

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: February 26, 2021

(When You Integrate Out..) Given

$$p(x|\eta) = N(x|0, \eta) = \frac{1}{\sqrt{2\pi\eta}} \exp\left(-\frac{x^2}{2\eta}\right)$$

and

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

Now calculating $p(x|\gamma)$

$$p(x|\gamma) = \int_0^\infty p(x|\eta)p(\eta|\gamma)d\eta \quad (1)$$

$$= \int_0^\infty \frac{1}{\sqrt{2\pi\eta}} \exp\left(-\frac{x^2}{2\eta}\right) \frac{\gamma^2}{2} \exp\left(-\frac{\gamma^2\eta}{2}\right) d\eta \quad (2)$$

$$= \int_0^\infty \frac{\gamma^2}{2\sqrt{2\pi\eta}} \exp\left(-\frac{x^2}{2\eta} - \frac{\gamma^2\eta}{2}\right) d\eta \quad (3)$$

To calculate the above integarm lets calculate its (MGF) Moment generating function of the above equation.

$$\mathbf{M_x}(t) = \mathbb{E}[\exp(tx)] = \int_{-\infty}^\infty \exp(tx) \mathbf{f_x}(x) dx$$

By substituting $\mathbf{f_x}(x)$ with the equation 3 we get

$$\mathbf{M_x}(t) = \int_{-\infty}^\infty \exp(tx) \int_0^\infty \frac{\gamma^2}{2\sqrt{2\pi\eta}} \exp\left(-\frac{x^2}{2\eta} - \frac{\gamma^2\eta}{2}\right) d\eta \quad dx \quad (4)$$

$$= \int_0^\infty \frac{\gamma^2}{2\sqrt{2\pi\eta}} \int_{-\infty}^\infty \exp\left(-\frac{x^2}{2\eta} + tx - \frac{\gamma^2\eta}{2}\right) dx \quad d\eta \quad (5)$$

As we know from the *ref*

$$\int_{-\infty}^\infty \exp(-ax^2 + bx + c) dx = \sqrt{\frac{\pi}{a}} \exp\left(\frac{b^2}{4a} + c\right) \quad (6)$$

By using the above formula equation 5 will evaluate to

$$\begin{aligned} \mathbf{M_x}(t) &= \int_0^\infty \frac{\gamma^2}{2\sqrt{2\pi\eta}} \left\{ \sqrt{2\eta\pi} \exp\left(\frac{t^2\eta}{2} - \frac{\gamma^2\eta}{2}\right) \right\} d\eta \\ &= \frac{\gamma^2}{2} \int_0^\infty \exp\left(\left(\frac{t^2 - \gamma^2}{2}\right)\eta\right) d\eta \\ &= \frac{1}{1 - \frac{t^2}{\gamma^2}} \end{aligned} \quad t^2 < \gamma^2 \text{ and } \gamma > 0 \implies |t| < \gamma$$

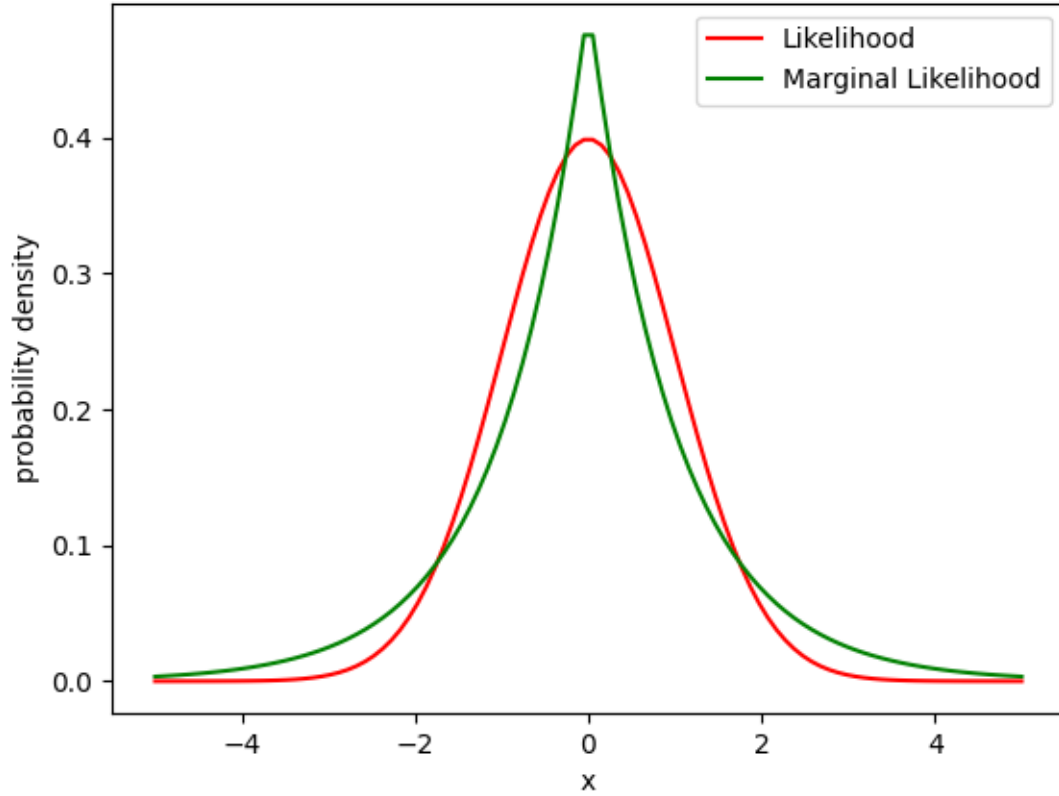


Figure 1: Likelihood and Marginal Likelihood

By comparing the above MGF with the reference table we found out that its the MGF of the Laplace distribution with $\mu = 0$ and $b = \frac{1}{\gamma}$ where $Laplace(\mu, b) = \frac{1}{2b} \exp(-\frac{|x-\mu|}{b})$ Therefore

$$p(x|\gamma) = Laplace(0, \gamma^{-1})$$

The marginal likelihood means that its the likelihood weighted averaged over all of its parameters.

The plots 1 for the $P(x|\eta)$ and $p(x|\gamma)$ considering $\eta = 1$ and $\gamma = 1$

we can observe that marginal likelihood is peaked around the zero(mean) than the likelihood

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: February 26, 2021

(It Gets Better..) Bayesian linear regression model
 likelihood :

$$p(y|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^T \mathbf{x}, \beta^{-1})$$

prior :

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I})$$

posterior:

$$p(y_*|\mathbf{x}_*) = \mathcal{N}(\boldsymbol{\mu}_N^T \mathbf{x}_*, \beta^{-1} + \mathbf{x}_*^T \boldsymbol{\Sigma}_N \mathbf{x}_*) = \mathcal{N}(\boldsymbol{\mu}_N^T \mathbf{x}_*, \sigma_N^2(x_*))$$

$$\boldsymbol{\mu}_N = \boldsymbol{\Sigma}_N (\beta \sum_{n=1}^N y_n \mathbf{x}_n)$$

$$\begin{aligned} \boldsymbol{\Sigma}_N &= (\beta \boldsymbol{\Sigma}_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T + \lambda \mathbf{I})^{-1} \\ &= \frac{1}{\beta} (\boldsymbol{\Sigma}_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T + \frac{\lambda}{\beta} \mathbf{I})^{-1} \\ &= \frac{1}{\beta} M^{-1} \end{aligned}$$

$$\text{where } M = \boldsymbol{\Sigma}_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T + \frac{\lambda}{\beta} \mathbf{I}$$

Here N denotes number of training data points.

We know that $\boldsymbol{\Sigma}_N$ is positive semidefinite matrix i.e $\forall \mathbf{x}_*, \mathbf{x}_*^T \boldsymbol{\Sigma}_N \mathbf{x}_* > 0$

Lets focus on variance how it changes with more data :

consider the following a new point \mathbf{v} added to the dataset then the new variance is

$$\begin{aligned} \boldsymbol{\Sigma}_{N+1} &= \frac{1}{\beta} (\boldsymbol{\Sigma}_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T + \mathbf{v} \mathbf{v}^T + \frac{\lambda}{\beta} \mathbf{I})^{-1} \\ \boldsymbol{\Sigma}_{N+1} &= \frac{1}{\beta} (M + \mathbf{v} \mathbf{v}^T)^{-1} \end{aligned}$$

To find the relation between the variances lets calculate

$$\begin{aligned} \sigma_N^2(\mathbf{x}_*) - \sigma_{N+1}^2(\mathbf{x}_*) &= (\mathbf{x}_*^T) \beta^{-1} M^{-1} (\mathbf{x}_*) - (\mathbf{x}_*^T) \beta^{-1} (M + \mathbf{v} \mathbf{v}^T)^{-1} (\mathbf{x}_*) \\ &= \beta^{-1} (\mathbf{x}_*^T) M^{-1} (\mathbf{x}_*) - \beta^{-1} (\mathbf{x}_*^T) M^{-1} (\mathbf{x}_*) + \beta^{-1} \frac{(\mathbf{x}_*^T) (\mathbf{M}^{-1} \mathbf{v}) (\mathbf{v}^T \mathbf{M}^{-1}) (\mathbf{x}_*)}{1 + \mathbf{v}^T \mathbf{M}^{-1} \mathbf{v}} \\ &= \beta^{-1} \frac{(\mathbf{x}_*^T \mathbf{M}^{-1} \mathbf{v}) (\mathbf{x}_*^T \mathbf{M}^{-1} \mathbf{v})^T}{1 + \mathbf{v}^T \mathbf{M}^{-1} \mathbf{v}} \\ &= \beta^{-1} \frac{\|\mathbf{x}_*^T \mathbf{M}^{-1} \mathbf{v}\|_2^2}{1 + \mathbf{v}^T \mathbf{M}^{-1} \mathbf{v}} \\ &\geq 0 \quad \text{as } l_2 \text{ norm is positive and } M^{-1} \text{ is PSD and } \beta^{-1} > 0 \\ \sigma_{N+1}^2(\mathbf{x}_*) &\leq \sigma_N^2(\mathbf{x}_*) \end{aligned}$$

Therefore the variance of the predictive posterior decreases or remains same as the number of data points increases.

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: February 26, 2021

(Distribution of Empirical Mean of Gaussian Observations)

Given

$$x_1 \dots x_N \sim \mathcal{N}(\mu, \sigma^2)$$

and its empirical mean as

$$\begin{aligned} \bar{x} &= \frac{1}{N} \sum_{n=1}^N x_n \\ &= \left[\frac{1}{N}, \dots, \frac{1}{N} \right]^T [x_1 \dots x_N] \\ &= \mathbf{a}^T \mathbf{b} \end{aligned} \quad \begin{aligned} \mathbf{a} &= \left[\frac{1}{N}, \dots, \frac{1}{N} \right] \\ \mathbf{b} &= [x_1 \dots x_N] \end{aligned}$$

As $x_1 \dots x_N \sim \mathcal{N}(\mu, \sigma^2)$ and are iid \mathbf{b} is also a multi dimensional random variable with $\boldsymbol{\mu} = [\mu \dots \mu]^T$ and $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$

since the \bar{x} is linear transformation of \mathbf{b} Gaussian random variable it is also a Gaussian with mean and co-variance as follows

$$\begin{aligned} \mathbb{E}[\bar{x}] &= \mathbb{E}[\mathbf{a}^T \mathbf{b}] \\ &= \mathbf{a}^T \mathbb{E}[\mathbf{b}] \\ &= \mathbf{a}^T \boldsymbol{\mu} \\ &= \sum_{n=1}^N \frac{1}{N} \mu \\ &= \mu \\ \text{var}(\bar{x}) &= \text{var}(\mathbf{a}^T \mathbf{b}) \\ &= \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a} \\ &= \sigma^2 \mathbf{a}^T \mathbf{a} \\ &= \frac{\sigma^2}{N} \end{aligned}$$

The distribution of \bar{x} is $\mathcal{N}(\mu, \frac{\sigma^2}{N})$

The result seems to be obvious that as the \bar{x} is mean, whose values are taken from the normal distribution. so the mean of the \bar{x} distribution should also be μ

\bar{x} also depends on the number of values taken i.e N so our variance of \bar{x} decreases as N increases. which mean that the \bar{x} values will be close to μ as the N values increases which reflects in our results also.

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: February 26, 2021

(Benefits of Probabilistic Joint Modeling-1)

Consider a dataset of test-scores of students from M schools in a district:

$\mathbf{X} : \{\mathbf{x}^{(m)}\}_{m=1}^M = \{x_1^{(m)} \dots x_{N_m}^{(m)}\}_{m=1}^M$ and $x_n^{(m)} \sim \mathcal{N}(\mu_m, \sigma^2)$ are independent.
 and also the means $\mu_1 \dots \mu_M$ such that $\mu_m \sim \mathcal{N}(\mu_0, \sigma_0^2)$

Part 1 :

Assume the hyperparameters μ_0 and σ_0^2 to be known
 posterior distribution of μ_m

$$\begin{aligned} p(\mu_m | \mathbf{x}^{(m)}, \mu_0, \sigma^2) &= \frac{p(\mathbf{x}^{(m)} | \mu_m, \sigma^2) p(\mu_m | \mu_0, \sigma_0^2)}{p(\mathbf{x}^{(m)} | \mu_0, \sigma_0^2)} \\ &= \frac{p(\mathbf{x}^{(m)} | \mu_m, \sigma^2) p(\mu_m | \mu_0, \sigma_0^2)}{\int p(\mathbf{x}^{(m)} | \mu_m, \sigma^2) p(\mu_m | \mu_0, \sigma_0^2)} \\ &\propto \prod_{n=1}^{N_m} \mathcal{N}(x_n^{(m)} | \mu_m, \sigma^2) \mathcal{N}(\mu_m | \mu_0, \sigma_0^2) \\ &= \exp\left(-\frac{\sum_{n=1}^{N_m} (x_n^{(m)} - \mu_m)^2}{2\sigma^2} - \frac{(\mu_m - \mu_0)^2}{2\sigma_0^2}\right) \\ &= \exp\left(-\frac{1}{2}\left(\mu_m^2\left(\frac{N_m}{\sigma^2} + \frac{1}{\sigma_0^2}\right) - \mu_m\left(\frac{2N_m\bar{x} - N_m}{\sigma^2} - \frac{2\mu_0}{\sigma_0^2}\right) + \left(\frac{\sum_{n=1}^{N_m} x_n^2}{\sigma^2} + \frac{\mu_0^2}{\sigma_0^2}\right)\right)\right) \end{aligned}$$

by completing the squares and bring it to normal form

$$p(\mu_m | \mathbf{x}^{(m)}, \mu_0, \sigma^2) = \mathcal{N}(\mu_m | \mu_{H_m}, \sigma_{H_m}^2)$$

where

$$\begin{aligned} \frac{1}{\sigma_{H_m}^2} &= \frac{1}{\sigma_0^2} + \frac{N_m}{\sigma^2} \\ \mu_{H_m} &= \frac{\sigma^2 \mu_0 + N_m \sigma_0^2 \bar{x}^{(m)}}{\sigma^2 + N_m \sigma_0^2} \quad ; \quad \bar{x}^{(m)} = \frac{\sum_{n=1}^{N_m} x_n^{(m)}}{N_m} \end{aligned}$$

Part 2 : Assume μ_0 to be unknown and σ_0^2 known.

Marginal Likelihood

$$\begin{aligned}
p(\mathbf{X}|\mu_0, \sigma_0^2, \sigma^2) &= \int p(\mathbf{X}|\mu, \sigma^2) p(\mu|\mu_0, \sigma_0^2) d\mu \\
&= \prod_{m=1}^M \int_{-\infty}^{\infty} p(\mathbf{x}^{(m)}|\mu_m, \sigma^2) p(\mu|\mu_0, \sigma_0^2) d\mu_m \\
&= \prod_{m=1}^M \int_{-\infty}^{\infty} \prod_{n=1}^{N_m} p(\mathbf{x}_n^{(m)}|\mu_m, \sigma^2) p(\mu|\mu_0, \sigma_0^2) d\mu_m \\
&= \prod_{m=1}^M \int_{-\infty}^{\infty} \prod_{n=1}^{N_m} \mathcal{N}(\mathbf{x}_n^{(m)}|\mu_m, \sigma^2) \mathcal{N}(\mu|\mu_0, \sigma_0^2) d\mu_m \\
&= \prod_{m=1}^M \int_{-\infty}^{\infty} \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^{N_m} \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{\sum_{i=1}^{N_m} (x_i^{(m)} - \mu_m)^2}{2\sigma^2} - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) d\mu_m \\
&= \prod_{m=1}^M \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^{N_m} \frac{1}{\sqrt{2\pi\sigma_0^2}} \int_{-\infty}^{\infty} \exp\left(-\left(\frac{N_m}{2\sigma^2} + \frac{1}{2\sigma_0^2}\right)\mu_m^2 \right. \\
&\quad \left. + \left(\frac{\bar{x}^{(m)}N_m}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)\mu_m \right. \\
&\quad \left. + \left(-\frac{N_m\bar{x}^{(m)}}{2\sigma^2} - \frac{\mu_0^2}{2\sigma_0^2}\right)\right) d\mu_m
\end{aligned}$$

using the equation 6

$$= \prod_{m=1}^M \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^{N_m} \frac{1}{\sqrt{2\pi\sigma_0^2}} \sqrt{\frac{\pi}{\left(\frac{N_m}{2\sigma^2} + \frac{1}{2\sigma_0^2}\right)}} \exp\left(\frac{\left(\frac{\bar{x}^{(m)}N_m}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)^2}{4\left(\frac{N_m}{2\sigma^2} + \frac{1}{2\sigma_0^2}\right)} - \frac{N_m\bar{x}^{(m)}}{2\sigma^2} - \frac{\mu_0^2}{2\sigma_0^2}\right)$$

Now derivation for MLE-II finding

$$\begin{aligned}
&\arg \max_{\mu_0} \log p(\mathbf{X}|\mu_0, \sigma_0^2, \sigma^2) \\
&= \arg \max_{\mu_0} \log \left[\prod_{m=1}^M \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^{N_m} \frac{1}{\sqrt{2\pi\sigma_0^2}} \sqrt{\frac{\pi}{\left(\frac{N_m}{2\sigma^2} + \frac{1}{2\sigma_0^2}\right)}} \exp\left(\frac{\left(\frac{\bar{x}^{(m)}N_m}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)^2}{4\left(\frac{N_m}{2\sigma^2} + \frac{1}{2\sigma_0^2}\right)} - \frac{N_m\bar{x}^{(m)}}{2\sigma^2} - \frac{\mu_0^2}{2\sigma_0^2}\right) \right] \\
&= \arg \max_{\mu_0} \sum_{m=1}^M \log \left[C + \left(\frac{\left(\frac{\bar{x}^{(m)}N_m}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)^2}{2\left(\frac{N_m}{\sigma^2} + \frac{1}{\sigma_0^2}\right)} - \frac{\mu_0^2}{2\sigma_0^2}\right) \right]
\end{aligned}$$

where C is constant wrt μ_0 Now derivating wrt μ_0 and equating to 0

$$\begin{aligned}
&\sum_{m=1}^M \frac{\left(\frac{\bar{x}^{(m)}N_m}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)}{\left(\frac{N_m}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\sigma_0^2} - \frac{\mu_0}{\sigma_0^2} = 0 \\
&\sum_{m=1}^M \frac{\frac{\bar{x}^{(m)}N_m}{\sigma^2}}{\left(\frac{\sigma_0^2 N_m + \sigma^2}{\sigma^2 \sigma_0^2}\right)} + \frac{\mu_0}{\left(\frac{\sigma_0^2 N_m + \sigma^2}{\sigma^2 \sigma_0^2}\right)\sigma_0^2} - \mu_0 = 0 \\
&\sum_{m=1}^M \frac{\sigma_0^2 \bar{x}^{(m)} N_m + \sigma^2 \mu_0}{\sigma_0^2 N_m + \sigma^2} - \mu_0 = 0
\end{aligned}$$

$$\mu_0 = \frac{\sum_{m=1}^M \frac{\bar{x}^{(m)} N_m}{\sigma_0^2 N_m + \sigma^2}}{\sum_{m=1}^M \frac{N_m}{\sigma_0^2 N_m + \sigma^2}}$$

The above is the required MLE-II estimate of the μ_0

Part 3: Advantage of the using MLE-II estimate of μ_0 Once we substitute back our μ_0 resulted by MLE - II in our equation of part 1 we get

$$\mu_{H_m} = \frac{\sigma^2 \frac{\sum_{m=1}^M \frac{\bar{x}^{(m)} N_m}{\sigma_0^2 N_m + \sigma^2}}{\sum_{m=1}^M \frac{N_m}{\sigma_0^2 N_m + \sigma^2}} + N_m \sigma_0^2 \bar{x}^{(m)}}{\sigma^2 + N_m \sigma_0^2}$$

Now our posteriror mean depends on the data of all the schools which was previously depended only on the data of perticular school. It is kind of generalized mean now. And also if any schools having less data, then this estimate would be better.

We used training data to choose the best value of hyperparameter μ_0 using MLE -II. This kind of choosing hyperparameter helps us to take out some data for cross validation to choose the hyperparameters

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: February 26, 2021

(Benefits of Probabilistic Joint Modeling-2)

student data from M schools where N_m denotes the number of students in school m
 For student n in school m , the response variable $y_n^{(m)} \in \mathbb{R}$ and features $\mathbf{x}_n^{(m)} \in \mathbb{R}^D$
 Linear regression model for these scores

$$\begin{aligned} \text{likelihood:} \quad & p(\mathbf{y}^{(m)} | \mathbf{X}^{(m)}, \boldsymbol{\omega}_m) = \mathcal{N}(\mathbf{y}^{(m)} | \mathbf{X}^{(m)} \boldsymbol{\omega}_m, \beta^{-1} \mathbf{I}_N) \\ \text{prior:} \quad & p(\boldsymbol{\omega}_m) = \mathcal{N}(\boldsymbol{\omega}_m | \boldsymbol{\omega}_0, \lambda^{-1} \mathbf{I}_D) \end{aligned}$$

the expression for the log of the MLE-II objective for estimating $\boldsymbol{\omega}_0$

$$\begin{aligned} & \arg \max_{\boldsymbol{\omega}_0} \log P(\mathbf{Y} | \mathbf{X}, \boldsymbol{\omega}_0) \\ &= \arg \max_{\boldsymbol{\omega}_0} \sum_{m=1}^M \log \int_{\boldsymbol{\omega}_m} P(\mathbf{y}^{(m)} | \mathbf{X}^{(m)}, \boldsymbol{\omega}_m) P(\boldsymbol{\omega}_m | \boldsymbol{\omega}_0) d\boldsymbol{\omega}_m \\ &= \arg \max_{\boldsymbol{\omega}_0} \sum_{m=1}^M \log \int_{\boldsymbol{\omega}_m} \mathcal{N}(\mathbf{y}^{(m)} | \mathbf{X}^{(m)} \boldsymbol{\omega}_m, \beta^{-1} \mathbf{I}_N) \mathcal{N}(\boldsymbol{\omega}_m | \boldsymbol{\omega}_0, \lambda^{-1} \mathbf{I}_D) d\boldsymbol{\omega}_m \\ & \text{by using the linear gaussian model results} \\ & \int_{\boldsymbol{\omega}_m} \mathcal{N}(\mathbf{y}^{(m)} | \mathbf{X}^{(m)} \boldsymbol{\omega}_m, \beta^{-1} \mathbf{I}_N) \mathcal{N}(\boldsymbol{\omega}_m | \boldsymbol{\omega}_0, \lambda^{-1} \mathbf{I}_D) d\boldsymbol{\omega}_m \\ &= \mathcal{N}(\mathbf{y}^{(m)} | \mathbf{X}^{(m)} \boldsymbol{\omega}_0, \mathbf{X}^{(m)} \lambda^{-1} \mathbf{I}_D \mathbf{X}^{(m)T} + \beta^{-1} \mathbf{I}_{N_m}) \end{aligned}$$

so by using the above result our MLE - II objective becomes

$$\arg \max_{\boldsymbol{\omega}_0} \sum_{m=1}^M \log \mathcal{N}(\mathbf{y}^{(m)} | \mathbf{X}^{(m)} \boldsymbol{\omega}_0, \mathbf{X}^{(m)} \lambda^{-1} \mathbf{I}_D \mathbf{X}^{(m)T} + \beta^{-1} \mathbf{I}_{N_m})$$

If our goal is to learn the school-specific weight vectors w_1, \dots, w_M .

MLE - II method uses the data from all the schools to get the optimal value of the hyperparameter $\boldsymbol{\omega}_0$. which helps us to best fit of our model. And also if any schools having less data, then this estimate would be better. where as in case of fixing the value of $\boldsymbol{\omega}_0$ to some fixed value we have to compromise on the optimal value of $\boldsymbol{\omega}_0$

Bayesian Linear Regression

Part 1

For each k , computed the posterior of w and showing a plot with 10 random functions drawn from the inferred posterior along with original training examples

Part 2 we can see the them in 2 3 4 5 Computed the mean and mean ± 2 *std in figures 6 7 8

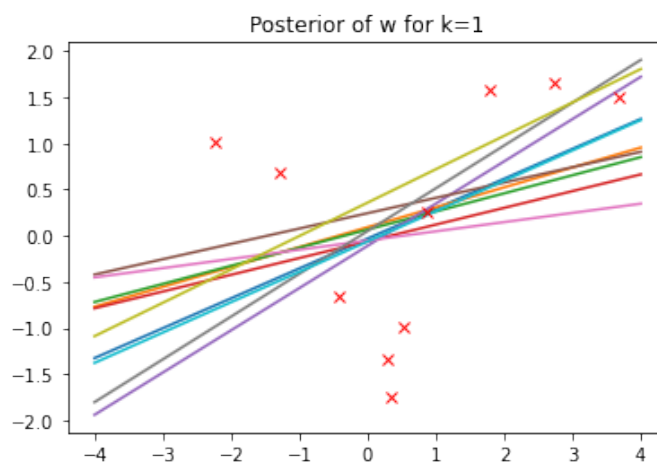


Figure 2: Posterior for $k = 1$

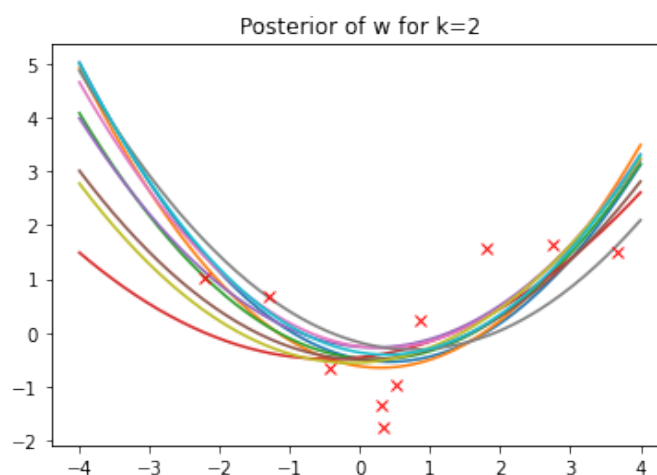


Figure 3: Posterior for $k = 2$

9 Part 3

Log marginal likelihood of data as follows.

$k=1$ the marginal likelihood value is $= -32.325577919349065$

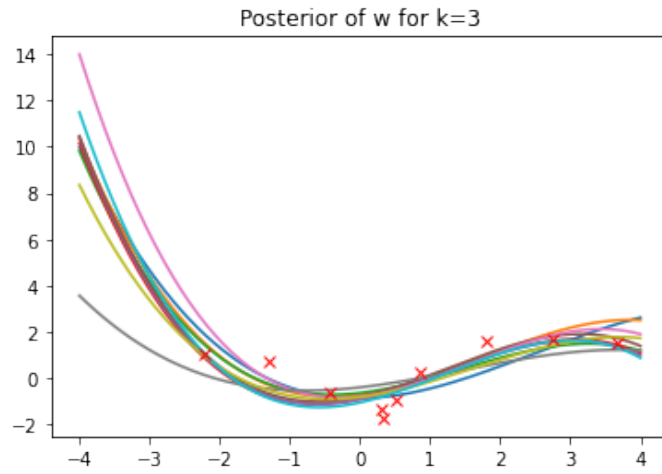


Figure 4: Posterior for $k = 3$

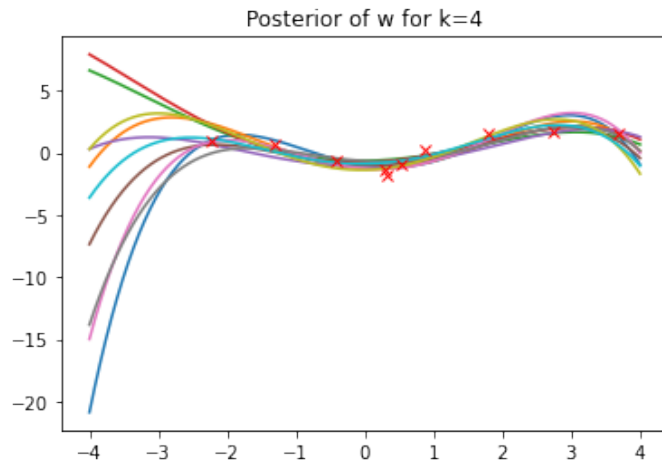


Figure 5: Posterior for $k = 4$

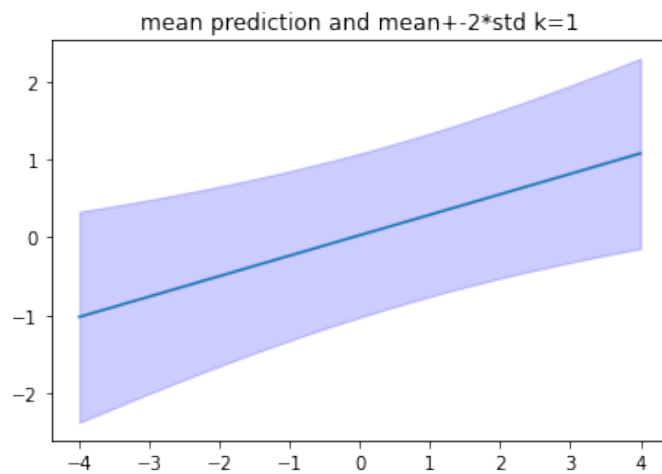


Figure 6: mean $\pm 2 \cdot \text{std}$ $k = 1$

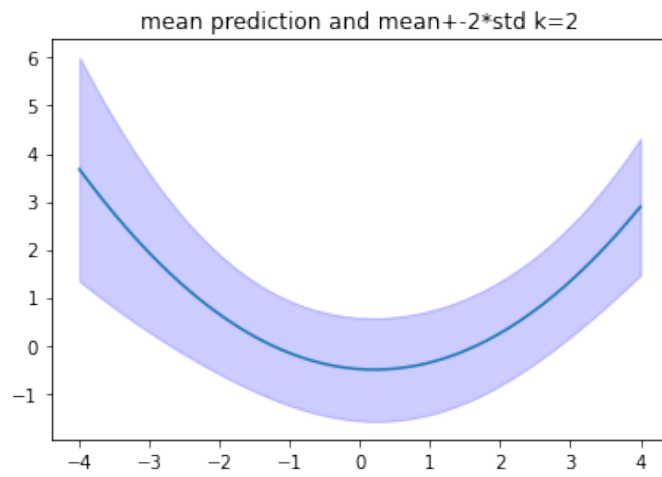


Figure 7: $\text{mean} \pm 2 \cdot \text{std}$ for $k = 2$

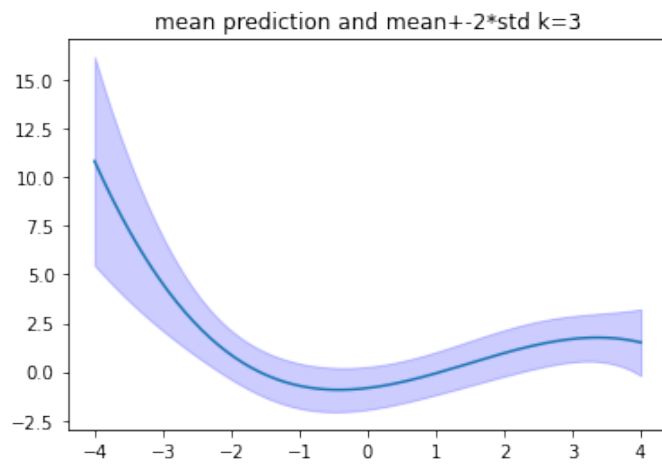


Figure 8: $\text{mean} \pm 2 \cdot \text{std}$ for $k = 3$

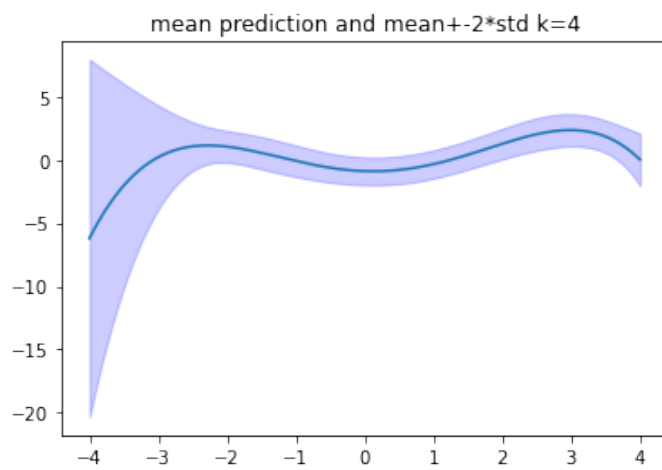


Figure 9: $\text{mean} \pm 2 \cdot \text{std}$ for $k = 4$

k= 2 the marginal likelihood value is = -23.126889656893887
k= 3 the marginal likelihood value is = -22.12287495542345
k= 4 the marginal likelihood value is = -22.605375642340185
k = 3 has highest log marginal likelihood, so it is the best model.

Part 4

The likelihood using map solution k= 1 the likelihood value is = -28.061373030860793

k= 2 the likelihood value is = -15.70039023251804

k= 3 the likelihood value is = -10.915688165826634

k= 4 the likelihood value is = -7.371110916347485

k= 4 is the best model as it is having highest value.

we can see that marginal likelihood and likelihood is not giving the same best model and it is appropriate to choose maximum marginal likelihood as the criteria for model selection as it will take is kind of averaged likelihood over all values of the ω

Part 5

Our best model is for k=3. As the maximum uncertainty in model for k=3 is between [-4,-3] from the plots of part-2, therefore, including an additional training input x' (along with its output y') in this region will improve the model by reducing the uncertainty.