

Student Name: Dishay Mehta

Roll Number: 200341

Date: March 30, 2023

$$a) \text{Gamma}(x|a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}$$

Let us approximate Gamma distribution to a Gaussian using Laplace's Approximation.

$\text{Gamma}(x|a, b) \approx \mathcal{N}(x_{map}, H^{-1})$ where x_{map} is the value of mode and H is the Hessian matrix. From slides,

$$x_{map} = \arg\max_x \text{Gamma}(x|a, b)$$

$$= \arg\max_x \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}$$

Differentiation wrt x and equating to 0,

$$0 = (a-1)x^{a-2}e^{-bx} + x^{a-1}(-b)e^{-bx}$$

$$\implies x_{map} = \frac{a-1}{b}$$

$$H = -\nabla^2 \log \text{Gamma}(x|a, b) \Big|_{x=\frac{a-1}{b}}$$

$$= -\nabla^2 (a \log b - \log \Gamma(a) + (a-1) \log x - bx) \Big|_{x=\frac{a-1}{b}}$$

$$= -\nabla \left(\frac{a-1}{x} - b \right) \Big|_{x=\frac{a-1}{b}}$$

$$= \left(\frac{a-1}{x^2} \right) \Big|_{x=\frac{a-1}{b}}$$

$$H = \frac{b^2}{a-1}$$

$$\implies H^{-1} = \frac{a-1}{b^2}$$

$$\text{Thus, } \text{Gamma}(x|a, b) \approx \mathcal{N}\left(\frac{a-1}{b}, \frac{a-1}{b^2}\right)$$

From the prob stats refresher slides,

$$\text{Mean of Gamma distribution } \text{Gamma}(x|a, b) = \frac{a}{b}$$

$$\text{Variance of Gamma distribution } \text{Gamma}(x|a, b) = \frac{a}{b^2}$$

Thus approximating $\text{Gamma}(x|a, b)$ by a Gaussian whose mean and variance are equal to the mean and variance of $\text{Gamma}(x|a, b)$ is $\text{Gamma}(x|a, b) \approx \mathcal{N}\left(\frac{a}{b}, \frac{a}{b^2}\right)$

Thus $\mathcal{N}\left(\frac{a-1}{b}, \frac{a-1}{b^2}\right)$ and $\mathcal{N}\left(\frac{a}{b}, \frac{a}{b^2}\right)$ will behave same if a tends to ∞ , since

$$\frac{a-1}{b} \approx \frac{a}{b} \text{ and } \frac{a-1}{b^2} \approx \frac{a}{b^2}$$

Hence, when a is large, then the approximation will be the same.

b) Using Laplace's Approximation, we get

$$\mathcal{N}\left(\frac{a-1}{b}, \frac{a-1}{b^2}\right) \approx \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}$$

$$\Gamma(a) \approx \frac{b^a x^{a-1} e^{-bx}}{\mathcal{N}\left(\frac{a-1}{b}, \frac{a-1}{b^2}\right)}$$

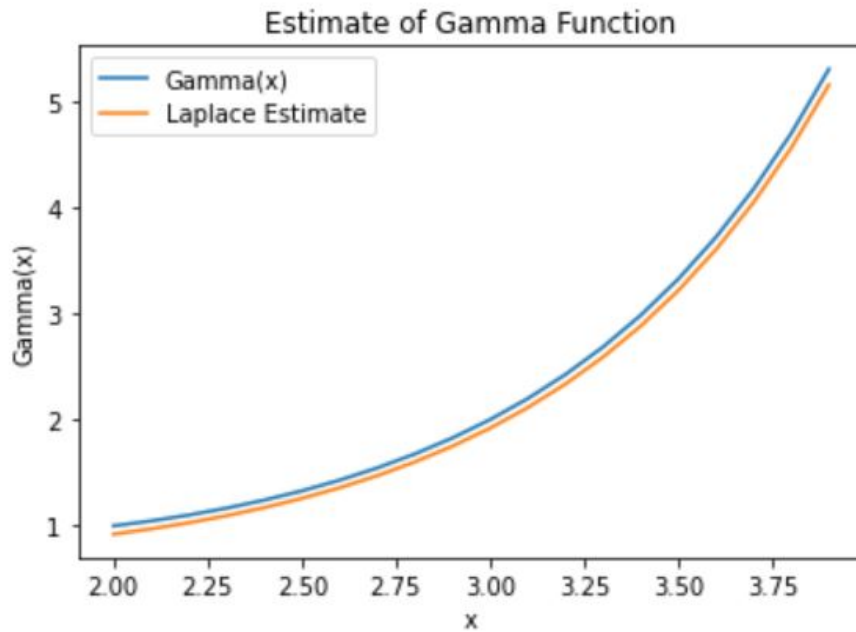
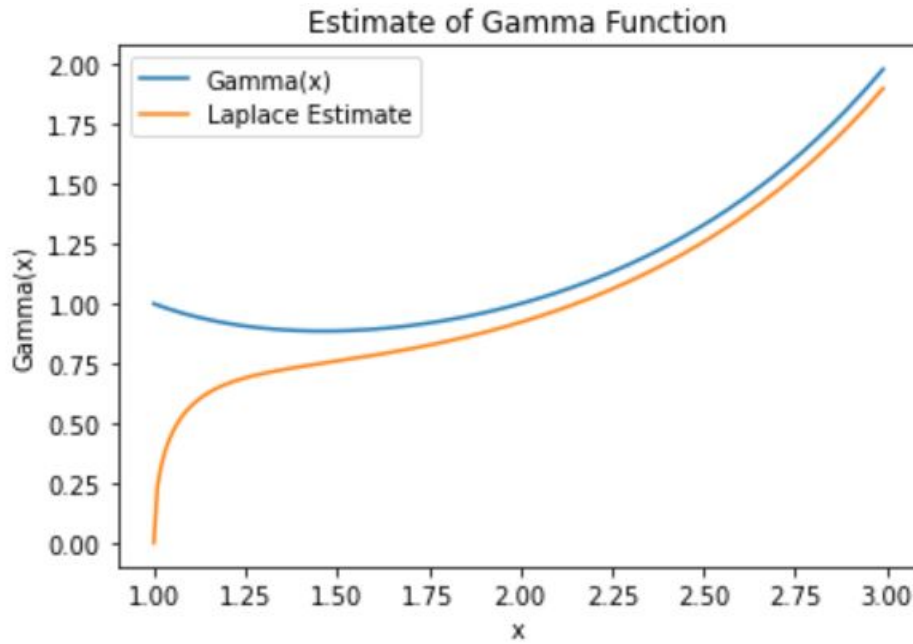
Solving this gets us,

$$\Gamma(a) = \sqrt{2\pi(a-1)} b^{a-1} x^{a-1} e^{-\frac{(bx-(a-1))^2}{2(a-1)} - bx}$$

I have put the value of x_{map} from above since the gaussian and gamma are maximum at x_{map} and we are approximating both, hence we equate their max values.

Putting the value of $x = \frac{a-1}{b}$, we get

$$\Gamma(a) = \sqrt{2\pi(a-1)}\left(\frac{a-1}{e}\right)^{a-1}$$



This is the plot of $\Gamma(a)$ using a standard library function (blue) and the approximated form (orange).

The first plot has values of x bounded by $[1,3]$ and a step of 0.01 and the second plot has values of x bounded by $[2,4]$ and a step of 0.1.

As we can see here the Laplace approximation is very close to the actual plot made by the some standard library function, indicating that laplace approximation is **highly accurate** in nature.

Probabilistic Machine Learning (CS772A), Spring 2023
Indian Institute of Technology Kanpur
Homework Assignment Number 1

QUESTION

2

Student Name: Dishay Mehta

Roll Number: 200341

Date: March 30, 2023

We have N scalar iid observations, $x_1, x_2, x_3, \dots, x_N$

$$\mu \sim \mathcal{N}(\mu_0, s_0)$$

$$\beta \sim \text{Gamma}(\beta|a, b)$$

Finding the conditional posterior, (Assuming $p(\mathbf{X})$ to be a constant $\frac{1}{A}$)

$$\begin{aligned} \bullet p(\mu|\mathbf{X}, \beta) &= \frac{p(\mathbf{X}|\mu, \beta)p(\mu)}{p(\mathbf{X})} \\ &= A \prod_{n=1}^N e^{-\frac{1}{2}\beta(x_n - \mu)^2} e^{-\frac{(\mu - \mu_0)^2}{2s_0}} \\ &= A e^{-\frac{(\mu - \mu_N)^2}{2\sigma_N^2}} \end{aligned}$$

Using the results from slides,

▪ The posterior distribution for the unknown mean parameter μ

On conditioning side, skipping all fixed params and hyperparams from the notation

$$p(\mu|\mathbf{X}) = \frac{p(\mathbf{X}|\mu)p(\mu)}{p(\mathbf{X})} \propto \prod_{n=1}^N \exp\left[-\frac{(x_n - \mu)^2}{2\sigma^2}\right] \times \exp\left[-\frac{(\mu - \mu_0)^2}{2s_0^2}\right]$$

▪ Easy to see that the above will be prop. to exp of a quadratic function of μ . Simplifying:

$$p(\mu|\mathbf{X}) \propto \exp\left[-\frac{(\mu - \mu_N)^2}{2\sigma_N^2}\right]$$

Gaussian posterior's precision is the sum of the prior's precision and sum of the noise precisions of all the observations

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

Gaussian posterior's mean is a convex combination of prior's mean and the MLE solution

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\bar{x} \quad (\text{where } \bar{x} = \frac{\sum_{n=1}^N x_n}{N})$$

Gaussian posterior (not a surprise since the chosen prior was conjugate to the likelihood)

Also the MLE solution for μ

Contribution from the prior

Contribution from the data

here the expressions for σ_N^2 and μ_N are

$$\frac{1}{\sigma_N^2} = N\beta + \frac{1}{s_0}$$

$$\mu_N = \frac{1}{N\beta s_0 + 1}\mu_0 + \frac{N\beta s_0}{N\beta s_0 + 1}\bar{x} \quad \text{where } (\bar{x} = \frac{\sum_{n=1}^N x_n}{N})$$

Thus,

$$p(\mu|\mathbf{X}, \beta) = \mathcal{N}(\mu_N, \sigma_N^2)$$

$$\begin{aligned} \bullet p(\beta|\mathbf{X}, \mu) &= \frac{p(\mathbf{X}|\mu, \beta)p(\beta)}{p(\mathbf{X})} \\ &= A \prod_{n=1}^N e^{-\frac{1}{2}\beta(x_n - \mu)^2} \frac{b^a}{\Gamma(a)} \beta^{a-1} e^{-b\beta} \end{aligned}$$

Using the results from slides,

- Often, it is easier to work with the precision ($=1/\text{variance}$) rather than variance

$$p(x_n|\mu, \lambda^{-1}) = \mathcal{N}(x|\mu, \lambda^{-1}) = \sqrt{\frac{\lambda}{2\pi}} \exp\left[-\frac{\lambda}{2}(x_n - \mu)^2\right]$$

- If mean is known, for precision, $\text{Gamma}(\alpha, \beta)$ is a conjugate prior to Gaussian lik.

Gamma prior on the precision

$$p(\lambda) \propto (\lambda)^{(\alpha-1)} \exp[-\beta\lambda] \quad (\text{Note: mean of } \text{Gamma}(\alpha, \beta) = \frac{\alpha}{\beta})$$

α and β are the shape and rate params, resp. of the Gamma distribution

- (Verify) The posterior $p(\lambda|\mathbf{X})$ will be $\text{Gamma}(\alpha + \frac{N}{2}, \beta + \frac{\sum_{n=1}^N (x_n - \mu)^2}{2})$

Thus,

$$p(\beta|\mathbf{X}, \mu) = \text{Gamma}(\beta|a + \frac{N}{2}, b + \frac{\sum_{n=1}^N (x_n - \mu)^2}{2})$$

Gibbs sampling can be used to get a sampling-based approximation of a multiparameter posterior. Gibbs sampler iteratively draws random samples from CPs. When run long enough, the sampler produces samples from the joint posterior.

This is the case of two parameter (μ, β) , the Gibbs Sampler looks like this

1) Initialize $\beta^{(0)}$

2) For $s=1,2,3,\dots,S$

- Draw a random sample for μ as $\mu^{(s)} \sim p(\mu|\mathbf{X}, \beta^{(s-1)})$
- Draw a random sample for μ as $\beta^{(s)} \sim p(\beta|\mathbf{X}, \mu^{(s)})$

These S samples $(\mu^s, \beta^s)_{s=1}^S$ represent the joint posterior $p(\mu, \beta|\mathbf{X})$

Student Name: Dishay Mehta

Roll Number: 200341

Date: March 30, 2023

We have the following distributions,

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

$$\mathbf{w}|\lambda \sim \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$$

$$\mathbf{w}|\mathbf{y}, \mathbf{X}, \lambda, \beta \sim \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \Sigma)$$

where $\boldsymbol{\mu}$ and Σ are defined as follows

$$\Sigma = (\beta\mathbf{X}^\top\mathbf{X} + \lambda\mathbf{I}_D)^{-1} \quad \boldsymbol{\mu} = (\mathbf{X}^\top\mathbf{X} + \frac{\lambda}{\beta}\mathbf{I}_D)^{-1}\mathbf{X}^\top\mathbf{y}$$

The model parameter is λ and β so it will be a global variable that can be calculated by maximizing $\mathbb{E}[CLL]$ in the M step. \mathbf{w} is a latent variable whose conditional parameter will be computed in the E step. Therefore, we will use EM Algorithm to alternatively model the parameters and the hyperparameters. With EM, can treat \mathbf{w} as "latent var", and λ, β as "parameters".

EM algorithm

This is the whole algorithm with the working also shown, the algorithmic sketch is enclosed inside the box

Step 1:

Initialize λ as λ_0 and β as β_0 and set $t=1$

Step 2: (Expectation step)

- Compute the posterior of \mathbf{w}_t given the current parameters, λ_{t-1} and β_{t-1}

$$\begin{aligned} p(\mathbf{w}_t|\mathbf{y}, \mathbf{X}, \lambda_{t-1}, \beta_{t-1}) &= \mathcal{N}(\boldsymbol{\mu}_{t-1}, \Sigma_{t-1}) \\ \Sigma_{t-1} &= (\beta_{t-1}\mathbf{X}^\top\mathbf{X} + \lambda_{t-1}\mathbf{I}_D)^{-1} \\ \boldsymbol{\mu}_{t-1} &= (\mathbf{X}^\top\mathbf{X} + \frac{\lambda_{t-1}}{\beta_{t-1}}\mathbf{I}_D)^{-1}\mathbf{X}^\top\mathbf{y} \end{aligned}$$

- Now we will calculate the complete log likelihood - $\log(p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \beta, \lambda))$

Using chain rule of probability,

$$p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \beta, \lambda) = p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w}|\lambda)$$

$$\therefore \log(p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \beta, \lambda)) = \log(p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)) + \log(p(\mathbf{w}|\lambda))$$

$$\implies CLL = \frac{1}{2}[N\log\beta + D\log\lambda - \beta(\mathbf{y} - \mathbf{X}\mathbf{w})^\top(\mathbf{y} - \mathbf{X}\mathbf{w}) - \lambda\mathbf{w}^\top\mathbf{w} - (N + D)\log 2\pi]$$

Since the posterior of \mathbf{w} is a normal distribution, we can directly write $\mathbb{E}[\mathbf{w}] = \boldsymbol{\mu}$

- Now we will compute the expectations as follows

$$\mathbb{E}[\mathbf{w}^\top] = \mathbb{E}[\mathbf{w}]^\top = \boldsymbol{\mu}^\top$$

$$\mathbb{E}[\mathbf{w}\mathbf{w}^\top] = Cov(\mathbf{w}) + \mathbb{E}[\mathbf{w}]\mathbb{E}[\mathbf{w}]^\top = \Sigma + \boldsymbol{\mu}\boldsymbol{\mu}^\top$$

$$\mathbb{E}[\mathbf{w}^\top R \mathbf{w}] = \text{Tr}(R \mathbb{E}[\mathbf{w}\mathbf{w}^\top]) = \text{Tr}(R(\Sigma + \boldsymbol{\mu}\boldsymbol{\mu}^\top))$$

- Now we will calculate the expectation of CLL

Using the results from above,

$$\mathbb{E}[CLL] = \frac{1}{2}[N \log \beta_{t-1} + D \log \lambda_{t-1} - \beta_{t-1}(\mathbf{y}^\top \mathbf{y} - \mathbb{E}[\mathbf{w}^\top] \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \mathbb{E}[\mathbf{w}] + \mathbb{E}[\mathbf{w}^\top (\mathbf{X}^\top \mathbf{X} + \frac{\lambda_{t-1}}{\beta_{t-1}} \mathbf{I}_D) \mathbf{w}]) - \text{const}]$$

$$\mathbb{E}[CLL] = \frac{1}{2}[N \log \beta_{t-1} + D \log \lambda_{t-1} - \beta_{t-1}(\mathbf{y}^\top \mathbf{y} - \boldsymbol{\mu}_{t-1}^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \boldsymbol{\mu}_{t-1} + \text{Tr}(\mathbf{X}^\top \mathbf{X}(\Sigma_{t-1} + \boldsymbol{\mu}_{t-1} \boldsymbol{\mu}_{t-1}^\top))) - \lambda_{t-1} \text{Tr}(\Sigma_{t-1} + \boldsymbol{\mu}_{t-1} \boldsymbol{\mu}_{t-1}^\top) - \text{const}]$$

Step 3: Maximization step

Now, we can do MLE for the parameters by maximising the $\mathbb{E}[CLL]$ as in the standard EM algorithm.

$$(\lambda_t, \beta_t) = \underset{\lambda, \beta}{\text{argmax}} \mathbb{E}[CLL]$$

$$= \underset{\lambda, \beta}{\text{argmax}} \frac{1}{2}[N \log \beta_{t-1} + D \log \lambda_{t-1} - \beta_{t-1}(\mathbf{y}^\top \mathbf{y} - \boldsymbol{\mu}_{t-1}^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \boldsymbol{\mu}_{t-1} + \text{Tr}(\mathbf{X}^\top \mathbf{X}(\Sigma_{t-1} + \boldsymbol{\mu}_{t-1} \boldsymbol{\mu}_{t-1}^\top))) - \lambda_{t-1} \text{Tr}(\Sigma_{t-1} + \boldsymbol{\mu}_{t-1} \boldsymbol{\mu}_{t-1}^\top) - \text{const}]$$

We need the point estimate of β_t and λ_t

MLE Estimate of β

We simply require to maximize the Expectation of CLL wrt β

$$\beta_t^{-1} = \frac{1}{N}(\mathbf{y}^\top \mathbf{y} - \boldsymbol{\mu}_{t-1}^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \boldsymbol{\mu}_{t-1} + \text{Tr}(\mathbf{X}^\top \mathbf{X}(\Sigma_{t-1} + \boldsymbol{\mu}_{t-1} \boldsymbol{\mu}_{t-1}^\top)))$$

MLE Estimate of λ

We simply require to maximize the Expectation of CLL wrt λ

$$\lambda_t^{-1} = \frac{1}{D}(\text{Tr}(\Sigma_{t-1} + \boldsymbol{\mu}_{t-1} \boldsymbol{\mu}_{t-1}^\top))$$

where D is the dimension of the vector space

Step 4:

If λ and β are not yet converged then set $t=t+1$ and repeat from step 2.

Student Name: Dishay Mehta
 Roll Number: 200341
 Date: March 30, 2023

We have the following distributions,

$$y_n|z_n \sim \mathbb{I}(z_n > 0)$$

$$z_n|\mathbf{w}, \mathbf{x}_n \sim \mathcal{N}(z_n|\mathbf{w}^\top \mathbf{x}_n, 1) \implies \mathbf{z}|\mathbf{w}, \mathbf{X} \sim \mathcal{N}(\mathbf{z}|\mathbf{X}\mathbf{w}, \mathbf{I}_N)$$

The model parameter is \mathbf{w} so it will be a global variable that can be calculated by maximizing the $\mathbb{E}[C_{LL}]$ in the M step. z_n is a latent variable whose conditional posterior will be computed in the E step. Therefore, we will use EM Algorithm to alternatively model the parameters and the hyperparameters. With EM, can treat z_n as "latent var", and \mathbf{w} as "parameter".

EM algorithm

This is the whole algorithm with the working also shown, the algorithmic sketch is enclosed inside the box

Step 1:

Initialize \mathbf{w} as \mathbf{w}^0 and set $t=1$

Step 2: (Expectation step)

- Compute the posterior of z_n^t given the current parameters, \mathbf{w}^{t-1}

$$p(z_n|x_n, y_n, \mathbf{w}) \propto p(z_n|x_n, \mathbf{w})p(y_n|z_n) \propto \text{prior} * \text{likelihood}$$

$$p(z_n^t|\mathbf{w}^{t-1}, y_n, \mathbf{x}_n) = \begin{cases} \mathcal{N}(z_n^t|\mathbf{w}^{(t-1)\top} \mathbf{x}_n, 1) \mathbb{I}(z_n^t > 0) & \text{if } y_n = 1 \\ \mathcal{N}(z_n^t|\mathbf{w}^{(t-1)\top} \mathbf{x}_n, 1) \mathbb{I}(z_n^t < 0) & \text{if } y_n = 0 \end{cases}$$

Since we constrained the variable of gaussian to be within certain limits, it is called as **Truncated Gaussian**.

- Now we will calculate the complete log likelihood

$$\begin{aligned} \log(p(\mathbf{z}^t, \mathbf{y}|\mathbf{X}, \mathbf{w}^{t-1})) &= \log(p(\mathbf{y}|\mathbf{z}^t)) + \log(p(\mathbf{z}^t|\mathbf{X}, \mathbf{w}^{t-1})) \\ &= \sum_{n=1}^N \log(p(y_n|z_n^t)) - \frac{1}{2}(\mathbf{z}^t - \mathbf{X}\mathbf{w}^{t-1})^\top (\mathbf{z}^t - \mathbf{X}\mathbf{w}^{t-1}) + \text{const} \end{aligned}$$

- Now we will calculate the expectation of CLL

Since CLL depends linearly on \mathbf{z}^t , we just used to calculate $\mathbb{E}[z_n^t|\mathbf{w}^{t-1}, y_n, \mathbf{x}_n]$ to calculate $\mathbb{E}[C_{LL}]$. Since $p(z_n^t|\mathbf{w}^{t-1}, y_n, \mathbf{x}_n)$ is a truncated gaussian, so we just need to compute the posterior mean as \mathbf{w}^{t-1} and z_n^t share a linear relation and hence we need to compute the expectation of $z_n^t|\mathbf{x}_n, y_n, \mathbf{w}^{t-1}$

$$\mathbb{E}[z_n^t|\mathbf{w}^{t-1}, y_n, \mathbf{x}_n] = \begin{cases} \mathbf{w}^{(t-1)\top} \mathbf{x}_n + \frac{\phi(\mathbf{w}^{(t-1)\top} \mathbf{x}_n)}{1 - \Phi(-\mathbf{w}^{(t-1)\top} \mathbf{x}_n)} & \text{if } y_n = 1 \\ \mathbf{w}^{(t-1)\top} \mathbf{x}_n + \frac{\phi(\mathbf{w}^{(t-1)\top} \mathbf{x}_n)}{\Phi(-\mathbf{w}^{(t-1)\top} \mathbf{x}_n)} & \text{if } y_n = 0 \end{cases}$$

where $\phi(\cdot)$ is the standard normal probability distribution function and $\Phi(\cdot)$ is the Probit function, which is given on the wikipedia link given in the assignment.

- Now we calculate the expected complete log likelihood

$$\mathbb{E}[CLL] = \text{const} - \mathbb{E}\left[\frac{(\mathbf{z}^t - \mathbf{X}\mathbf{w}^{t-1})^\top (\mathbf{z}^t - \mathbf{X}\mathbf{w}^{t-1})}{2}\right] \quad (\text{the const term is indep of } \mathbf{w}^{t-1})$$

- Expected CLL in EM is given by (assume observations are i.i.d.)

$$\begin{aligned} \mathcal{Q}(\Theta, \Theta^{\text{old}}) &= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \Theta^{\text{old}})} [\log p(\mathbf{x}_n, \mathbf{z}_n | \Theta)] \\ &= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \Theta^{\text{old}})} [\log p(\mathbf{x}_n | \mathbf{z}_n, \Theta) + \log p(\mathbf{z}_n | \Theta)] \end{aligned}$$

- If $p(\mathbf{z}_n | \Theta)$ and $p(\mathbf{x}_n | \mathbf{z}_n, \Theta)$ are exp-family distributions, $\mathcal{Q}(\Theta, \Theta^{\text{old}})$ has a very simple form
- In resulting expressions, replace terms containing \mathbf{z}_n 's by their respective expectations, e.g.,
 - \mathbf{z}_n replaced by $\mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \Theta)}[\mathbf{z}_n]$
 - $\mathbf{z}_n \mathbf{z}_n^\top$ replaced by $\mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \Theta)}[\mathbf{z}_n \mathbf{z}_n^\top]$

Thus from the above results,

$$\mathbb{E}[CLL] = \text{const} - \frac{\|\mathbb{E}[\mathbf{z}^t] - \mathbf{X}\mathbf{w}^{t-1}\|^2}{2}$$

Step 3: Maximization Step

Now, we can do MLE for the parameters by maximising the $\mathbb{E}[CLL]$ as in the standard EM algorithm.

$$\begin{aligned} \mathbf{w}^t &= \underset{\mathbf{w}}{\text{argmax}} \mathbb{E}[CLL] \\ &= \underset{\mathbf{w}}{\text{argmax}} \left(\text{const} - \frac{\|\mathbb{E}[\mathbf{z}^t] - \mathbf{X}\mathbf{w}^{t-1}\|^2}{2} \right) \end{aligned}$$

MLE Estimate of \mathbf{w}

We require to maximize the Expectation of CLL wrt \mathbf{w}

$$\mathbf{w}^t = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}[\mathbf{z}^t]$$

$$\text{where } \mathbb{E}[\mathbf{z}^t] = \begin{bmatrix} \mathbb{E}[z_1^t], \mathbb{E}[z_2^t], \dots, \mathbb{E}[z_N^t] \end{bmatrix}^\top$$

Step 4:

If \mathbf{w}^t are not yet converged then set $t=t+1$ and repeat from step 2.

Student Name: Dishay Mehta

Roll Number: 200341

Date: March 30, 2023

Part 1

Given:

$$p(\mathbf{f}) = \text{GP}(0, \kappa) = \mathcal{N}(0, \kappa)$$

$$p(y_n | x_n, f) = \mathcal{N}(y_n | f(x_n), \sigma^2)$$

For posterior,

$$p(\mathbf{f} | \mathbf{y}) = \frac{p(\mathbf{f})p(\mathbf{y} | \mathbf{f})}{p(\mathbf{y})}$$

here $p(\mathbf{y})$ is independent of \mathbf{f} , so the posterior becomes,

$$p(\mathbf{f} | \mathbf{y}) \propto p(\mathbf{f})p(\mathbf{y} | \mathbf{f})$$

$$\propto \mathcal{N}(0, \kappa) \prod_{n=1}^N \mathcal{N}(y_n | f(x_n), \sigma^2)$$

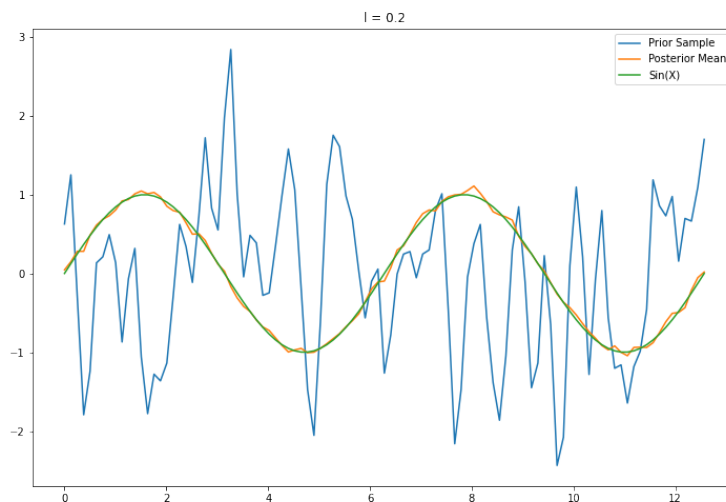
$$\propto \mathcal{N}(0, \kappa) \mathcal{N}(\mathbf{y} | \mathbf{f}, \sigma^2 \mathbf{I}_N)$$

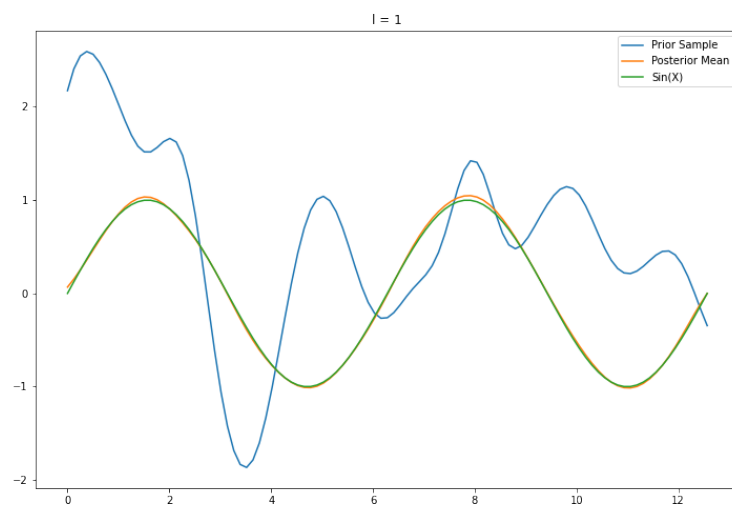
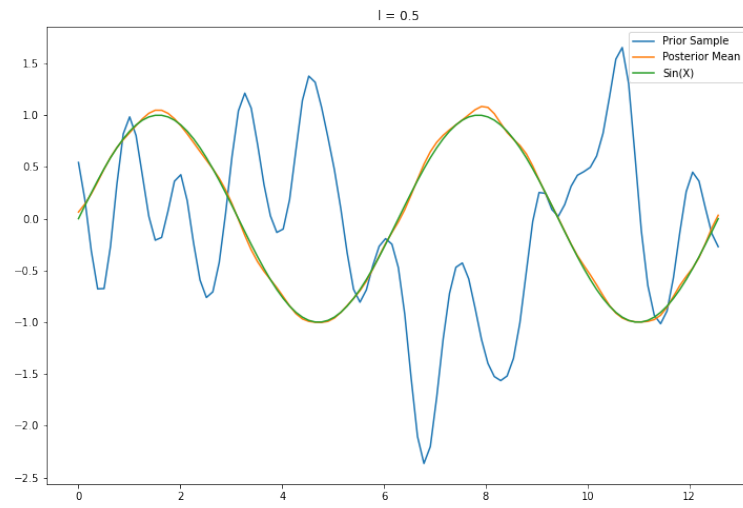
$$p(\mathbf{f} | \mathbf{y}) \propto \exp\left[-\frac{1}{2\sigma^2}(\mathbf{f} - \mathbf{y})^\top(\mathbf{f} - \mathbf{y}) - \frac{1}{2}\mathbf{f}^\top \kappa^{-1} \mathbf{f}\right] \therefore p(\mathbf{f} | \mathbf{y}) = \mathcal{N}(\mathbf{f} | \mu, \Sigma)$$

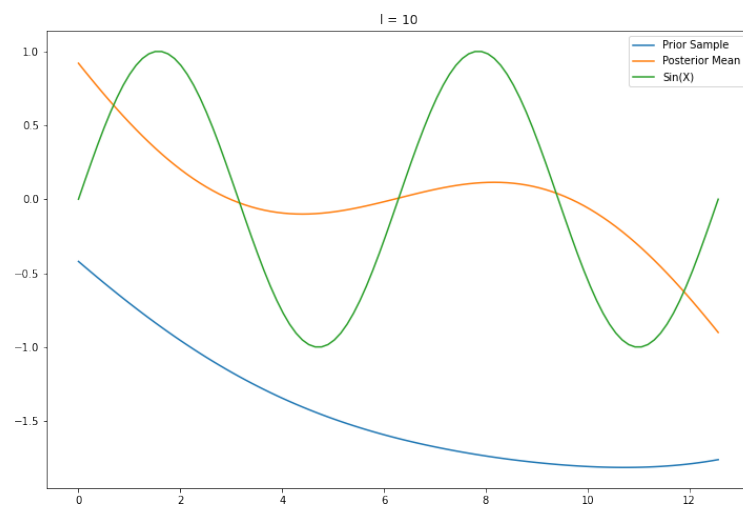
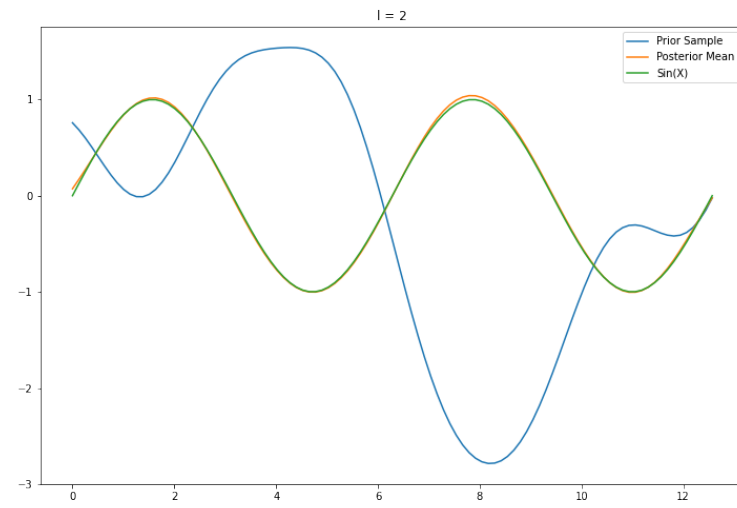
where $\Sigma = \sigma^2(\sigma^2 \mathbf{I}_N + \kappa)^{-1} \kappa$ and $\mu = \kappa(\sigma^2 \mathbf{I}_N + \kappa)^{-1} \mathbf{y}$

Part 2

Visualising the GP Priors and Posteriors for Regression







Inference:

The difference between the plots generated using different values of the parameter l is mainly in the smoothness and amplitude of the GP prior and posterior functions, as well as in the uncertainty estimates. Here are some specific observations:

- For smaller values of l , the GP functions tend to be wiggly and have higher frequency variations and for larger values of l , the GP functions tend to be smoother and have lower frequency variations.
- The GP posterior functions tend to be smoother than the GP prior functions because they are conditioned on the observed data and thus have reduced uncertainty.
- As l increases the number of peaks and valleys in prior decreases. The posterior however is very close to the true function as we increase l from 0.2 to 2 but then for $l=10$, there is a significant difference between the posterior mean and true function, since the kernel becomes large as l increases.