

## COMP4670/8600: Statistical Machine Learning

### Contribution statement:

Ziyang Chen (u6908560) mainly contributed on Section 1.

Han Zhang (u7235649) mainly contributed on Section 2.

### Answer to Question 1.1

Because  $X$  is continuous, therefore, we can get following equation in distribution  $\mathcal{D}$  according to Markov's Inequality and the given equation (1.1) in question sheet.

$$\mathcal{R}_{\mathcal{D}}[f] = \mathbb{E}_{\mathcal{D}}[\ell(f(X), Y)] = \iint \ell(f(x), y) p(x, y) dx dy \quad (1)$$

where  $q(x, y)$  is the density function of  $\mathcal{D}$ .

Combining the given equation  $0 \leq \ell(y, y') \leq 1$ , we can get

$$\iint \ell(f(x), y) p(x, y) dx dy \leq \iint p(x, y) dx dy \quad (2)$$

Therefore, we can get following equation according to equation 1 and equation 2.

$$\mathcal{R}_{\mathcal{D}}[f] \leq \iint p(x, y) dx dy \quad (3)$$

Similarly, we can get the following equation in distribution  $\tilde{\mathcal{D}}$

$$\mathcal{R}_{\tilde{\mathcal{D}}}[f] \leq \iint \tilde{p}(x, y) dx dy \quad (4)$$

where  $\tilde{p}(x, y)$  is the density function of  $\tilde{\mathcal{D}}$ .

Combining equation 3 and equation 4, we can get

$$|\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f]| \leq \iint |p(x, y) - \tilde{p}(x, y)| dx dy \quad (5)$$

Therefore, according to the given equation (1.6) about the definition of  $D_{TV}(p, \tilde{p})$ , we can get

$$|\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f]| \leq D_{TV}(p, \tilde{p}) \quad (6)$$

### Answer to Question 1.2

Using the triangle inequality, we can get:

$$\left| \mathcal{R}_{\mathcal{D}}[f] - \hat{\mathcal{R}}_{\tilde{\mathcal{S}}}[f] \right| = \left| \mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f] + \mathcal{R}_{\tilde{\mathcal{D}}}[f] - \hat{\mathcal{R}}_{\tilde{\mathcal{S}}}[f] \right| \leq \left| \mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f] \right| + \left| \mathcal{R}_{\tilde{\mathcal{D}}}[f] - \hat{\mathcal{R}}_{\tilde{\mathcal{S}}}[f] \right| \quad (7)$$

According to the above proof in Question 1.1, we can have  $|\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f]| \leq D_{TV}(p, \tilde{p})$ , therefore, combining equation 7, we can get

$$\left| \mathcal{R}_{\mathcal{D}}[f] - \hat{\mathcal{R}}_{\tilde{\mathcal{S}}}[f] \right| \leq D_{TV}(p, \tilde{p}) + \left| \mathcal{R}_{\tilde{\mathcal{D}}}[f] - \hat{\mathcal{R}}_{\tilde{\mathcal{S}}}[f] \right| \quad (8)$$

According to the equation (A, 1):  $\Pr\{|\hat{\mu} - \mu_Z| \geq \epsilon\} \leq 2 \exp(-2N\epsilon^2)$  in Proposition 1,  $\mu_Z$  is the true mean and  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N Z_i$  is the sample mean. Therefore,  $\mu_Z$  and  $\hat{\mu}$  correspond to  $\mathcal{R}_{\mathcal{D}}[f]$  and  $\hat{\mathcal{R}}_{\mathcal{S}}[f]$  in equation 8.

Therefore, equation (A, 1) can be written

$$\Pr\{|\hat{\mathcal{R}}_{\mathcal{S}}[f] - \mathcal{R}_{\mathcal{D}}[f]| \geq \epsilon\} \leq 2 \exp(-2N\epsilon^2) \quad (9)$$

It also can be written

$$\Pr\{|\mathcal{R}_{\mathcal{D}}[f] - \hat{\mathcal{R}}_{\mathcal{S}}[f]| \leq \epsilon\} \geq 1 - 2 \exp(-2N\epsilon^2) \quad (10)$$

Let  $\delta = 2 \exp(-2N\epsilon^2)$ , we can get  $\epsilon = \sqrt{\frac{\log(\frac{2}{\delta})}{2N}}$ . Through losing bound, we can get  $\epsilon \leq \sqrt{\frac{\log(\frac{1}{\delta})}{N}}$ . Therefore, the equation 10 can be changed into

$$\Pr\{|\mathcal{R}_{\mathcal{D}}[f] - \hat{\mathcal{R}}_{\mathcal{S}}[f]| \leq \sqrt{\frac{\log(\frac{1}{\delta})}{N}}\} \geq 1 - \delta \quad (11)$$

According to the model complexity and uniform convergence, we can get, with probability  $1 - \delta$ ,  $\forall f \in \mathcal{H} : \Delta_{gen}(f) \leq \sqrt{\frac{\ln(|\mathcal{H}|) + \ln(\frac{1}{\delta})}{N}}$ .

So, we can have

$$\Pr\{|\mathcal{R}_{\mathcal{D}}[f] - \hat{\mathcal{R}}_{\mathcal{S}}[f]| \leq \sqrt{\frac{\ln(|\mathcal{H}|) + \ln(\frac{1}{\delta})}{N}}\} \geq 1 - \delta \quad (12)$$

Therefore, adding  $D_{TV}(p, \tilde{p})$  into two sides of the equation within curly bracket, we can get

$$\Pr\{D_{TV}(p, \tilde{p}) + |\mathcal{R}_{\mathcal{D}}[f] - \hat{\mathcal{R}}_{\mathcal{S}}[f]| \leq D_{TV}(p, \tilde{p}) + \sqrt{\frac{\ln(|\mathcal{H}|) + \ln(\frac{1}{\delta})}{N}}\} \geq 1 - \delta \quad (13)$$

Combining equation 8 and equation 13, we can get

$$\Pr\{|\mathcal{R}_{\mathcal{D}}[f] - \hat{\mathcal{R}}_{\mathcal{S}}[f]| \leq D_{TV}(p, \tilde{p}) + \sqrt{\frac{\ln(|\mathcal{H}|) + \ln(\frac{1}{\delta})}{N}}\} \geq 1 - \delta \quad (14)$$

### Answer to Question 1.3

Comparing with the equation 14 and the given equation without noisy data, we can see the noisy data will increase the generalisation bound and the difference of them will be  $D_{TV}(p, \tilde{p})$ , which is the risk distance between true data and noisy data. The rate will increase when it decreases with the number of data points.

### Answer to Question 1.4

For a noisy point  $y \in \mathcal{Y}$ , it may be got from 2 situations.

1. it was  $-y$  before, when the label noise happened, it changed into  $y$ .
2. it was  $y$  before and there is no change when the label noise happened;

For the situation 1, we can get the probability is

$$\Pr(\tilde{Y} = y | Y = -y) \cdot \Pr(Y = -y | X) \quad (15)$$

For the situation 2, we can get the probability is

$$\Pr(\tilde{Y} = y|Y = y) \cdot \Pr(Y = y|X) \quad (16)$$

Therefore,

$$\begin{aligned} \Pr(\tilde{Y} = y|X) &= \Pr(\tilde{Y} = y|Y = -y) \cdot \Pr(Y = -y|X) + \Pr(\tilde{Y} = y|Y = y) \cdot \Pr(Y = y|X) \\ &= \Pr(\tilde{Y} = y|Y = -y) \cdot (1 - \Pr(Y = y|X)) + \Pr(\tilde{Y} = y|Y = y) \cdot \Pr(Y = y|X) \end{aligned} \quad (17)$$

Because the label flipping probability  $\Pr(\tilde{Y} = -y|Y = y) = \sigma_y$ , when  $-y$  changes into  $y$ , the label flipping probability in this situation is  $\Pr(\tilde{Y} = y|Y = -y) = \sigma_{-y}$ .

Therefore, equation 17 can be written as

$$\Pr(\tilde{Y} = y|X) = \sigma_{-y} \cdot (1 - \Pr(Y = y|X)) + (1 - \sigma_y) \cdot \Pr(Y = y|X) \quad (18)$$

### Answer to Question 1.5

According to the Question 1.2, we already have equation 14. Combining the equation (1.6) in question sheet, we can get

$$\Pr\left\{\left|\mathcal{R}_{\mathcal{D}}[f] - \widehat{\mathcal{R}}_{\tilde{\mathcal{S}}}[f]\right| \leq \iint |p(x, y) - \tilde{p}(x, y)| \, dx dy + \sqrt{\frac{\ln(|\mathcal{H}|) + \ln(\frac{1}{\delta})}{N}}\right\} \geq 1 - \delta \quad (19)$$

In this question, we can get

$$p(x, y) = \Pr(X = x, Y = y) = \Pr(X)\Pr(Y = y|X = x) = \Pr(X)\Pr(Y = y|X) \quad (20)$$

$$\tilde{p}(x, y) = \Pr(\tilde{X} = x, \tilde{Y} = y) = \Pr(\tilde{X})\Pr(\tilde{Y} = y|\tilde{X} = x) = \Pr(\tilde{X})\Pr(\tilde{Y} = y|\tilde{X}) \quad (21)$$

According to the **LN2** property, we can know for all  $x \in \mathcal{X}$

$$\Pr(\tilde{X} = x) = \Pr(X = x) \quad \text{similarly} \quad \Pr(\tilde{X}) = \Pr(X) \quad (22)$$

which can imply

$$\Pr(\tilde{X} = x|\tilde{Y} = y) = \Pr(X = x|\tilde{Y} = y) \quad \text{similarly} \quad \Pr(\tilde{X}|\tilde{Y} = y) = \Pr(X|\tilde{Y} = y) \quad (23)$$

Using the Bayes theorem and combining equation 22 and equation 23, we can get

$$\Pr(\tilde{Y} = y|\tilde{X}) = \frac{\Pr(\tilde{Y} = y) \cdot \Pr(\tilde{X}|\tilde{Y} = y)}{\Pr(\tilde{X})} = \frac{\Pr(\tilde{Y} = y) \cdot \Pr(X|\tilde{Y} = y)}{\Pr(X)} = \Pr(\tilde{Y} = y|X) \quad (24)$$

Therefore, combining equation 20, equation 21, equation 23, equation 24 and the equation 18 in

question 1.4, and  $\sigma_y = \sigma_{-y} = \sigma$  in the symmetric label noise, we can get

$$\begin{aligned}
p(x, y) - \tilde{p}(x, y) &= \Pr(X) \cdot \Pr(Y = y|X) - \Pr(\tilde{X}) \cdot \Pr(\tilde{Y} = y|\tilde{X}) \\
&= \Pr(X) \cdot \Pr(Y = y|X) - \Pr(X) \cdot \Pr(\tilde{Y} = y|X) \\
&= \Pr(X) \cdot \Pr(Y = y|X) - \Pr(X) \cdot [\sigma_{-y} \cdot (1 - \Pr(Y = y|X)) + (1 - \sigma_y) \cdot \Pr(Y = y|X)] \\
&= \Pr(X) \cdot \Pr(Y = y|X) - \Pr(X) \cdot [\sigma \cdot (1 - \Pr(Y = y|X)) + (1 - \sigma) \cdot \Pr(Y = y|X)] \\
&= \Pr(X) \cdot (2\sigma \Pr(Y = y|X) - \sigma) \\
&= 2\sigma \Pr(X) \Pr(Y = y|X) - \sigma \Pr(X)
\end{aligned} \tag{25}$$

Because  $0 \leq \Pr \leq 1$  and  $0 \leq \Pr(Y = y|X) \leq 1$ . Thus, we can get

$$p(x, y) - \tilde{p}(x, y) = 2\sigma \Pr(X) \Pr(Y = y|X) - \sigma \Pr(X) \leq 2\sigma \tag{26}$$

Therefore, according to  $0 \leq \sigma \leq 1$ , we can get

$$\iint |p(x, y) - \tilde{p}(x, y)| \, dx dy \leq \iint 2\sigma \, dx dy = 2\sigma \tag{27}$$

Combining the equation 19, we can get

$$\Pr\left\{\left|\mathcal{R}_{\mathcal{D}}[f] - \widehat{\mathcal{R}}_{\tilde{\mathcal{S}}}[f]\right| \leq 2\sigma + \sqrt{\frac{\ln(|\mathcal{H}|) + \ln(\frac{1}{\delta})}{N}}\right\} \geq 1 - \delta \tag{28}$$

### Answer to Question 2.1

The *Matern32* and *Matern52* functions in *kernels.py* have been implemented and submitted.

### Answer to Question 2.2

The *optimisation()*, *fit()*, *predict()* and *log\_marginal\_likelihood()* in *gp.py* have been implemented and submitted.

### Answer to Question 2.3

According to the question sheet, we can get

$$f(\mathbf{x}) = \mu(\mathbf{x}) + \sigma(\mathbf{x})Z \tag{29}$$

where  $\mu(\mathbf{x})$  is the mean of the  $f(\mathbf{x})$  and  $\sigma(\mathbf{x})$  is the standard deviation of the  $f(\mathbf{x})$ .

Seeing the equation (2.11), we can get

$$EI(\mathbf{x}) = \mathbb{E}(I(\mathbf{x})) = \int_{-\infty}^{+\infty} I(\mathbf{x}) \Pr(Z) dZ = \int_{-\infty}^{+\infty} I(\mathbf{x}) \phi(Z) dZ \tag{30}$$

where  $\phi(\cdot)$  is the normal probability distribution functions, which is

$$\phi(Z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{Z^2}{2}} \tag{31}$$

According to the equation (2.8), let  $I(\mathbf{x}) = f(\mathbf{x}) - f(\mathbf{x}^+) - \xi = 0$ . Combining to equation 29, we can get

$$Z_0 = \frac{f(\mathbf{x}^+) + \xi - \mu(\mathbf{x})}{\sigma(\mathbf{x})} \tag{32}$$

Continuing in equation 30, we can get

$$\begin{aligned}\mathbb{E}(I(\mathbf{x})) &= \int_{-\infty}^{+\infty} I(\mathbf{x})\phi(Z)dZ \\ &= \int_{-\infty}^{Z_0} I(\mathbf{x})\phi(Z)dZ + \int_{Z_0}^{+\infty} I(\mathbf{x})\phi(Z)dZ\end{aligned}\quad (33)$$

Because  $I(\mathbf{x}) = 0$ ,  $\int_{-\infty}^{Z_0} I(\mathbf{x})\phi(Z)dZ = 0$ . Therefore, we can write equation 33 as

$$\begin{aligned}\mathbb{E}(I(\mathbf{x})) &= \int_{Z_0}^{+\infty} I(\mathbf{x})\phi(Z)dZ \\ &= \int_{Z_0}^{+\infty} (\mu(\mathbf{x}) + \sigma(\mathbf{x})Z - f(\mathbf{x}^+) - \xi)\phi(Z)dZ \\ &= \int_{Z_0}^{+\infty} (\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi)\phi(Z)dZ + \int_{Z_0}^{+\infty} (\sigma(\mathbf{x})Z)\phi(Z)dZ \\ &= (\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi) \int_{Z_0}^{+\infty} \phi(Z)dZ + \sigma(\mathbf{x}) \int_{Z_0}^{+\infty} Z\phi(Z)dZ \\ &= (\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi)(1 - \Phi(Z_0)) + \sigma(\mathbf{x}) \int_{Z_0}^{+\infty} Z \frac{1}{\sqrt{2\pi}} e^{-\frac{Z^2}{2}} dZ \\ &= (\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi)(1 - \Phi(Z_0)) + \frac{\sigma(\mathbf{x})}{\sqrt{2\pi}} \int_{Z_0}^{+\infty} Z e^{-\frac{Z^2}{2}} dZ \\ &= (\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi)(1 - \Phi(Z_0)) + \frac{\sigma(\mathbf{x})}{\sqrt{2\pi}} [e^{-\frac{Z^2}{2}}]_{Z_0}^{+\infty} \\ &= (\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi)(1 - \Phi(Z_0)) + \sigma(\mathbf{x})\phi(Z_0)\end{aligned}\quad (34)$$

where  $\Phi(\cdot)$  is the normal cumulative distribution functions.

According to the equation (B.6) in question sheet, we can get  $1 - \Phi(-Z) = \Phi(Z)$ . Because  $\Phi(\cdot)$  is the normal cumulative distribution functions, we can also get  $\Phi(Z) = \Phi(-Z)$ . Therefore, we can get the following equation from equation 34

$$\mathbb{E}(I(\mathbf{x})) = (\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi)\Phi(Z_0) + \sigma(\mathbf{x})\phi(Z_0)\quad (35)$$

Combining all the information and equation (2.10) in question sheet, we can get

$$\mathbb{E}(I(\mathbf{x})) = \begin{cases} (\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi)\Phi(Z) + \sigma(\mathbf{x})\phi(Z) & \text{if } \sigma(\mathbf{x}) > 0 \\ 0 & \text{if } \sigma(\mathbf{x}) = 0 \end{cases}\quad (36)$$

where

$$Z = \begin{cases} \frac{\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi}{\sigma(\mathbf{x})} & \text{if } \sigma(\mathbf{x}) > 0 \\ 0 & \text{if } \sigma(\mathbf{x}) = 0 \end{cases}\quad (37)$$

#### Answer to Question 2.4

The Probability of Improvement (PI) and Expected Improvement (EI) acquisitions functions in *acquisitions.py* have been implemented and submitted.

#### Answer to Question 2.5

#### Answer to Question 2.6

It is usually difficult to find the optimal hyperparameter combination among complex hyperparameters when the dimensions are high. The hyperparameter optimization approach can make this process easy.