

Student Name: Musale Krushna Pavan  
 Roll Number: 20111268  
 Date: April 20, 2021

### Batch Bayesian Active Learning

#### Overview

Introduced the new Bayesian batch active learning approach that overcomes standard greedy procedures problems. In this algorithm produces diverse batches that enable efficient active learning at scale. We derive interpretable closed-form solutions akin to existing active learning procedures for linear models, and generalize to arbitrary models using random projections.

#### The key idea

Select the next **batch of data that best approximates the Expected complete data log posterior**.

Expected complete data log posterior is defined as follows:

$$\begin{aligned} \mathbb{E}_{\mathcal{Y}_p} [\log p(\boldsymbol{\theta} | \mathcal{D}_0 \cup (\mathcal{X}_p, \mathcal{Y}_p))] \\ = \log p(\boldsymbol{\theta} | \mathcal{D}_0) + \underbrace{\sum_{m=1}^M \left( \mathbb{E}_{\mathbf{y}_p} [\log p(\mathbf{y}_p | \mathbf{x}_p, \boldsymbol{\theta})] + \mathbb{H}(\mathbf{y}_p | \mathbf{x}_p, \mathcal{D}_0) \right)}_{\mathcal{L}_m(\boldsymbol{\theta})} \end{aligned} \quad (1)$$

This means that expectation of the parameters given all the data. Now the idea is to find the batch of data  $\mathcal{D}'$  in a way that updated log posterior  $P(\boldsymbol{\theta} | \mathcal{D}_0 \cup \mathcal{D}')$  best approximates the complete data log posterior  $P(\boldsymbol{\theta} | \mathcal{D}_0 \cup \mathcal{D}_p)$ .

As the first term depends only on  $\mathcal{D}_0$ . We can choose the batch that best approximates  $\mathcal{L} = \sum_m \mathcal{L}_m$ . where  $\mathcal{L}_m : \boldsymbol{\Theta} \rightarrow \mathcal{R}$  considered as vector in function space. We can achieve this above idea by the following sparse objective.

$$\boldsymbol{\omega}^* = \arg \min_{\boldsymbol{\omega}} \|\mathcal{L} - \mathcal{L}((\boldsymbol{\omega}))\|^2 \text{ such that } \omega_m \in \{0, 1\} \quad \forall m, \sum_m \omega_m \leq b \quad (2)$$

#### Relaxed objective

As solving Eq. (2) is generally intractable, To approximately solve the optimization problem Eq.(2) by relaxing the **binary weight constraint to be non-negative** and replace the **cardinality constraint with a polytope** constraint in Hilbert space.

Let  $\sigma_m = \|\mathcal{L}_m\|$ ,  $\sigma = \sum_m \omega_m \sigma_m$  and  $\mathbf{K} \in \mathbb{R}^{M \times M}$  be a kernel matrix with  $K_{mn} = \langle \mathcal{L}_m, \mathcal{L}_n \rangle$ . Then the relaxed objective can be written as follows:

$$\arg \min_{\boldsymbol{\omega}} (1 - \boldsymbol{\omega})^T \mathbf{K} (1 - \boldsymbol{\omega}) \text{ subject to } \omega_m \geq 0 \quad \forall m, \sum_m \omega_m \sigma_m = \sigma \quad (3)$$

where we used  $\|\mathcal{L} - \mathcal{L}((\boldsymbol{\omega}))\|^2 = (1 - \boldsymbol{\omega})^T \mathbf{K} (1 - \boldsymbol{\omega})$ . The polytope has vertices  $\{\sigma / \sigma_m \mathbf{1}_m\}_{m=1}^M$  and contains the point  $\boldsymbol{\omega} = [1, 1, \dots, 1]^T$ . Now we can solve equation 3 using Frank-Wolfe algorithm which intern solves our objective Eq2.

The optimum of Eq(3) is attained at vertices of polytope. Frank-Wolfe algorithm runs for  $b$

number of iterations to get  $\omega^*$ , the weights are iteratively updated along the f-th vertex of polytope. Frank-Wolfe algorithm requires the choice of inner product1  $\langle \mathcal{L}_m \mathcal{L}_n \rangle$  and the norm  $\sigma_n = \|\mathcal{L}_n\|$   
choice1: Fisher inner product

$$\langle \mathcal{L}_m \mathcal{L}_n \rangle_{\hat{\pi}, \mathcal{F}} = \mathbb{E}_{\hat{\pi}}[\nabla_{\theta} \mathcal{L}_n(\theta)^T \nabla_{\theta} \mathcal{L}_m(\theta)] \quad (4)$$

Advantage of this choice of inner product is that for specific models this leads to **simple, interpretable** expressions.

Choice2: d Euclidean inner product

$$\langle \mathcal{L}_m \mathcal{L}_n \rangle_{\hat{\pi}, 2} = \mathbb{E}_{\hat{\pi}}[\mathcal{L}_n(\theta) \mathcal{L}_m(\theta)] \quad (5)$$

Advantage of this choice of inner product is that only **requires tractable likelihood** computations. For this we can use black-box method

Once we get the continuous weight from the Frank-Wolfe algorithm we project them we project them back to the feasible space using

$$\text{if } \omega_m^* > 0 : \quad (6)$$

$$\hat{\omega}_m^* = 1 \quad (7)$$

$$\text{else} \quad (8)$$

$$\hat{\omega}_m^* = 0 \quad (9)$$

#### Models having closed form expression:

By using the weighted Fisher inner product from we get the closed form expression for the two of the following models:

1. Bayesian linear regression
2. Probit regression

For the models for which we can't compute the gradient we used **random feature projection** which approximates the key quantities required. This method which resolves gradient issue and also linearly complex with the pool set  $|\mathcal{P}|$ . For this approach we use projections for the weighted Euclidean inner product.

The appropriate projection which can be calculated for any model with tractable likelihood

$$\hat{\mathcal{L}}_n = \frac{1}{\sqrt{J}}[\mathcal{L}_n(\theta_1), \dots, \mathcal{L}_n(\theta_J)]^T, \quad \theta_j \sim \hat{\pi} \quad (10)$$

$\hat{\mathcal{L}}_n$  represents the J-dimensional projection of  $\mathcal{L}_n$  in Euclidean space. Using above approximate we can approximate inner products as dot products between vectors

$$\langle \mathcal{L}_m \mathcal{L}_n \rangle_{\hat{\pi}, 2} = \hat{\mathcal{L}}_n^T \hat{\mathcal{L}}_m \quad (11)$$

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: April 20, 2021

Local Conjugacy and Gibbs Sampling  
 Given

$$x_1, x_2 \dots x_N \sim \mathcal{N}(x|\mu, \beta^{-1}) \quad (12)$$

$$\mu \sim \mathcal{N}(\mu|\mu_0, s_0) \quad (13)$$

$$\beta \sim \text{Gamma}(\beta|a, b) \quad (14)$$

Now finding the conditional posterior of  $\mu$  i.e  $P(\mu|\mathbf{X}, \beta, \mu_0, s_0)$

$$P(\mu|\mathbf{X}, \beta, \mu_0, s_0) \propto p(\mathbf{X}|\mu, \beta)p(\mu|\mu_0, s_0) \quad (15)$$

$$\propto \mathcal{N}(\mathbf{X}|\mu, \beta^{-1})\mathcal{N}(\mu|\mu_0, s_0) \quad (16)$$

$$\text{by conjugacy} \quad (17)$$

$$P(\mu|\mathbf{X}, \beta, \mu_0, s_0) = \mathcal{N}(\mu_N, s_N^{-1}) \quad (18)$$

$$\text{where} \quad (19)$$

$$\mu_N = \frac{\mu_0 + s_0\beta \sum_{n=1}^N x_n}{\beta N s_0 + 1} \quad (20)$$

$$s_N^{-1} = \frac{1}{s_0} + N\beta \quad (21)$$

Conditional posterior of  $\beta$  i.e  $P(\beta|\mathbf{X}, \mu, a, b)$

$$P(\beta|\mathbf{X}, \mu, a, b) \propto p(\mathbf{X}|\mu, \beta)P(\beta|a, b) \quad (22)$$

$$\propto \mathcal{N}(\mathbf{X}|\mu, \beta^{-1})\text{Gamma}(\beta|a, b) \quad (23)$$

$$\text{by conjugacy} \quad (24)$$

$$P(\beta|x, \mu, a, b) = \text{Gamma}(a + \frac{N}{2}, b + \frac{\sum_{n=1}^N (x_n - \mu)^2}{2}) \quad (25)$$

Now using Gibbs Sampling algorithm for finding the combined posterior  
 Steps

**Result:** Set of samples  $(\mu^{(s)}, \beta^{(s)})_{s=1}^S$

Intialize an estimate of  $\mu^0, \beta^0$  ;

**for**  $s=1$  to  $T$  **do**

    Draw a random sample  $\beta^{(s)} \sim \text{Gamma}(a + \frac{N}{2}, b + \frac{\sum_{n=1}^N (x_n - \mu^{(s-1)})^2}{2})$

    Draw a random sample  $\mu^{(s)} \sim \mathcal{N}(\mu_N(s), s_N^{-1}(s))$

**end**

**Algorithm 1:** Gibbs Sampling

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

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: April 20, 2021

**EM for Sparse Modeling** Given:

A linear regression model  $\mathbf{y} = \mathbf{X}\boldsymbol{\omega} + \boldsymbol{\epsilon}$

$$\omega_d \sim \mathcal{N}(0, \sigma^2 k_{\gamma_d}) \quad \text{where } k_{\gamma_d} = \gamma_d v_1 + (1 - \gamma_d) v_0 \text{ and } v_1 > v_0 > 0 \quad (26)$$

$$\boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{K}) \quad \text{where } \mathbf{K} = \text{diag}(k_{\gamma_1}, k_{\gamma_2}, \dots, k_{\gamma_d}) \quad (27)$$

$$\theta \sim \text{Beta}(a_0, b_0) \quad (28)$$

$$\gamma_d \sim \text{Bernoulli}(\theta) \quad (29)$$

$$\sigma^2 \sim \text{IG}(v/2, v\lambda/2) \quad (30)$$

$$\epsilon_n \sim \mathcal{N}(0, \sigma^2) \quad (31)$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N) \quad (32)$$

### Part 1

The effect of assuming the above prior on  $\boldsymbol{\omega}$  induces sparsity on the weight vector  $\boldsymbol{\omega}$  although the sparsity is present for only some parameters.

We have chosen two types of variances in the prior of the weights. which mean that we are taking different amount of precision on different weight parameters. This type of prior first classifies the weight paramerts into 2 categories based on their importance.

### Part 2

**EM Algorithm:**

#### step 1:

Initialize the following parameters randomly and set  $t = 1$

$$\gamma^{(0)}, \quad \sigma^{2(0)}, \quad \theta^{(0)} \quad (33)$$

#### step2:

computing the posterior over  $\boldsymbol{\omega}$

$$p(\boldsymbol{\omega}^{(t)} | \mathbf{y}, \mathbf{X}, \gamma^{(t-1)}, \sigma^{2(t-1)}, \theta^{(t-1)}) \propto p(\mathbf{y} | \boldsymbol{\omega}, \mathbf{X}, \sigma^{2(t-1)}) p(\boldsymbol{\omega} | \sigma^{2(t-1)}, \gamma^{(t-1)}) \quad (34)$$

$$\propto \mathcal{N}(\mathbf{y} | \mathbf{X}\boldsymbol{\omega}, \sigma^{2(t-1)} \mathbf{I}_N) \mathcal{N}(\boldsymbol{\omega} | \mathbf{0}, \sigma^{2(t-1)} \mathbf{K}^{(t-1)}) \quad (35)$$

$$\text{using the gaussian properties} \quad (36)$$

$$= \mathcal{N}(\boldsymbol{\omega}^{(t)} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (37)$$

$$\boldsymbol{\mu}^{(t)} = (\mathbf{K}^{(t-1)^{-1}} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (38)$$

$$\boldsymbol{\Sigma}^{(t)} = \sigma^{2(t)} (\mathbf{K}^{(t-1)^{-1}} + \mathbf{X}^T \mathbf{X})^{-1} \quad (39)$$

#### step3:

Maximize the expected complete log likelihood

$\mathbb{E}[\mathbf{CLL}]$

$$\mathbb{E}[\log p(\boldsymbol{\omega}^{(t)}, \mathbf{y}|\gamma, \theta, \sigma^2)] = \mathbb{E}[\log p(\mathbf{y}|\boldsymbol{\omega}^{(t)}, \gamma, \sigma^2) + \log p(\boldsymbol{\omega}^{(t)}|\gamma, \sigma^2)] \quad (40)$$

$$= \mathbb{E}[\log \mathcal{N}(\mathbf{y}|\mathbf{X}\boldsymbol{\omega}^{(t)}, \sigma^2 \mathbf{I}_N) + \log \mathcal{N}(\boldsymbol{\omega}^{(t)}|\mathbf{0}, \sigma^2 \mathbf{K})] \quad (41)$$

$$= \mathbb{E}\left[-\frac{N+D}{2} \log(2\pi\sigma^2) - \frac{\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \boldsymbol{\omega}^{(t)} + \boldsymbol{\omega}^{(t)T} \mathbf{X}^T \mathbf{X} \boldsymbol{\omega}^{(t)}}{2\sigma^2} - \frac{1}{2} \sum_d \log k_{\gamma_d} - \frac{\boldsymbol{\omega}^{(t)T} \mathbf{K}^{-1} \boldsymbol{\omega}^{(t)}}{2\sigma^2}\right] \quad (42)$$

$$= -\frac{N+D}{2} \log(2\pi\sigma^2) - \frac{1}{2} \sum_d \log k_{\gamma_d} - \frac{\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbb{E}[\boldsymbol{\omega}^{(t)}] + \mathbb{E}[\boldsymbol{\omega}^{(t)T} (\mathbf{X}^T \mathbf{X} + \mathbf{K}^{-1}) \boldsymbol{\omega}^{(t)}]}{2\sigma^2} \quad (43)$$

where

$$\mathbb{E}[\boldsymbol{\omega}^{(t)}] = \boldsymbol{\mu}^{(t)}$$

$$\mathbb{E}[\boldsymbol{\omega}^{(t)T} (\mathbf{X}^T \mathbf{X} + \mathbf{K}^{-1}) \boldsymbol{\omega}^{(t)}] = \text{Tr}((\mathbf{X}^T \mathbf{X} + \mathbf{K}^{-1}) \boldsymbol{\Sigma}^{(t)}) + \boldsymbol{\mu}^{(t)T} (\mathbf{X}^T \mathbf{X} + \mathbf{K}^{-1}) \boldsymbol{\mu}^{(t)}$$

Ref: Matrixcookbook Eq(328)

#### Step 4:

maximization for MAP

$$\gamma^{(t)}, \sigma^{2(t)}, \theta^{(t)} = \arg \max_{\gamma, \sigma^2, \theta} \mathbb{E}[\mathbf{CLL}] + \log p(\gamma, \sigma^2, \theta) \quad (44)$$

$$= \arg \max_{\gamma, \sigma^2, \theta} \mathbb{E}[\mathbf{CLL}] + \log p(\gamma|\theta) + \log p(\sigma^2) + \log p(\theta) \quad (45)$$

$$= \arg \max_{\gamma, \sigma^2, \theta} \mathbb{E}[\mathbf{CLL}] + \sum_d \log p(\gamma_d|\theta) + \log p(\sigma^2) + \log p(\theta) \quad (46)$$

$$\text{where} \quad (47)$$

$$\log p(\gamma_d|\theta) = \gamma_d \log \theta + (1 - \gamma_d) \log(1 - \theta) \quad (48)$$

$$\log p(\sigma^2) = v/2 \log(v\lambda/2) - \log(\Gamma(v/2)) - (v/2 + 1) \log \sigma^2 - \frac{v\lambda}{2\sigma^2} \quad (49)$$

$$\log p(\theta) = (a_0 - 1) \log \theta + (b_0 - 1) \log(1 - \theta) + \log\left(\frac{\Gamma(a_0)\Gamma(b_0)}{\Gamma(a_0 + b_0)}\right) \quad (50)$$

Derivating wrt  $\sigma^2$  and equating to 0

$$0 = -\frac{N+D}{2\sigma^2} + \frac{\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbb{E}[\boldsymbol{\omega}^{(t)}] + \mathbb{E}[\boldsymbol{\omega}^{(t)T} (\mathbf{X}^T \mathbf{X} + \mathbf{K}^{-1}) \boldsymbol{\omega}^{(t)}]}{2\sigma^4} - \frac{v+2}{2\sigma^2} + \frac{v\lambda}{2\sigma^4} \quad (51)$$

$$\sigma^{2(t)} = \frac{\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbb{E}[\boldsymbol{\omega}^{(t)}] + \mathbb{E}[\boldsymbol{\omega}^{(t)T} (\mathbf{X}^T \mathbf{X} + \mathbf{K}^{-1}) \boldsymbol{\omega}^{(t)}] + v\lambda}{N+D+v+2} \quad (52)$$

As the  $\gamma_d \in 0, 1$  we can directly evaluate it.  $\arg \max_{\gamma_d}$  can be written (having the terms only related to  $\gamma_d$ ) as follows:

$$\gamma_d^{(t)} = \arg \max_{\gamma_d \in \{0,1\}} \left\{ -\frac{1}{2} \log k_{\gamma_d} - \frac{\mathbb{E}[\boldsymbol{\omega}^{(t)T} (\mathbf{X}^T \mathbf{X} + \mathbf{K}^{-1}) \boldsymbol{\omega}^{(t)}]}{2\sigma^{2(t)}} \right. \quad (53)$$

$$\left. + \gamma_d \log \theta^{(t-1)} + (1 - \gamma_d) \log(1 - \theta^{(t-1)}) \right\} \quad (54)$$

We can directly substitute the values of  $\gamma_d = 0, 1$  and then find out the value for which is has the above equation has max value

Derivating wrt  $\theta$  and equating to 0

$$0 = \frac{\sum_d \gamma_d^{(t)}}{\theta} - \frac{D - \sum_d \gamma_d^{(t)}}{1 - \theta} + \frac{a_0 - 1}{\theta} - \frac{b_0 - 1}{1 - \theta} \quad (55)$$

$$\theta^{(t)} = \frac{\sum_d \gamma_d^{(t)} + a_0 - 1}{D + a_0 + b_0 - 2} \quad (56)$$

**Step5 :**

if not converged set  $t = t + 1$  and go to step2

else return  $\gamma^{(t)}, \sigma^{2(t)}, \theta^{(t)}$  and  $p(\boldsymbol{\omega}^{(t)} | \mathbf{y}, \mathbf{X}, \gamma^{(t-1)}, \sigma^{2(t-1)}, \theta^{(t-1)})$

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: April 20, 2021

### Part 1: GP Posterior

Given

$$\text{likelihood : } p(\mathbf{y}|\mathbf{f}) : \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}_N) \quad (57)$$

$$\text{prior : } p(\mathbf{f}|\mathbf{X}) : \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}) \quad (58)$$

$$\text{where } \mathbf{K}_{ij} = k(x_i, x_j) \quad (59)$$

By using the gaussian properties from the lecture 4.1

The posterior

$$p(\mathbf{f}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}) \quad (60)$$

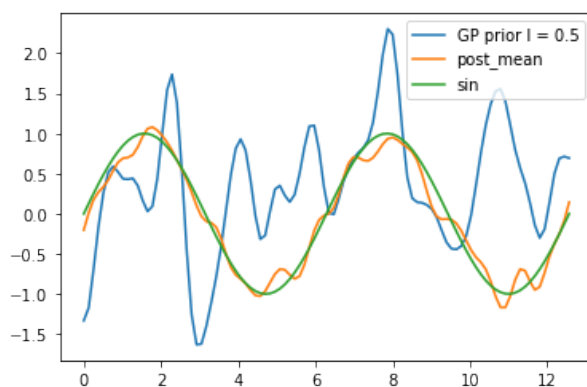
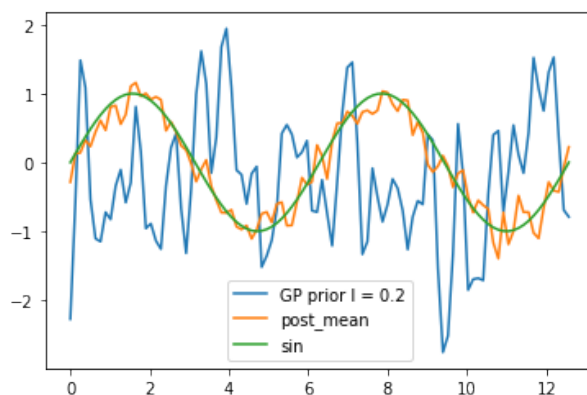
$$= \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N) \quad (61)$$

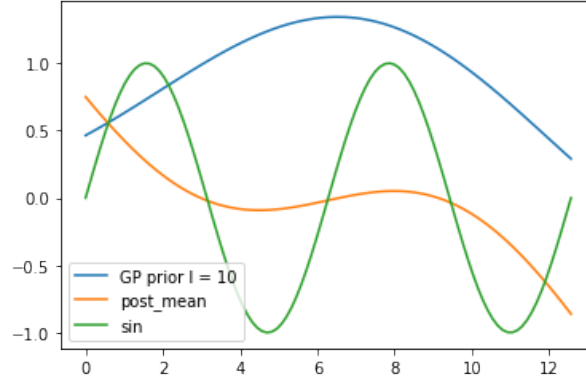
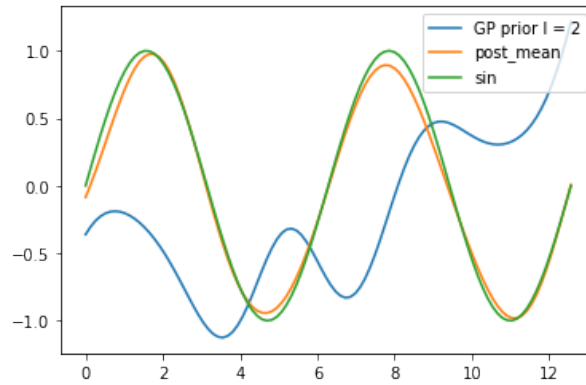
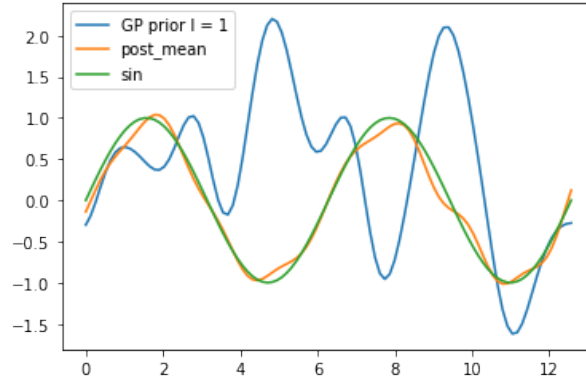
$$\text{where} \quad (62)$$

$$\boldsymbol{\mu}_N = (\sigma^2 \mathbf{K}^{-1} + \mathbf{I}_N)^{-1} \mathbf{y} \quad (63)$$

$$\boldsymbol{\Sigma}_N = (\mathbf{K}^{-1} + \sigma^{-2} \mathbf{I}_N)^{-1} \quad (64)$$

### Part 2: Visualizing GP Priors and Posteriors for Regression





$l$  signifies how close  $x$  and  $x'$  should be to influence each other significantly. Lower  $l$  values create more wiggles in function with wide uncertainty between training data points. As the  $l$  value increases we can see that prior function of the  $GP$  becomes smoother. Where as the posterior mean tries to overfit the data for small values of  $l$  but later it gets closer to the original sin function but after larger values of  $l$  it also becomes oversmooth and underfits the  $\sin$  curve



Student Name: Musale Krushna Pavan  
 Roll Number: 20111268  
 Date: April 20, 2021

**Speeding up Gaussian Processes**  
**Part 1 Given:**

training data: (65)

$$(\mathbf{X}, \mathbf{f}) = \{x_n, f_n\}_{n=1}^N \quad (66)$$

$$y_n = f(\mathbf{x}_n) = f_n \quad (67)$$

$$\mathbf{f} \sim \mathcal{GP}(0, k) \quad (68)$$

posterior predictive distribution (69)

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{f}) = \mathcal{N}(f_* | \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{f}, k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*) \quad (70)$$

Let's reduce this cost to make GPs more scalable. To do this, suppose there are

pseudo training inputs (71)

$$\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_M\} \quad M \ll N \quad (72)$$

pseudo training outputs (73)

$$\mathbf{t} = \{t_1, \dots, t_M\} \quad (74)$$

modeled by (75)

$$t_m = f(\mathbf{z}_m) \quad (76)$$

Note: that  $(\mathbf{Z}, \mathbf{t})$  are NOT known

Let likelihood for each training output  $f_n$  with PPD as follows

$$p(f_* | \mathbf{x}_*, \mathbf{Z}, \mathbf{t}) = \mathcal{N}(f_* | \tilde{\mathbf{k}}_*^T \tilde{\mathbf{K}}^{-1} \mathbf{t}, k(\mathbf{x}_*, \mathbf{x}_*) - \tilde{\mathbf{k}}_*^T \tilde{\mathbf{K}}^{-1} \tilde{\mathbf{k}}_*) \quad (77)$$

The likelihood can be shown as

$$p(\mathbf{f} | \mathbf{X}, \mathbf{Z}, \mathbf{t}) = \prod_{n=1}^N p(f_n | \mathbf{x}_n, \mathbf{Z}, \mathbf{f}) \quad (78)$$

$$= N(\mathbf{f} | \mathbf{K}_{NM} \tilde{\mathbf{K}}^{-1} \mathbf{t}, \Lambda) \quad (79)$$

where  $\Lambda$  is a diagonal matrix with (80)

$$\Lambda_{*,*} = k(\mathbf{x}_*, \mathbf{x}_*) - \tilde{\mathbf{k}}_*^T \tilde{\mathbf{K}}^{-1} \tilde{\mathbf{k}}_* \text{ and } [\mathbf{K}_{NM}]_{nm} = k(x_n, x_m) \quad (81)$$

finding the prior over pseudo data

$$p(\mathbf{t} | \mathbf{Z}) = \mathcal{N}(\mathbf{t} | \mathbf{0}, \tilde{\mathbf{K}}) \quad (82)$$

Here, we assume the pseudo-inputs  $\mathbf{Z}$  are known. finding the prior over pseudo data

$$p(\mathbf{t}|\mathbf{Z}, \mathbf{X}, \mathbf{f}) \propto p(\mathbf{f}|\mathbf{X}, \mathbf{Z}, \mathbf{t})p(\mathbf{t}|\mathbf{Z}) \quad (83)$$

$$\text{using gaussian properties} \quad (84)$$

$$= \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t) \quad (85)$$

$$\text{where} \quad (86)$$

$$\boldsymbol{\Sigma}_t = (\tilde{\mathbf{K}}^{-1} + (\mathbf{K}_{NM}\tilde{\mathbf{K}}^{-1})^T \Lambda^{-1} \mathbf{K}_{NM}\tilde{\mathbf{K}}^{-1})^{-1} \quad (87)$$

$$\boldsymbol{\mu}_t = \boldsymbol{\Sigma}_t \tilde{\mathbf{K}}^{-1} \mathbf{K}_{NM}^T \Lambda^{-1} \mathbf{f} \quad (88)$$

expression of the posterior predictive distribution for the output  $y_*$  of a new input  $x_*$

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{f}, \mathbf{Z}) = \int p(y_*|\mathbf{x}_*, \mathbf{t}, \mathbf{Z})p(\mathbf{t}|\mathbf{Z}, \mathbf{X}, \mathbf{f})d\mathbf{t} \quad (89)$$

$$= \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \quad (90)$$

$$\text{where} \quad (91)$$

$$\boldsymbol{\mu}_* = \tilde{\mathbf{k}}_*^T \tilde{\mathbf{K}}^{-1} \boldsymbol{\mu}_t \quad (92)$$

$$\boldsymbol{\Sigma}_* = \tilde{\mathbf{k}}_*^T \tilde{\mathbf{K}}^{-1} \boldsymbol{\Sigma}_t \tilde{\mathbf{K}}^{-1} \tilde{\mathbf{k}}_* + k(\mathbf{x}_*, \mathbf{x}_*) - \tilde{\mathbf{k}}_*^T \tilde{\mathbf{K}}^{-1} \tilde{\mathbf{k}}_* \quad (93)$$

### Time Complexity:

We can see from the above equation that calculation  $\boldsymbol{\mu}_*$  takes the most time in which  $\mathbf{K}_{NM}^T \Lambda^{-1} \mathbf{K}_{NM}$  takes  $\mathcal{O}(M^2 N)$ . Through the  $\tilde{\mathbf{K}}^{-1}$  and some other calculations takes  $\mathcal{O}(M^3)$  its less than  $\mathcal{O}(M^2 N)$  as  $M \ll N$ .

So the time complexity using this pseudo training data is  $\mathcal{O}(M^2 N)$ .

## Part 2

Marginal Likelihood

$$p(\mathbf{f}|\mathbf{X}, \mathbf{Z}) = \int p(\mathbf{f}|\mathbf{X}, \mathbf{Z}, \mathbf{t})p(\mathbf{t}|\mathbf{Z})d\mathbf{t} \quad (94)$$

$$= \mathcal{N}(\boldsymbol{\mu}_f, \boldsymbol{\Sigma}_f) \quad (95)$$

$$\text{where} \quad (96)$$

$$\boldsymbol{\mu}_f = 0 \quad (97)$$

$$\boldsymbol{\Sigma}_f = \mathbf{K}_{NM} \tilde{\mathbf{K}}^{-1} \mathbf{K}_{NM}^T + \Lambda \quad (98)$$

Now we can use MLE-II for getting the values  $\mathbf{Z}$

$$\tilde{\mathbf{Z}} = \arg \max_{\mathbf{Z}} \log p(\mathbf{f}|\mathbf{X}, \mathbf{Z}) \quad (99)$$

$$= \arg \min_{\mathbf{Z}} (\log |\boldsymbol{\Sigma}_f| + \mathbf{f}^T \boldsymbol{\Sigma}_f^{-1} \mathbf{f}) \quad (100)$$

The above objective function can be solved using gradient ascent.