in Equation 5

Part B

As discussed in class, we do MLE-II inference by performing MLE on the marginal likelihood $\mathbb{P}[\mathbf{y} | X, \beta, \boldsymbol{\lambda}]$, and then approximate the posterior using the conditional posterior of \mathbf{w} ($\mathbb{P}[\mathbf{w} | \mathbf{y}, X, \hat{\beta}, \hat{\boldsymbol{\lambda}}]$).

Note. We define $\boldsymbol{\lambda}$ and Λ the same as in part A

Taking a cue from the MLE-II derivation given in PRML 3.5 [2], we represent the marginal in a different manner, and then try to maximize it for the purpose of MLE.

$$\begin{aligned}
\mathbb{P}[\mathbf{y} | X, \beta, \boldsymbol{\lambda}] &= \int_{\mathbf{w}} \mathbb{P}[\mathbf{y} | X, \mathbf{w}, \beta] \mathbb{P}[\mathbf{w} | X, \beta, \boldsymbol{\lambda}] d\mathbf{w} \\
&= \left(\frac{\beta}{2\pi} \right)^{\frac{N}{2}} \sqrt{\frac{|\Lambda|}{(2\pi)^D}} \int_{\mathbf{w}} \exp(-E(\mathbf{w})) d\mathbf{w}
\end{aligned}$$

where $E : \mathbf{w} \mapsto \mathbb{R}$ is the evidence function [2] and is defined as

$$E(\mathbf{w}) = \frac{\beta}{2} (\mathbf{y} - X\mathbf{w})^T (\mathbf{y} - X\mathbf{w}) + \frac{1}{2} \mathbf{w}^T \Lambda \mathbf{w}$$

We can also represent $E(\mathbf{w})$ using $\boldsymbol{\mu}_w$ which is the mean of the posterior as derived in part A.

$$\begin{aligned}
2E(\mathbf{w}) &= \beta (\mathbf{y} - X\boldsymbol{\mu}_w - X(\mathbf{w} - \boldsymbol{\mu}_w))^T (\mathbf{y} - X\boldsymbol{\mu}_w - X(\mathbf{w} - \boldsymbol{\mu}_w)) + \mathbf{w}^T \Lambda \mathbf{w} \\
&= \beta (\mathbf{y} - X\boldsymbol{\mu}_w)^T (\mathbf{y} - X\boldsymbol{\mu}_w) + (\mathbf{w} - \boldsymbol{\mu}_w)^T \beta X^T X (\mathbf{w} - \boldsymbol{\mu}_w) - 2(\mathbf{w} - \boldsymbol{\mu}_w)^T \beta X^T (\mathbf{y} - X\boldsymbol{\mu}_w)
\end{aligned}$$

We can write $\beta X^T (y - X\boldsymbol{\mu})$ as $\Lambda\boldsymbol{\mu}_w$ by expanding $\boldsymbol{\mu}_w$ as shown below.

$$\begin{aligned}\boldsymbol{\mu}_w &= \beta (\beta X^T X + \Lambda)^{-1} X^T \mathbf{y} \\ \Rightarrow \beta X^T (\mathbf{y} - X\boldsymbol{\mu}_w) &= \beta \left(\mathbf{I} - \beta X^T X (\beta X^T X + \Lambda)^{-1} \right) X^T \mathbf{y} = \beta \Lambda (\beta X^T X + \Lambda)^{-1} X^T \mathbf{y} = \Lambda \boldsymbol{\mu}_w\end{aligned}$$

<+++> Therefore

$$\begin{aligned}E(\mathbf{w}) &= \beta (\mathbf{y} - X\boldsymbol{\mu}_w)^T (\mathbf{y} - X\boldsymbol{\mu}_w) + (\mathbf{w} - \boldsymbol{\mu}_w)^T \beta X^T X (\mathbf{w} - \boldsymbol{\mu}_w) - 2(\mathbf{w} - \boldsymbol{\mu}_w)^T \Lambda \boldsymbol{\mu}_w + \mathbf{w}^T \Lambda \mathbf{w} \\ &= \beta (\mathbf{y} - X\boldsymbol{\mu}_w)^T (\mathbf{y} - X\boldsymbol{\mu}_w) + (\mathbf{w} - \boldsymbol{\mu}_w)^T \beta X^T X (\mathbf{w} - \boldsymbol{\mu}_w) - 2\mathbf{w}^T \Lambda \boldsymbol{\mu}_w + 2\boldsymbol{\mu}_w^T \Lambda \boldsymbol{\mu}_w + \mathbf{w}^T \Lambda \mathbf{w} \\ &= \beta (\mathbf{y} - X\boldsymbol{\mu}_w)^T (\mathbf{y} - X\boldsymbol{\mu}_w) + (\mathbf{w} - \boldsymbol{\mu}_w)^T (\beta X^T X + \Lambda) (\mathbf{w} - \boldsymbol{\mu}_w) + \boldsymbol{\mu}_w^T \Lambda \boldsymbol{\mu}_w \\ &= 2E(\boldsymbol{\mu}_w) + (\mathbf{w} - \boldsymbol{\mu}_w)^T \Sigma^{-1} (\mathbf{w} - \boldsymbol{\mu}_w)\end{aligned}$$

Therefore,

$$E(\mathbf{w}) = E(\boldsymbol{\mu}_w) + \frac{1}{2} (\mathbf{w} - \boldsymbol{\mu}_w)^T \Sigma^{-1} (\mathbf{w} - \boldsymbol{\mu}_w) \quad (6)$$

Therefore, we can now write the Marginal Likelihood in this form

$$\int_{\mathbf{w}} \exp(-E(w)) = \exp(-E(\boldsymbol{\mu}_w)) \int_{\mathbf{w}} \exp\left(-\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu}_w)^T \Sigma^{-1} (\mathbf{w} - \boldsymbol{\mu}_w)\right)$$

Since the integral is just over the exponent part of the posterior distribution of \mathbf{w} , we can write

$$\int_{\mathbf{w}} \exp(-E(w)) = \exp(-E(\boldsymbol{\mu}_w)) \sqrt{(2\pi)^D |\Sigma|}$$

Hence, the negative log marginal can be written as

$$-\log(\mathbb{P}[\mathbf{y} | X, \beta, \boldsymbol{\lambda}]) = \frac{1}{2} E(\boldsymbol{\mu}_w) - \frac{1}{2} \log(|\Sigma|) - \log(|\Lambda|) - \frac{N}{2} \log(\beta) + \text{constants} \quad (7)$$

Now, we can use iterative MLE-II scheme to solve for the hyperparamters.

Suppose we have estimates of β and $\boldsymbol{\lambda}$ as β^t and $\boldsymbol{\lambda}^t$ from the t^{th} iteration of the MLE-II scheme, then we can compute the conditional posterior of \mathbf{w} given these estimates. Note, we have already computed this posterior, as given in Equation 2. Therefore, we can write

$$\mathbb{P}[\mathbf{w}^{t+1} | \mathbf{y}, X, \beta^t, \boldsymbol{\lambda}^t] = \mathcal{N}(\mathbf{w}^{t+1} | \boldsymbol{\mu}_w^t, \Sigma_w^t)$$

where

$$\begin{aligned}\Sigma_w^t &= (\beta^t X^T X + \Lambda^t)^{-1} \\ \boldsymbol{\mu}_w^t &= \left(X^T X + \frac{1}{\beta^t} \Lambda^t \right)^{-1} X^T \mathbf{y}\end{aligned}$$

Now, we can continue with the point estimation step of the MLE-II method, *i.e.* estimating the update values of β and $\boldsymbol{\lambda}$.

$$\beta^{t+1}, \boldsymbol{\lambda}^{t+1} = \arg \min_{\beta, \boldsymbol{\lambda}} -\log(\mathbb{P}[\mathbf{y} | X, \beta, \boldsymbol{\lambda}])$$

Note, we do not directly estimate the MLE updates, instead use an iterative method (like Jacobi method). Therefore, instead of solving for the equations (partial derivatives equated to zero), we will try to find an equation which allows for easy updates of the hyperparamters.

Again, taking a cue from PRML [2], we will consider the estimation of λ (particular λ_d for some d) first.

From the properties of Determinants and Eigenvalues, we know $|A| = \prod_{k=1}^K \alpha_k$ where α_k is the k^{th} eigenvalue of A .

Let the eigenvalues of $(\beta X^T X + \Lambda^t - \lambda_d^{t+1})$ be represented by the set $\{\gamma_{d'}\}_{d'=1}^D$, then we can say that the eigen values of $(\beta X^T X + \Lambda^t)$ will be $\{\gamma_{d'} + \lambda_d^{t+1}\}_{d'=1}^D$. Therefore, we can now put this in Equation 7

Note. $|\Lambda| = \prod_{d'=1}^D \lambda_{d'}$ as Λ is just a diagonal matrix

$$-\log(\mathbb{P}[\mathbf{y} | X, \beta, \lambda]) = \frac{1}{2} E(\mu_w^t) - \frac{1}{2} \sum_{d'=1}^D \log(\gamma_{d'} + \lambda_d^{t+1}) - \frac{1}{2} \sum_{d'=1}^D \log(\lambda_d^{t+1}) - \frac{N}{2} \log(\beta)$$

Although we cannot compute $\gamma_{d'}$ as we do not yet know the value of λ_d^{t+1} , however, using the idea of the Jacobi method, we can estimate it by putting λ_d^t instead of λ_d^{t+1} . Now, differentiating with respect to λ_d^{t+1} , and equating to 0, we get

$$0 = \frac{1}{2\lambda_d^{t+1}} - \frac{1}{2} \{\mu_{d,w}^t\}^2 - \frac{1}{2} \sum_{d'=1}^D \frac{1}{\gamma_{d'} + \lambda_d^{t+1}}$$

Note. We are putting in the value of λ_d^{t+1} instead of λ_d^t in the term for Λ and hence we get the derivative of $E(\mu_w)$ as what is shown.

$$\lambda_d^{t+1} = \frac{1}{\{\mu_{d,w}^t\}^2} \left(1 - \sum_{d'=1}^D \frac{\lambda_d^{t+1}}{\gamma_{d'} + \lambda_d^{t+1}} \right)$$

In the fraction part, again using the idea from Jacobi method, we can estimate λ_d^{t+1} by λ_d^t , and therefore we can write

$$\lambda_d^{t+1} = \frac{1}{\{\mu_{d,w}^t\}^2} \left(1 - \sum_{d'=1}^D \frac{\lambda_d^t}{\gamma_{d'} + \lambda_d^t} \right) \triangleq \frac{\alpha_d}{\{\mu_{d,w}^t\}^2}$$

If for all $d = 1 \dots D$, $\lambda_d = \lambda$, we cannot directly write this, since all $\lambda_{d'}$ are not independent. The only change in that case will be the derivatives with respect to λ . Now,

$$\begin{aligned} 0 &= \frac{D}{2} \lambda^{t+1} - \frac{1}{2} \{\mu_w^t\}^T \mu_w^t - \frac{1}{2} \sum_{d'=1}^D \frac{1}{\gamma_{d'} + \lambda^{t+1}} \\ \Rightarrow \lambda^{t+1} &= \frac{\alpha}{\{\mu_w^t\}^T \mu_w^t} \end{aligned}$$

We now give the MLE update for β^{t+1} . Since in this case we can use all the estimates $\{\lambda_d^t\}_{d=1}^D$, we can simply define γ'_d to be the eigenvalues of $\beta X^T X$.

Since for some eigenvector \mathbf{u}_d , $\beta X^T X \mathbf{u}_d = \gamma'_d \mathbf{u}_d$, we can say

$$\begin{aligned} \frac{d(\gamma'_d)}{d(\beta)} \mathbf{u}_d &= X^T X \mathbf{u}_d \\ \Rightarrow \frac{d(\gamma'_d)}{d(\beta)} \beta \mathbf{u}_d &= \beta X^T X \mathbf{u}_d = \gamma'_d \mathbf{u}_d \\ \Rightarrow \frac{d(\gamma'_d)}{d(\beta)} &= \frac{\gamma'_d}{\beta} \end{aligned}$$

Therefore, we can write

$$0 = \frac{N}{2\beta} - \frac{1}{2}(\mathbf{y} - X\boldsymbol{\mu}_w^t)^T (\mathbf{y} - X\boldsymbol{\mu}_w^t) - \sum_{d=1}^D \frac{\beta}{2} \langle ++ \rangle$$

$\langle ++ \rangle$

Therefore, we can give an iterative algorithm for this (as discussed in class)

Algorithm 1: MLE-II for Hyperparameter Estimation

1. Assume some initial estimates for the hyperparameters, β^0 and $\boldsymbol{\Lambda}^0$
2. Repeat until convergence for $t = 1, 2, \dots$
 - (a) Estimate the posterior over \mathbf{w}^t as

$$\mathbb{P}[\mathbf{w}^t | \mathbf{y}, X, \beta^{t-1}, \boldsymbol{\Lambda}^{t-1}] = \mathcal{N}(\mathbf{w}^t | \boldsymbol{\mu}_w^{t-1}, \Sigma_w^{t-1})$$

where

$$\begin{aligned} \Sigma_w^t &= (\beta^t X^T X + \Lambda^t)^{-1} \\ \boldsymbol{\mu}_w^t &= \left(X^T X + \frac{1}{\beta^t} \Lambda^t \right)^{-1} X^T \mathbf{y} \end{aligned}$$

- (b) Re-estimate the hyperparameters as

$$\begin{aligned} \boldsymbol{\Lambda}^t &= \\ \beta^t &= \end{aligned}$$

References

- [1] Matrix Cookbook. *Kaare Brandt Petersen, Michael Syskind Pedersen*
<http://www.matrixcookbook.com>
- [2] Pattern Recognition and Machine Learning. 2006. *Christopher M. Bishop*
- [3] Order of Integration (n.d.) *Wikipedia*
[https://en.wikipedia.org/wiki/Order_of_integration_\(calculus\)](https://en.wikipedia.org/wiki/Order_of_integration_(calculus))
- [4] Moment Generating Function *Wikipedia*
https://en.wikipedia.org/wiki/Moment-generating_function