

Student Name: Vishweshwar Tyagi

Roll Number: 191173

Date: March 1, 2021

Solution 1

We are given that $x \sim N(0, \eta)$, hence,

$$p(x | \eta) = \frac{1}{\sqrt{2\pi\eta}} \exp\left(-\frac{1}{2\eta}x^2\right) \quad x \in \mathbb{R} \quad (1)$$

Also, it is given that $\eta \sim \text{Exp}(\frac{\gamma^2}{2})$ with $\gamma > 0$, hence,

$$p(\eta | \gamma) = \frac{\gamma^2}{2} \exp\left(-\frac{\gamma^2}{2}\eta\right) \quad \eta > 0 \quad (2)$$

Hence, we have the pdf. for marginal distribution of x given by

$$\begin{aligned} p(x | \gamma) &= \int_{\eta>0} p(x, \eta | \gamma) d\eta \\ &= \int_{\eta>0} p(x, | \eta) p(\eta | \gamma) d\eta \\ &= \int_{\eta>0} \frac{\gamma^2}{2\sqrt{2\pi\eta}} \exp\left(-\frac{1}{2\eta}x^2 - \frac{\gamma^2}{2}\eta\right) d\eta \end{aligned} \quad (3)$$

This integral is hard to compute. We will try to identify the marginal distribution of x using the method of **Moment Generating Function** (MGF) and using the fact that MGFs are unique for each distribution.

The Moment Generating Function for the marginal distribution of x is given by

$$\begin{aligned} M_X(t) &= E_X(e^{tx}) \\ &= \int_{x=-\infty}^{\infty} e^{tx} p(x | \gamma) dx \\ &= \int_{x=-\infty}^{\infty} e^{tx} \left\{ \int_{\eta>0} \frac{\gamma^2}{2\sqrt{2\pi\eta}} \exp\left(-\frac{1}{2\eta}x^2 - \frac{\gamma^2}{2}\eta\right) d\eta \right\} dx \quad (\text{using (3)}) \\ &= \int_{x=-\infty}^{\infty} \int_{\eta>0} \frac{\gamma^2}{2\sqrt{2\pi\eta}} \exp\left(-\frac{\gamma^2}{2}\eta\right) \exp\left(-\frac{1}{2\eta}(x^2 - 2x t\eta)\right) d\eta dx \\ &= \int_{x=-\infty}^{\infty} \int_{\eta>0} \frac{\gamma^2}{2\sqrt{2\pi\eta}} \exp\left(-\frac{\eta(\gamma^2 - t^2)}{2}\right) \exp\left(-\frac{1}{2\eta}(x - t\eta)^2\right) d\eta dx \\ &= \int_{\eta>0} \frac{\gamma^2}{2} \exp\left(-\frac{\eta(\gamma^2 - t^2)}{2}\right) \left\{ \int_{x=-\infty}^{\infty} \frac{1}{\sqrt{2\pi\eta}} \exp\left(-\frac{1}{2\eta}(x - t\eta)^2\right) dx \right\} d\eta \end{aligned}$$

Note that

$$\int_{x=-\infty}^{\infty} \frac{1}{\sqrt{2\pi\eta}} \exp\left(-\frac{1}{2\eta}(x - t\eta)^2\right) dx = 1 \quad (4)$$

$$\int_{\eta>0} \left(\frac{\gamma^2 - t^2}{2}\right) \exp\left(-\frac{\eta(\gamma^2 - t^2)}{2}\right) d\eta = 1 \quad (5)$$

Hence, we get,

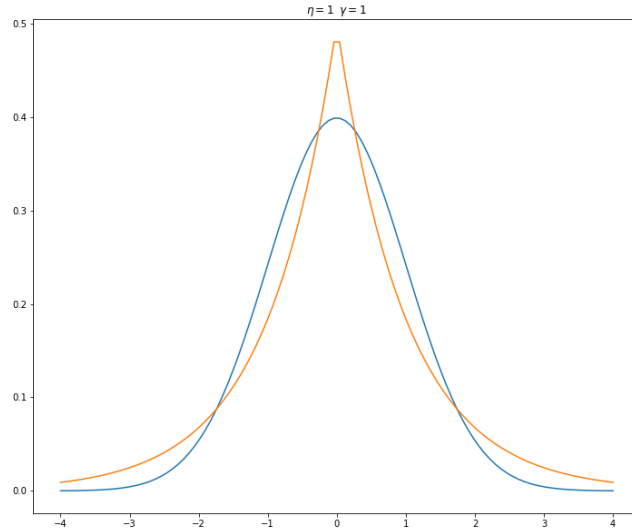
$$\begin{aligned} M_X(t) &= \int_{\eta>0} \frac{\gamma^2}{2} \exp\left(-\frac{\eta(\gamma^2 - t^2)}{2}\right) \left\{ \int_{x=-\infty}^{\infty} \frac{1}{\sqrt{2\pi\eta}} \exp\left(-\frac{1}{2\eta}(x - t\eta)^2\right) dx \right\} d\eta \\ &= \frac{\gamma^2}{2} \left(\frac{\gamma^2 - t^2}{2}\right)^{-1} \int_{\eta>0} \left(\frac{\gamma^2 - t^2}{2}\right) \exp\left(-\frac{\eta(\gamma^2 - t^2)}{2}\right) d\eta \quad (\text{using (4)}) \\ &= \frac{\gamma^2}{2} \left(\frac{\gamma^2 - t^2}{2}\right)^{-1} \text{ where } \frac{\gamma^2 - t^2}{2} > 0 \quad (\text{using (5)}) \\ &= \frac{1}{1 - (t/\gamma)^2} \text{ where } |t| < \gamma \end{aligned}$$

This compares with the MGF of Laplace distribution $L(\mu, b)$ where

$$M_L(t) = \frac{e^{t\mu}}{1 - b^2 t^2} \text{ where } |t| < \frac{1}{b} \text{ with } \mu = 0 \text{ and } b = \frac{1}{\gamma}$$

and the required marginal distribution of x is Laplacian distribution $L(x | 0, \frac{1}{\gamma})$ with

$$p(x | \gamma) = \frac{\gamma}{2} \exp(-\gamma |x|) \quad \forall x \in \mathbb{R} \quad (6)$$



While one would expect the marginal likelihood to follow Gaussian, since the likelihood itself follows Gaussian distribution, this is not the case here due to non conjugacy of Exponential prior with Gaussian likelihood. Also the calculations were fairly involved and generally we look for conjugate priors for this reason.

Student Name: Vishweshwar Tyagi

Roll Number: 191173

Date: March 1, 2021

Solution 2

Let $M \in \mathbf{R}^{D \times D}$, $\mathbf{v} \in \mathbf{R}^{D \times 1}$. We are given that,

$$(M + \mathbf{v}\mathbf{v}^T)^{-1} = M^{-1} - \frac{(M^{-1}\mathbf{v})(\mathbf{v}^T M^{-1})}{1 + \mathbf{v}^T M^{-1}\mathbf{v}} \quad (1)$$

Note that,

$$\begin{aligned} \Sigma_N &= (\beta \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T + \lambda I_D)^{-1} \\ &= \frac{1}{\lambda} \left(\frac{\beta}{\lambda} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T + I_D \right)^{-1} \end{aligned} \quad (2)$$

$$= \frac{1}{\lambda} \left(\frac{\beta}{\lambda} \mathbf{x}_N^T \mathbf{x}_N + \frac{\beta}{\lambda} \sum_{i=1}^{N-1} \mathbf{x}_i \mathbf{x}_i^T + \lambda I_D \right)^{-1} \quad (3)$$

Let $f(N) = (\frac{\beta}{\lambda} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T + I_D)^{-1}$, then, using (1), we can continue from (3) as follows

$$\Sigma_N = \frac{1}{\lambda} \left[f(N-1) - \frac{[f(N-1)\mathbf{x}_N] [\mathbf{x}_N^T f(N-1)]}{1 + \mathbf{x}_N^T f(N-1)\mathbf{x}_N} \right] \quad (4)$$

Hence, using (4) and (2), we can further obtain

$$\begin{aligned} \Sigma_N - \Sigma_{N-1} &= \frac{1}{\lambda} \left[f(N-1) - \frac{[f(N-1)\mathbf{x}_N] [\mathbf{x}_N^T f(N-1)]}{1 + \mathbf{x}_N^T f(N-1)\mathbf{x}_N} \right] - \frac{1}{\lambda} \left(\frac{\beta}{\lambda} \sum_{i=1}^{N-1} \mathbf{x}_i \mathbf{x}_i^T + I_D \right)^{-1} \\ &= \frac{1}{\lambda} \left[f(N-1) - \frac{[f(N-1)\mathbf{x}_N] [\mathbf{x}_N^T f(N-1)]}{1 + \mathbf{x}_N^T f(N-1)\mathbf{x}_N} \right] - \frac{1}{\lambda} f(N-1) \\ &= -\frac{1}{\lambda} \left[\frac{[f(N-1)\mathbf{x}_N] [\mathbf{x}_N^T f(N-1)]}{1 + \mathbf{x}_N^T f(N-1)\mathbf{x}_N} \right] \\ &= -\frac{1}{\lambda} \left[\frac{[\mathbf{x}_N^T f(N-1)]^T [\mathbf{x}_N^T f(N-1)]}{1 + \mathbf{x}_N^T f(N-1)\mathbf{x}_N} \right] \end{aligned} \quad (5)$$

Now note that $f(N) = (\frac{\beta}{\lambda} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T + I_D)^{-1} = \lambda \Sigma_N \quad \forall N \in \mathbb{N}$

Σ_N being a covariance matrix is symmetric, positive semi-definite and since $\lambda > 0$, we have $f(N)$ as symmetric, positive semi-definite $\quad \forall N \in \mathbb{N}$ (6)

Using (5) and (6), we have

$$\Sigma_N - \Sigma_{N-1} = -\frac{1}{\lambda} \left[\frac{[\mathbf{x}_N^T f(N-1)]^T [\mathbf{x}_N^T f(N-1)]}{1 + \mathbf{x}_N^T f(N-1) \mathbf{x}_N} \right] \quad (7)$$

Now, $f(N-1)$ is positive semi-definite from (6), hence, for arbitrary $\mathbf{x}_N \in \mathbb{R}^{D \times 1}$, we have,

$$1 + \mathbf{x}_N^T f(N-1) \mathbf{x}_N \geq 1 \quad (8)$$

Also, it is easy to see that $[\mathbf{x}_N^T f(N-1)]^T [\mathbf{x}_N^T f(N-1)]$ ($= Y$, say) is positive semi-definite, since for arbitrary $\mathbf{x} \in \mathbb{R}^{D \times 1}$, we have

$$\begin{aligned} \mathbf{x}^T Y \mathbf{x} &= \mathbf{x}^T [\mathbf{x}_N^T f(N-1)]^T [\mathbf{x}_N^T f(N-1)] \mathbf{x} \\ &= [\mathbf{x}_N^T f(N-1) \mathbf{x}]^T [\mathbf{x}_N^T f(N-1) \mathbf{x}] \\ &= \|\mathbf{x}_N^T f(N-1) \mathbf{x}\|_2^2 \geq 0 \end{aligned}$$

From (8) and (9),

$$\frac{[\mathbf{x}_N^T f(N-1)]^T [\mathbf{x}_N^T f(N-1)]}{1 + \mathbf{x}_N^T f(N-1) \mathbf{x}_N} \text{ is positive semi-definite}$$

Since $\frac{-1}{\lambda} < 0$, we have, $-\frac{1}{\lambda} \left[\frac{[\mathbf{x}_N^T f(N-1)]^T [\mathbf{x}_N^T f(N-1)]}{1 + \mathbf{x}_N^T f(N-1) \mathbf{x}_N} \right]$ as negative semi-definite

Thus, using (7), $\Sigma_N - \Sigma_{N-1}$ is negative semi-definite (10)

Now, we have for arbitrary $\mathbf{x}_* \in \mathbb{R}^{D \times 1}$

$$\begin{aligned} \sigma_N^2(\mathbf{x}_*) - \sigma_{N-1}^2(\mathbf{x}_*) &= [\beta^{-1} + \mathbf{x}_*^T \Sigma_N \mathbf{x}_*] - [\beta^{-1} + \mathbf{x}_*^T \Sigma_{N-1} \mathbf{x}_*] \\ &= \mathbf{x}_*^T [\Sigma_N - \Sigma_{N-1}] \mathbf{x}_* \leq 0 \end{aligned} \quad (\text{using (10)})$$

Thus,

$$\beta^{-1} \leq \sigma_N^2(\mathbf{x}_*) \leq \sigma_{N-1}^2(\mathbf{x}_*) \quad (11)$$

Thus, the variance of the predictive posterior decreases as the size of the training set, N , increases (12)

Also, we see that $\langle \sigma_N^2(\mathbf{x}_*) \rangle_N = \langle \beta^{-1} + \mathbf{x}_*^T \Sigma_N \mathbf{x}_* \rangle_N$ is a monotonically decreasing sequence bounded below by β^{-1} , with $\mathbf{x}_*^T \Sigma_N \mathbf{x}_* \rightarrow 0$ as $N \rightarrow \infty$ (because of (10)), thus,

$$\langle \sigma_N^2(\mathbf{x}_*) \rangle_N \rightarrow \beta^{-1} \text{ as } N \rightarrow \infty \quad (13)$$

Solution 3

We are given that $x_i \sim N(\mu, \sigma^2) \mid_{i=1}^N$ are independent and identically distributed, hence

$$E(x_i) = \mu \mid_{i=1}^N \quad (1)$$

$$V(x_i) = \sigma^2 \mid_{i=1}^N \quad (2)$$

$$\text{Cov}(x_i, x_j) = 0 \mid_{i \neq j} \quad (3)$$

We want to derive the distribution of $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$

We know that any **linear combination** of independent normal random variables is also normally distributed and since it is easy to see that \bar{x} is a linear combination of $x_i \mid_{i=1}^N$, hence,

$$\bar{x} \sim N(\eta, \gamma^2) \quad (4)$$

where $\eta = E(\bar{x})$ and $\sigma^2 = V(\bar{x})$ are to be evaluated as follows

$$\begin{aligned} \eta &= E(\bar{x}) \\ &= E\left(\frac{1}{N} \sum_{i=1}^N x_i\right) \\ &= \frac{1}{N} \sum_{i=1}^N E(x_i) \\ &= \frac{1}{N} \sum_{i=1}^N \mu \quad (\text{from (1)}) \\ &= \mu \end{aligned} \quad (5)$$

$$\begin{aligned}
\gamma^2 &= V(\bar{x}) \\
&= V\left(\frac{1}{N} \sum_{i=1}^N x_i\right) \\
&= \frac{1}{N^2} \left\{ \sum_{i=1}^N V(x_i) + \sum_{i \neq j} \text{Cov}(x_i, x_j) \right\} \\
&= \frac{1}{N^2} \left\{ \sum_{i=1}^N \sigma^2 + 0 \right\} && \text{(from (2) and (3))} \\
&= \frac{\sigma^2}{N} && (6)
\end{aligned}$$

Hence, using (4), (5) and (6), we obtain,

$$\bar{x} \sim N\left(\mu, \frac{\sigma^2}{N}\right) \quad (7)$$

This makes intuitive sense since we expect the sample mean to be concentrated near μ and as our sample size (N) increases, we become more and more certain, i.e, $\sigma^2/N \rightarrow 0$ as $N \rightarrow \infty$

Student Name: Vishweshwar Tyagi

Roll Number: 191173

Date: March 1, 2021

Solution 4

We are given that $\mathbf{x}^m = (x_1^m, x_2^m \dots x_{N_m}^m)$ satisfies for each $m \in \{1, 2 \dots M\}$, the following

$$x_j^m \mid \mu_m \stackrel{i.i.d}{\sim} N(\mu_m, \sigma^2) \quad \forall j \in \{1, 2 \dots N_m\} \quad (1)$$

Further, it is given that

$$\mu_m \stackrel{i.i.d}{\sim} N(\mu_0, \sigma_0^2) \quad \forall m \in \{1, 2 \dots M\} \quad (2)$$

Letting $\bar{\mathbf{x}}^m = \frac{1}{N_m} \sum_{i=1}^{N_m} x_i^m$, we have from (7) of Question No. 3,

$$\bar{\mathbf{x}}^m \mid \mu_m \sim N(\mu_m, \sigma^2/N_m) \quad \forall m \in \{1, 2 \dots M\} \quad (3)$$

Hence, with $\mathbf{x} = (\mathbf{x}^1, \mathbf{x}^2 \dots \mathbf{x}^M)$ and $\boldsymbol{\mu} = (\mu_1, \mu_2 \dots \mu_M)$, we have from (3)

$$\begin{aligned} p(\mathbf{x} \mid \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2) &= \prod_{m=1}^M p(\bar{\mathbf{x}}^m \mid \mu_m, \sigma^2, \mu_0, \sigma_0^2) \\ &= \prod_{m=1}^M N(\bar{\mathbf{x}}^m \mid \mu_m, \sigma^2/N_m) \end{aligned} \quad (4)$$

Also, using (2), we get

$$p(\boldsymbol{\mu} \mid \mu_0, \sigma_0^2) = \prod_{m=1}^M N(\mu_m \mid \mu_0, \sigma_0^2) \quad (5)$$

Hence,

$$\begin{aligned} p(\boldsymbol{\mu} \mid \mathbf{x}, \mu_0, \sigma_0^2) &\propto p(\mathbf{x} \mid \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2) p(\boldsymbol{\mu} \mid \mu_0, \sigma_0^2) \\ &\propto \prod_{m=1}^M N(\bar{\mathbf{x}}^m \mid \mu_m, \sigma^2/N_m) \prod_{m=1}^M N(\mu_m \mid \mu_0, \sigma_0^2) \\ &\propto \prod_{m=1}^M N(\bar{\mathbf{x}}^m \mid \mu_m, \sigma^2/N_m) N(\mu_m \mid \mu_0, \sigma^2) \\ &= \prod_{m=1}^M N(\mu_m \mid \mu_{N_m}, \sigma_{N_m}^2) \end{aligned} \quad (6)$$

where

$$\frac{1}{\sigma_{N_m}^2} = \frac{1}{\sigma_0^2} + \frac{N_m}{\sigma^2} \quad (7)$$

$$\mu_{N_m} = \frac{\sigma^2}{N_m \sigma_0^2 + \sigma^2} \mu_0 + \frac{N_m \sigma_0^2}{N_m \sigma_0^2 + \sigma^2} \bar{\mathbf{x}}^m \quad (8)$$

Therefore, the required posterior for μ_m is given by

$$N(\mu_m \mid \frac{\sigma^2}{N_m \sigma_0^2 + \sigma^2} \mu_0 + \frac{N_m \sigma_0^2}{N_m \sigma_0^2 + \sigma^2} \bar{\mathbf{x}}^m, (\frac{1}{\sigma_0^2} + \frac{N_m}{\sigma^2})^{-1}) \quad (9)$$

The marginal likelihood is given by

$$\begin{aligned} p(\mathbf{x} \mid \sigma^2, \mu_0, \sigma_0^2) &= \int_{\boldsymbol{\mu}} p(\mathbf{x} \mid \boldsymbol{\mu}, \sigma^2, \mu_0, \sigma_0^2) p(\boldsymbol{\mu} \mid \mu_0, \sigma_0^2) d\boldsymbol{\mu} \\ &= \prod_{m=1}^M \int_{\mu_m} p(\bar{\mathbf{x}}^m \mid \mu_m, \sigma^2, \mu_0, \sigma_0^2) p(\mu_m \mid \mu_0, \sigma_0^2) d\mu_m \\ &= \prod_{m=1}^M \int_{\mu_m} N(\bar{\mathbf{x}}^m \mid \mu_m, \sigma^2/N_m) N(\mu_m \mid \mu_0, \sigma_0^2) d\mu_m \\ &= \prod_{m=1}^M N(\bar{\mathbf{x}}^m \mid \mu_0, \sigma_0^2 + \sigma^2/N_m) \end{aligned} \quad (10)$$

Now, in order to use MLE 2 for estimation of μ_0 , we consider the log of marginal likelihood and try to maximise it w.r.t μ_0

$$\begin{aligned} \ln p(\mathbf{x} \mid \sigma^2, \mu_0, \sigma_0^2) &= \ln \prod_{m=1}^M N(\bar{\mathbf{x}}^m \mid \mu_0, \sigma_0^2 + \sigma^2/N_m) \\ &= \sum_{m=1}^M \ln N(\bar{\mathbf{x}}^m \mid \mu_0, \sigma_0^2 + \sigma^2/N_m) \\ &\propto - \sum_{m=1}^M \frac{(\bar{\mathbf{x}}^m - \mu_0)^2}{\sigma_0^2 + \sigma^2/N_m} = g(\mu_0) \quad (\text{say}) \end{aligned}$$

Then,

$$\begin{aligned} \nabla_{\mu_0} g(\mu_0) &= 0 \\ \rightarrow \sum_{m=1}^M \frac{2(\bar{\mathbf{x}}^m - \mu_0)}{\sigma_0^2 + \sigma^2/N_m} &= 0 \\ \rightarrow \mu_0 \sum_{m=1}^M \frac{1}{\sigma_0^2 + \sigma^2/N_m} &= \sum_{m=1}^M \frac{\bar{\mathbf{x}}^m}{\sigma_0^2 + \sigma^2/N_m} \end{aligned}$$

Observing that, $\nabla^2 g(\mu_0) = -2 \sum_{m=1}^M \frac{1}{\sigma_0^2 + \sigma^2/N_m} < 0$, we get the required estimate for μ_0 as

$$\hat{\mu}_0|_{MLE\ 2} = \left(\sum_{m=1}^M \frac{\bar{\mathbf{x}}^m}{\sigma_0^2 + \sigma^2/N_m} \right) / \left(\sum_{m=1}^M \frac{1}{\sigma_0^2 + \sigma^2/N_m} \right) \quad (11)$$

Using the MLE 2 estimate for μ_0 in posterior of μ_m by substituting (11) in (9), we observe the resultant posterior mean of μ_m is

$$\frac{\sigma^2}{N_m\sigma_0^2 + \sigma^2} \left(\sum_{m=1}^M \frac{\bar{\mathbf{x}}^m}{\sigma_0^2 + \sigma^2/N_m} \right) / \left(\sum_{m=1}^M \frac{1}{\sigma_0^2 + \sigma^2/N_m} \right) + \frac{N_m\sigma_0^2}{N_m\sigma_0^2 + \sigma^2} \bar{\mathbf{x}}^m \quad (12)$$

The benefit of using this, in comparison with the posterior mean of (9), is that the posterior mean in (12) incorporates the data from all M schools whereas the one in (9) doesn't.

Student Name: Vishweshwar Tyagi
Roll Number: 191173
Date: March 1, 2021

Solution 5

We have for $\mathbf{y}^m \in \mathbb{R}^{N_m \times 1}$ and $X^m \in \mathbb{R}^{N_m \times D}$

$$p(\mathbf{y}^m | X^m, \mathbf{w}_m) = N(\mathbf{y}^m | X^m \mathbf{w}_m, \beta^{-1} I_{N_m}) \quad \forall m \in \{1, 2 \dots M\} \quad (1)$$

Hence, with \mathbf{y} , X , and \mathbf{w} denoting the collective responses, inputs and weight vectors respectively, we have,

$$\begin{aligned} p(\mathbf{y} | X, \mathbf{w}) &= \prod_{m=1}^M p(\mathbf{y}^m | X^m, \mathbf{w}_m) \\ &= \prod_{m=1}^M N(\mathbf{y}^m | X^m \mathbf{w}_m, \beta^{-1} I_{N_m}) \end{aligned} \quad (2)$$

Further, it is also given that $p(\mathbf{w}_d) = N(\mathbf{w}_d | \mathbf{w}_0, \lambda^{-1} I_D)$, therefore,

$$\begin{aligned} p(\mathbf{w}) &= \prod_{m=1}^M p(\mathbf{w}_d) \\ &= \prod_{m=1}^M N(\mathbf{w}_d | \mathbf{w}_0, \lambda^{-1} I_D) \end{aligned} \quad (3)$$

We will first find the marginal likelihood $p(\mathbf{y} | X)$ in order to arrive at the required MLE 2 objective

$$\begin{aligned} p(\mathbf{y} | X) &= \int_{\mathbf{w}} p(\mathbf{y} | X, \mathbf{w}) p(\mathbf{w}) d\mathbf{w} \\ &= \prod_{m=1}^M \int_{\mathbf{w}^m} N(\mathbf{y}^m | X^m \mathbf{w}_m, \beta^{-1} I_{N_m}) N(\mathbf{w}_m | \mathbf{w}_0, \lambda^{-1} I_D) d\mathbf{w}^m \\ &= \prod_{m=1}^M N(\mathbf{y}^m | X^m \mathbf{w}_0, \lambda^{-1} X^m X^{mT} + \beta^{-1} I_{N_m}) \end{aligned}$$

Further,

$$\begin{aligned} \ln p(\mathbf{w} | \mathbf{y}, X) &= \ln \prod_{m=1}^M N(\mathbf{y}^m | X^m \mathbf{w}_0, \lambda^{-1} X^m X^{mT} + \beta^{-1} I_{N_m}) \\ &= \sum_{m=1}^M \ln N(\mathbf{y}^m | X^m \mathbf{w}_0, \lambda^{-1} X^m X^{mT} + \beta^{-1} I_{N_m}) \end{aligned} \quad (4)$$

Thus, the required MLE 2 objective is,

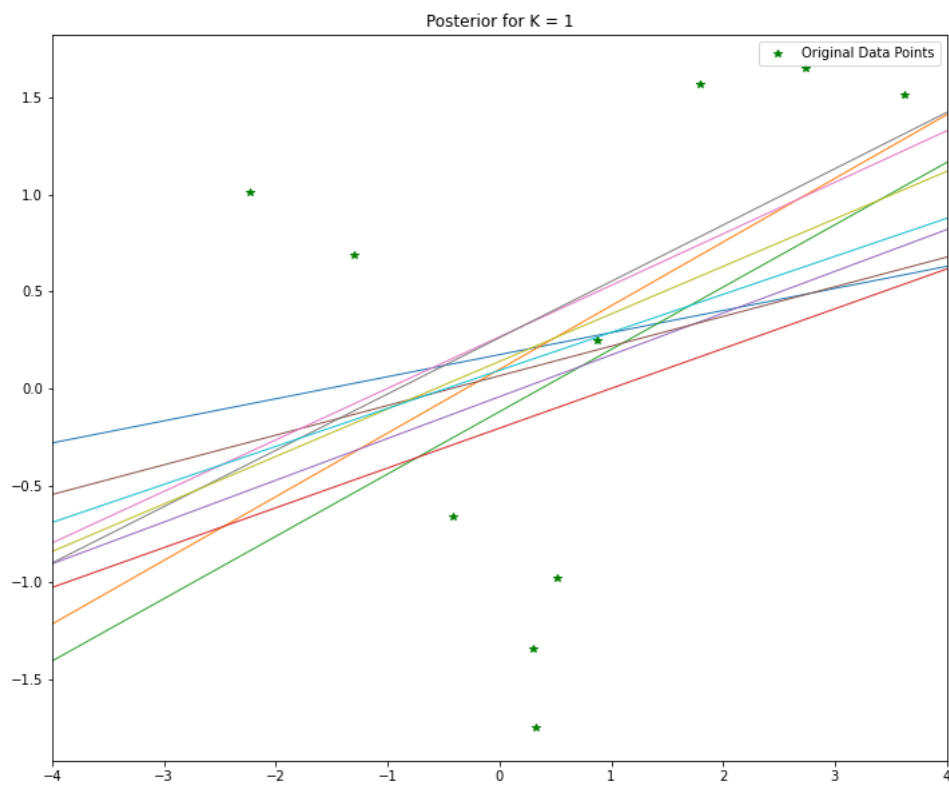
$$\hat{\boldsymbol{w}}_0|_{MLE} = \arg \max_{\boldsymbol{w}_0} \sum_{m=1}^M \ln N(\boldsymbol{y}^m | X^m \boldsymbol{w}_0, \lambda^{-1} X^m X^{mT} + \beta^{-1} I_{N_m}) \quad (5)$$

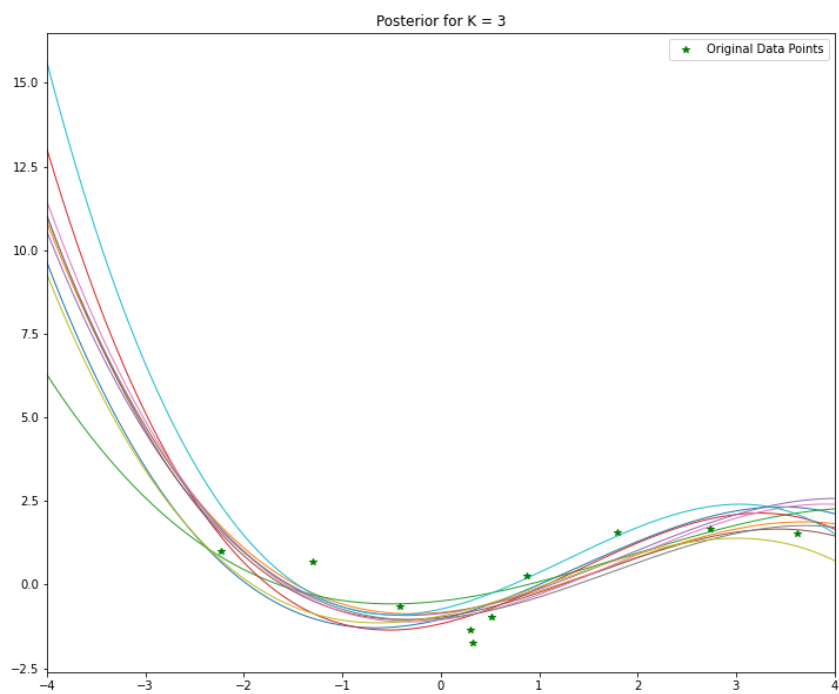
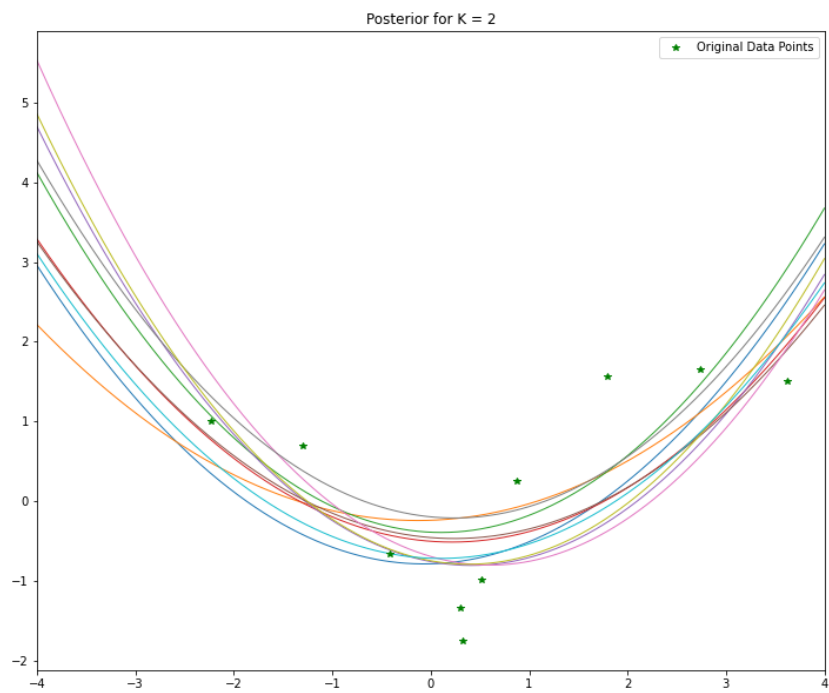
This has obvious benefit over fixing \boldsymbol{w}_0 to some value since this approach incorporates the knowledge of complete data, i.e, data from all schools, for estimation of \boldsymbol{w}_0 . Once having derived the posterior $p(\boldsymbol{w} | \boldsymbol{y}, X)$, we can plug in estimation of \boldsymbol{w}_0 , i.e, $\hat{\boldsymbol{w}}_0|_{MLE}$.

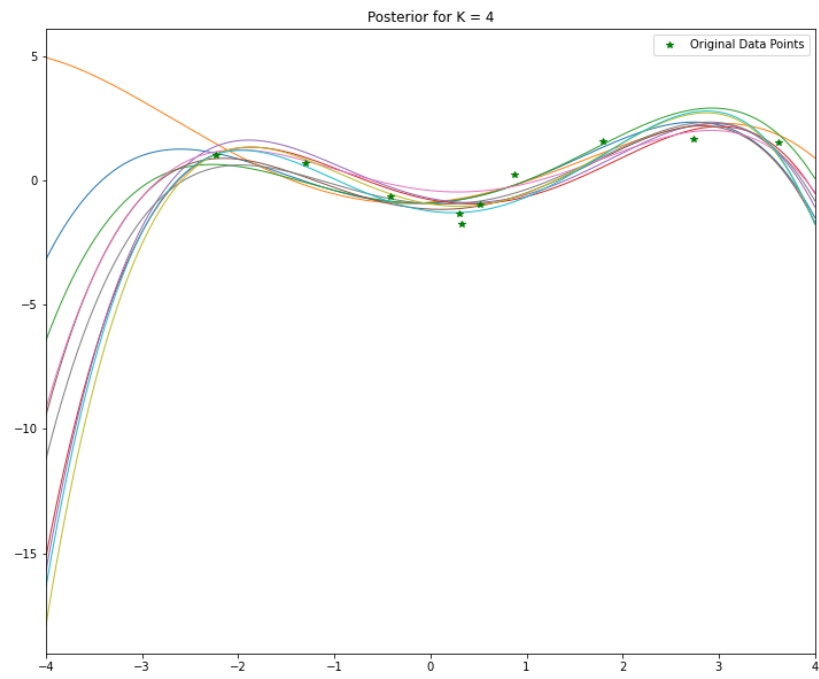
Student Name: Vishweshwar Tyagi
Roll Number: 191173
Date: March 1, 2021

Solution 6

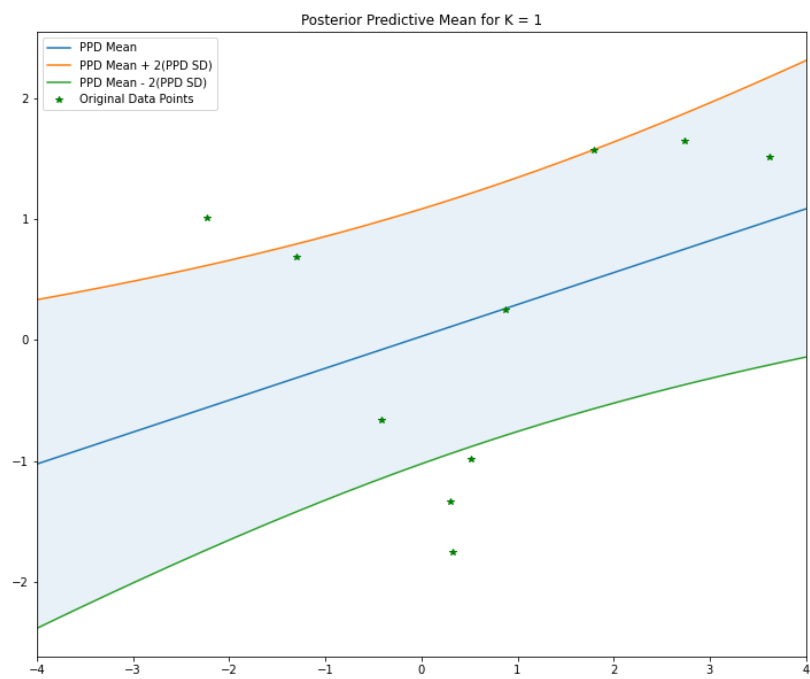
a)

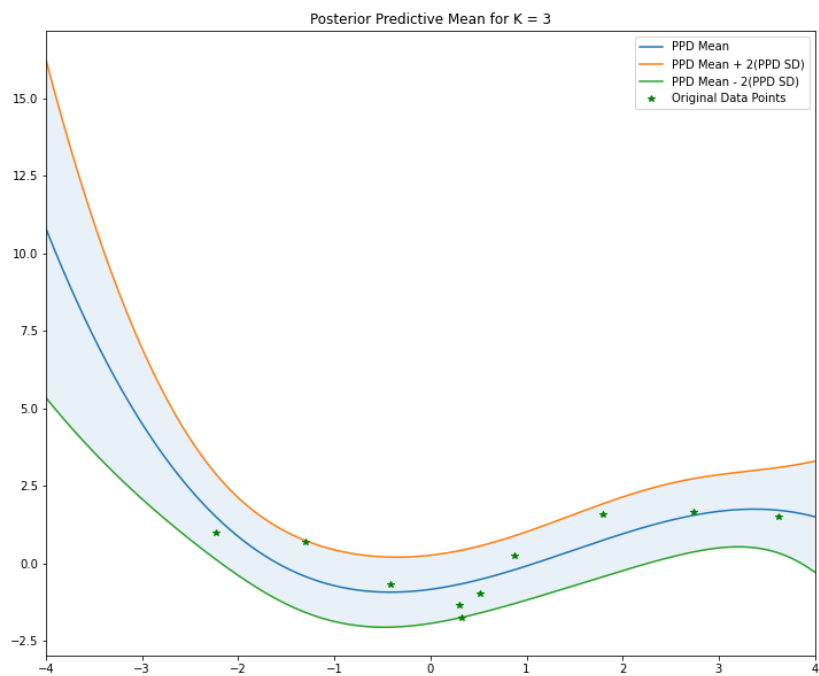
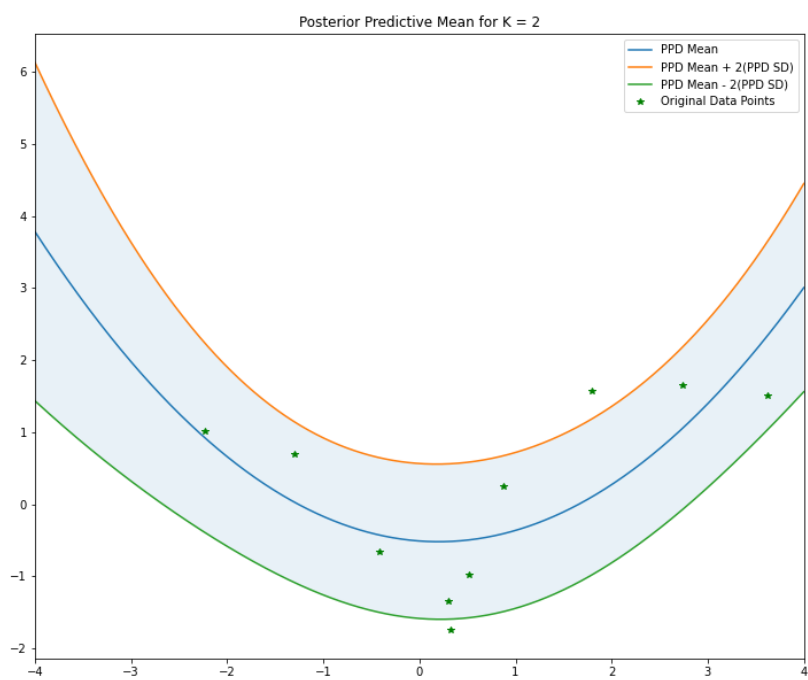


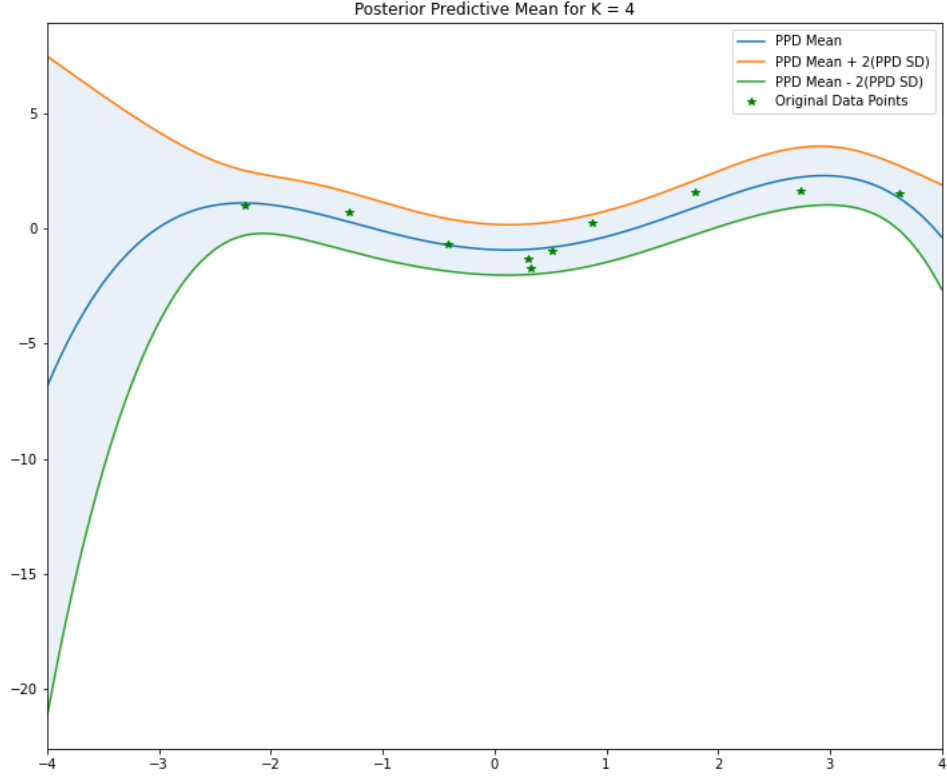




b)







c) We have the marginal likelihood given by

$$p(\mathbf{y} \mid \phi(\mathbf{x})) = N(\mathbf{y} \mid \mathbf{0}, \phi(\mathbf{x})\phi(\mathbf{x})^T + \beta^{-1}I_M)$$

where M denotes the total number of training examples.

We obtained the following values for log marginal likelihood

For $K = 1$, the log marginal likelihood is -32.3520152804452

For $K = 2$, the log marginal likelihood is -22.772153178782283

For $K = 3$, the log marginal likelihood is -22.079070642241916

For $K = 4$, the log marginal likelihood is -22.386776180349212

It is easy to see that the model with $\mathbf{K} = \mathbf{3}$ has the highest log marginal likelihood and therefore seems to best explain the data.

d) The posterior mean for \mathbf{w} is given by

$$(\phi(\mathbf{x})\phi(\mathbf{x})^T + \beta^{-1}I_D)^{-1}X^T\mathbf{y}$$

where $D = K + 1$ is the dimension of ϕ map.

Now, the MAP estimate for \mathbf{w} is given by the mode of posterior but since we know that the posterior follows normal distribution due to conjugacy, the mode is equal to the mean and hence,

$$\hat{\mathbf{w}}|_{MAP} = (\phi(\mathbf{x})\phi(\mathbf{x})^T + \beta^{-1}I_D)^{-1}X^T\mathbf{y}$$

We obtained the following values for log likelihood using **MAP** estimate for the weight vector $\hat{\mathbf{w}}_{MAP}$

For $K = 1$, the log likelihood is -28.094004379075553

For $K = 2$, the log likelihood is -15.360663659052214

For $K = 3$, the log likelihood is -10.935846883615742

For $K = 4$, the log likelihood is -7.225291259028564

Using the **MAP** estimate, we see that the log likelihood of model with $\mathbf{K} = 4$ is the highest. This result doesn't agree with what we obtained above using the log marginal likelihood.

Using log marginal likelihood in order to select the best model would be a better criteria as compared to only using the log likelihood since the former doesn't rely on point based estimate (MAP estimate for weight vector in this case $\hat{\mathbf{w}}_{MAP}$) and instead does prior averaging over the weight vector.

e) It can be observed from the graph 'Posterior Predictive Mean for $K = 3$ ', the 2σ uncertainty gap seems to be the widest for $x \in [-4, -3]$ which makes sense as it lacks training examples in this region.

Hence, I would prefer $x' \in [-4, -3]$ in order to improve the model by reducing the uncertainty gap in this region.