

Let $f \in \mathcal{H}$ with a bounded loss function (as per Eq. (1.3)). Prove that

$$|\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f]| \leq D_{TV}(p, \tilde{p}), \quad (1.5)$$

where p, \tilde{p} are the density functions of $\mathcal{D}, \tilde{\mathcal{D}}$, respectively; and

$$D_{TV}(p, \tilde{p}) \stackrel{\text{def}}{=} \iint |p(x, y) - \tilde{p}(x, y)| dx dy. \quad (1.6)$$

$$\begin{aligned} |\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f]| &= \left| \mathbb{E}_{(x, y) \sim \mathcal{D}} [\ell(f(x), y)] - \mathbb{E}_{(x, y) \sim \tilde{\mathcal{D}}} [\ell(f(x), y)] \right| \\ &= \left| \iint \ell(f(x), y) \cdot p(x, y) dx dy - \iint \ell(f(x), y) \cdot \tilde{p}(x, y) dx dy \right| \\ &= \left| \iint \ell(f(x), y) \cdot p(x, y) - \ell(f(x), y) \cdot \tilde{p}(x, y) dx dy \right| \\ &= \left| \iint \ell(f(x), y) \cdot (p(x, y) - \tilde{p}(x, y)) dx dy \right| \\ \because \ell(y, y) \in [0, 1] \quad &\left| \mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f] \right| \\ \therefore 0 = |0| \leqslant \left| \iint \ell(f(x), y) \cdot (p(x, y) - \tilde{p}(x, y)) dx dy \right| \leqslant \left| \iint (p(x, y) - \tilde{p}(x, y)) dx dy \right| \\ \because \forall (x, y) \in \mathcal{D} \quad |p(x, y) - \tilde{p}(x, y)| \geq p(x, y) - \tilde{p}(x, y) \\ \therefore \iint |p(x, y) - \tilde{p}(x, y)| dx dy \geq \iint p(x, y) - \tilde{p}(x, y) dx dy \\ \because \forall (x, y) \in \mathcal{D} \quad |p(x, y) - \tilde{p}(x, y)| \geq -(p(x, y) - \tilde{p}(x, y)) \\ \therefore \iint |p(x, y) - \tilde{p}(x, y)| dx dy \geq \iint -(p(x, y) - \tilde{p}(x, y)) dx dy = -\iint (p(x, y) - \tilde{p}(x, y)) dx dy \\ \therefore -\iint |p(x, y) - \tilde{p}(x, y)| dx dy \leq \iint |p(x, y) - \tilde{p}(x, y)| dx dy \leq \iint |p(x, y) - \tilde{p}(x, y)| dx dy \\ \therefore \left| \iint |p(x, y) - \tilde{p}(x, y)| dx dy \right| \leq \iint |p(x, y) - \tilde{p}(x, y)| dx dy = D_{TV}(p, \tilde{p}) \Rightarrow |\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f]| \leq D_{TV}(p, \tilde{p}) \end{aligned}$$

Question 1.2: Uniform Bound for Noisy Sample Gap

(15 Points)

Let \mathcal{H} be a finite hypothesis class. Prove that given a finite sample $\mathcal{S} \sim \mathcal{D}^N$, for $\delta > 0$

$$\Pr \left\{ \forall f \in \mathcal{H} : |\mathcal{R}_{\mathcal{D}}[f] - \widehat{\mathcal{R}}_{\mathcal{S}}[f]| \leq D_{TV}(p, \tilde{p}) + \sqrt{\frac{\ln |\mathcal{H}| + \ln(1/\delta)}{N}} \right\} \geq 1 - \delta. \quad (1.7)$$

$$P(\forall f \in \mathcal{H} : |\mathcal{R}_{\mathcal{D}}[f] - \widehat{\mathcal{R}}_{\mathcal{S}}[f]| \leq t) \geq 1 - \delta, \text{ To prove } t = D_{TV}(p, \tilde{p}) + \sqrt{\frac{\ln |\mathcal{H}| + \ln(1/\delta)}{N}}$$

$$|\mathcal{R}_{\mathcal{D}}[f] - \widehat{\mathcal{R}}_{\mathcal{S}}[f]| = |\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f] + \mathcal{R}_{\tilde{\mathcal{D}}}[f] - \widehat{\mathcal{R}}_{\mathcal{S}}[f]| \leq |\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f]| + |\mathcal{R}_{\tilde{\mathcal{D}}}[f] - \widehat{\mathcal{R}}_{\mathcal{S}}[f]|$$

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

$$\therefore |\mathcal{R}_{\mathcal{D}}[f] - \widehat{\mathcal{R}}_{\mathcal{S}}[f]| \leq D_{TV}(p, \tilde{p}) + |\mathcal{R}_{\tilde{\mathcal{D}}}[f] - \widehat{\mathcal{R}}_{\mathcal{S}}[f]|$$

From the Proposition 1, we know that $P(|\mathcal{R}_{\tilde{\mathcal{D}}}[f] - \widehat{\mathcal{R}}_{\mathcal{S}}[f]| \geq \epsilon) \leq 2 \exp(-2N\epsilon^2)$

Then, $P(|\mathcal{R}_{\tilde{\mathcal{D}}}[f] - \widehat{\mathcal{R}}_{\mathcal{S}}[f]| \geq \epsilon) \geq 1 - 2 \exp(-2N\epsilon^2)$

$$\therefore 1 - 2 \exp(-2N\epsilon^2) = 1 - \delta$$

$$\therefore \epsilon = \sqrt{\frac{\ln(2/\delta)}{2N}} \leq \sqrt{\frac{\ln 2 + \ln(1/\delta)}{N}}$$

From piazza, Josh Nguyen said that we can assume $|\mathcal{H}| \geq 2$ and $\delta \in (0, 1)$

$$\text{so } \ln |\mathcal{H}| \geq \ln 2$$

$$\therefore \epsilon \leq \sqrt{\frac{\ln 2 + \ln(1/\delta)}{N}} \leq \sqrt{\frac{\ln |\mathcal{H}| + \ln(1/\delta)}{N}} \quad \forall f \in \mathcal{H}$$

$$\therefore |R_{\mathcal{D}}[f] - \hat{R}_{\mathcal{S}}[f]| \leq \epsilon \leq \sqrt{\frac{\ln |\mathcal{H}| + \ln(1/\delta)}{N}} \quad \forall f \in \mathcal{H}$$

$$\text{so } |R_{\mathcal{D}}[f] - \hat{R}_{\mathcal{S}}[f]| \leq D_{\text{TV}}(P, \tilde{P}) + |R_{\mathcal{D}}[f] - \hat{R}_{\mathcal{S}}[f]| \leq D_{\text{TV}}(P, \tilde{P}) + \sqrt{\frac{\ln |\mathcal{H}| + \ln(1/\delta)}{N}} \quad \forall f \in \mathcal{H}$$

$$\text{so } P(\forall f \in \mathcal{H} : |R_{\mathcal{D}}[f] - \hat{R}_{\mathcal{S}}[f]| \leq D_{\text{TV}}(P, \tilde{P}) + \sqrt{\frac{\ln |\mathcal{H}| + \ln(1/\delta)}{N}}) \geq 1 - \delta$$

Question 1.3: Interpreting the Uniform Bound for Noisy Sample Gap (10 Points)

Compare Eq. (1.7) with the bound we found in the lecture in the absence of noisy data, namely, for any $\delta > 0$,

$$\Pr \left\{ \forall f \in \mathcal{H} : |R_{\mathcal{D}}[f] - \hat{R}_{\mathcal{S}}[f]| \leq \sqrt{\frac{\ln |\mathcal{H}| + \ln(1/\delta)}{N}} \right\} \geq 1 - \delta.$$

What is the effect of noisy data on the generalisation bound and the rate at which it decreases with the number of data points?

Compare the Uniform Bound for the Noisy Sample Gap with the bound found in lecture.

We find that, the noisy data increases the generalisation bound by $D_{\text{TV}}(P, \tilde{P})$, which measure the total difference between clean data distribution pdf and the noisy data distribution pdf. If the $\tilde{P} = P$, then Bound for the Noisy Sample Gap equals to the bound introduced in the lecture. If it has more noisy, the bound will get bigger and bigger.

The rate at which it decreases with the number of data points is $\sqrt{\frac{c}{N}}$, because $D_{\text{TV}}(P, \tilde{P})$ is constant and the term $\sqrt{\frac{\ln |\mathcal{H}| + \ln(1/\delta)}{N}}$ tells me $\sqrt{\frac{c}{N}}$, so the rate is $\sqrt{\frac{c}{N}}$

Question 1.4: Noisy Posterior for Label Noise

(5 Points)

Given a clean distribution \mathcal{D} and a noisy version of the distribution $\tilde{\mathcal{D}}$ under asymmetric label noise with flipping probabilities $\sigma_{-1}, \sigma_+ \in [0, 1]$, show that for all $y \in \mathcal{Y}$

$$\Pr(\tilde{Y} = y | X) = \sigma_{-y} \cdot (1 - \Pr(Y = y | X)) + (1 - \sigma_y) \cdot \Pr(Y = y | X). \quad (1.8)$$

From (LN3), we know $\tilde{Y} \perp\!\!\!\perp X | Y$, so we know $P(\tilde{Y}, x | Y) = P(\tilde{Y} | Y) P(x | Y)$ and $P(\tilde{Y} | x, Y) = P(\tilde{Y} | Y) - \text{(LN3.1)}$

for all $y \in \mathcal{Y} = \{-1, +1\}$

$$\begin{aligned} P(\tilde{Y} = y | X) &= \sum_{\tilde{y}} P(\tilde{Y} = y, Y = \tilde{y} | X) = P(\tilde{Y} = y, Y = y | X) + P(\tilde{Y} = y, Y = -y | X) \\ &= P(\tilde{Y} = y | Y = y) \cdot P(Y = y | X) + P(\tilde{Y} = y | Y = -y) \cdot P(Y = -y | X) \quad (\text{LN3.1}) \\ &= (1 - \varepsilon_y) \cdot P(Y = y | X) + \varepsilon_y \cdot P(Y = -y | X) \quad (\text{LN1}) \end{aligned}$$

\therefore We know Y only can have two values, namely, -1 or $+1$.

$$\therefore P(Y = y | X) + P(Y = -y | X) = 1 \equiv P(Y = -y | X) = 1 - P(Y = y | X)$$

$$\therefore P(\tilde{Y} = y | X) = (1 - \varepsilon_y) \cdot P(Y = y | X) + \varepsilon_y (1 - P(Y = y | X))$$

Question 1.5: Uniform Bound for Symmetric Label Noise

(10 Points)

Suppose that a noisy distribution $\tilde{\mathcal{D}}$ is given by symmetric label noise, with flipping $\sigma \in [0, 1]$, applied to clean distribution \mathcal{D} . Prove that given a finite sample $\mathcal{S} \sim \mathcal{D}^N$, for $\delta > 0$

$$\Pr \left\{ \forall f \in \mathcal{H} : |\mathcal{R}_{\mathcal{D}}[f] - \hat{\mathcal{R}}_{\tilde{\mathcal{S}}}[f]| \leq 2\sigma + \sqrt{\frac{\ln |\mathcal{H}| + \ln(1/\delta)}{N}} \right\} \geq 1 - \delta. \quad (1.9)$$

Hint: Start by bounding $|\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f]|$ like in Question 1.1. Then reuse the technique in Question 1.2.

First, try to bound $|\mathcal{R}_{\mathcal{D}}[f] - \mathcal{R}_{\tilde{\mathcal{D}}}[f]| =$

$$= \left| \iint l(f(x, y) \cdot p(x, y), dx dy - \iint l(f(\tilde{x}, \tilde{y}) \cdot p(\tilde{x}, \tilde{y}), d\tilde{x} d\tilde{y} \right| =$$

* (supplementary proof of $P(\tilde{Y} | \tilde{x}) = P(\tilde{Y} | x)$)
 LN2 says that $p(x | \tilde{x}) = p(x | \tilde{x})$
 so $P(\tilde{Y} | \tilde{x}) = \frac{P(\tilde{Y} | \tilde{x}) \cdot P(\tilde{x})}{P(\tilde{x})} = \frac{P(x | \tilde{x}) \cdot P(\tilde{x})}{P(x)}$ (LN2)
 $= P(\tilde{Y} | x)$
 $\therefore P(\tilde{Y} | \tilde{x}) = P(\tilde{Y} | x)$

$$= \left| \int l(f(x, +1) \cdot p(x, +1) + l(f(x, -1) \cdot p(x, -1), dx - \int l(f(\tilde{x}, +1) \cdot p(\tilde{x}, +1) + l(f(\tilde{x}, -1) \cdot p(\tilde{x}, -1), d\tilde{x} \right|$$

$$= \left| \int l(f(x, +1) \cdot p(y=+1 | x) p(x) + l(f(x, -1) \cdot p(y=-1 | x) p(x) dx - \dots \right.$$

$$\left. \dots \int l(f(\tilde{x}, +1) \cdot p(\tilde{y}=+1 | \tilde{x}) p(\tilde{x}) + l(f(\tilde{x}, -1) \cdot p(\tilde{y}=-1 | \tilde{x}) p(\tilde{x}) d\tilde{x} \right|$$

(LN2 and supplementary proof)

$$= \int p(x) [\ell(f(x), +1) \cdot p(y=+1|x) + \ell(f(x), -1) \cdot p(y=-1|x)] dx - \dots$$

$$\cdots \int p(x) [\ell(f(x), +1) [6 + (1-2\sigma) p(Y=+1|x)] + \ell(f(x), -1) [6 + (1-2\sigma) p(Y=-1|x)]] dx$$

(E.91.8 and 6-1-6+1=6)

$$= \int p(x) [\ell(f(x), +1)(2\sigma p(Y=+1|x) - \sigma) + \ell(f(x), -1)(2\sigma p(Y=-1|x) - \sigma)] dx$$

$$\therefore 2\sigma p(Y=+1|x) - \sigma \leq 2\sigma p(Y=+1|x) \quad (\text{because } 0 \leq \sigma \leq 1 \text{ and } 2\sigma p(Y=+1|x) \geq 0)$$

$$\therefore \left| \int p(x) [\ell(f(x), +1)(2\sigma p(Y=+1|x) - \sigma) + \ell(f(x), -1)(2\sigma p(Y=-1|x) - \sigma)] dx \right| \leq$$

$$\left| \int p(x) [\ell(f(x), +1) \cdot 2\sigma p(Y=+1|x) + \ell(f(x), -1) \cdot 2\sigma p(Y=-1|x)] dx \right|$$

$$\therefore \ell(y, y') \in [0, 1]$$

$$\therefore 0 \leq \ell(f(x), +1) \cdot 2\sigma p(Y=+1|x) \leq 2\sigma p(Y=+1|x) \text{ and } 0 \leq \ell(f(x), -1) \cdot 2\sigma p(Y=-1|x) \leq 2\sigma p(Y=-1|x)$$

$$\therefore \left| \int p(x) [\ell(f(x), +1) \cdot 2\sigma p(Y=+1|x) + \ell(f(x), -1) \cdot 2\sigma p(Y=-1|x)] dx \right| \leq \left| \int p(x) [2\sigma p(Y=+1|x) + 2\sigma p(Y=-1|x)] dx \right|$$

$$= \left| 2\sigma \int p(x) (p(Y=+1|x) + p(Y=-1|x)) dx \right| = |2\sigma| = 2\sigma \quad \left(\sum_y p(Y=y|x) = 1, \int p(x) dx = 1 \right)$$

$$\therefore |R_D[f] - R_B[f]| \leq 2\sigma$$

$$P(\forall f \in H : |R_D[f] - \hat{R}_3[f]| \leq t) \geq 1 - \delta, \text{ to prove } t = 2\sigma + \sqrt{\frac{\ln|H| + \ln(1/\delta)}{N}}$$

$$\therefore |R_D[f] - \hat{R}_3[f]| = |R_D[f] - R_B[f] + R_B[f] - \hat{R}_3[f]| \leq |R_D[f] - R_B[f]| + |R_B[f] - \hat{R}_3[f]|$$

$$\therefore |R_D[f] - R_B[f]| \leq 2\sigma$$

$$\therefore |R_D[f] - \hat{R}_3[f]| \leq 2\sigma + |R_B[f] - \hat{R}_3[f]|$$

From the Proposition 1, we know that $P(|R_B[f] - \hat{R}_3[f]| \geq \epsilon) \leq 2e^{-2N\epsilon^2}$

$$\text{Then } |R_D[f] - \hat{R}_3[f]| \leq \sqrt{\frac{\ln|H| + \ln(1/\delta)}{N}} \quad (\text{proved in Q1.2})$$

$$\text{so } |R_D[f] - \hat{R}_3[f]| \leq 2\sigma + |R_B[f] - \hat{R}_3[f]| \leq 2\sigma + \sqrt{\frac{\ln|H| + \ln(1/\delta)}{N}} \quad \forall f \in H$$

$$\text{so } P(\forall f \in H : |R_D[f] - \hat{R}_3[f]| \leq 2\sigma + \sqrt{\frac{\ln|H| + \ln(1/\delta)}{N}}) \geq 1 - \delta$$

Question 2.3: Expected Improvement

Derive mathematically the result of the Expected Improvement (EI) acquisition function in Equation 2.13. That is show that,

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

$$f(x) \geq f(x^+) - \xi$$

where Z is defined in Equation 2.14.

$$\textcircled{1} \quad \text{if } \sigma(x) > 0$$

$$f(x) \sim N(\mu(x), \sigma^2(x)|x)$$

$$I(x) = \max\{0, f(x) - f(x^+) - \xi\}$$

$$EI(x) = \mathbb{E}_{f(x) \sim N(\mu(x), \sigma^2(x)|x)}[I(x)]$$

$$\text{let } z_f = \frac{f(x) - \mu(x)}{\sigma(x)}, \text{ so } z_f \sim N(0, 1), z^+ = \frac{f(x^+) + \xi - \mu(x)}{\sigma(x)}, z = -z^+ = \frac{\mu(x) - f(x^+) - \xi}{\sigma(x)}$$

$$EI(x) = \mathbb{E}_{z_f \sim N(0, 1)}[I(x)]$$

$$EI(x) = \int_{-\infty}^{+\infty} I(x) \phi(z_f) dz_f$$

$$= \int_{z^+}^{+\infty} [\bar{z}_f \sigma(x) + \mu(x) - f(x^+) - \xi] \phi(z_f) dz_f + 0 \quad (f(x) = \bar{z}_f \sigma(x) + \mu(x))$$

$$= \int_{z^+}^{+\infty} [\mu(x) - f(x^+) - \xi] \phi(z_f) dz_f - \sigma(x) \int_{z^+}^{+\infty} \bar{z}_f \cdot \phi(z_f) dz_f$$

$$= [\mu(x) - f(x^+) - \xi] \cdot \int_{z^+}^{+\infty} \phi(z_f) dz_f - \sigma(x) \cdot \int_{z^+}^{+\infty} \bar{z}_f \cdot \phi(z_f) dz_f$$

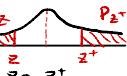
$$= [\mu(x) - f(x^+) - \xi] \cdot (1 - \Phi(z^+)) - \sigma(x) \cdot \int_{z^+}^{+\infty} \bar{z}_f \cdot \phi(z_f) dz_f \quad (1 - \Phi(z^+) = \Phi(z))$$

$$= [\mu(x) - f(x^+) - \xi] \cdot \Phi(z) - \sigma(x) \int_{z^+}^{+\infty} \bar{z}_f \cdot \phi(z_f) dz_f$$

$$\therefore \phi(z_f) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z_f^2}{2}}, \frac{d}{dz_f} \exp(-\frac{z_f^2}{2}) = -\bar{z}_f \cdot \exp(-\frac{z_f^2}{2})$$

$$= [\mu(x) - f(x^+) - \xi] \cdot \Phi(z) - \sigma(x) \int_{z^+}^{+\infty} \bar{z}_f \cdot \frac{1}{\sqrt{2\pi}} \exp(-\frac{z_f^2}{2}) dz_f$$

$$= [\mu(x) - f(x^+) - \xi] \cdot \Phi(z) + \frac{\sigma(x)}{\sqrt{2\pi}} \cdot \int_{-\infty}^z -\bar{z}_f \cdot \exp(-\frac{z_f^2}{2}) dz_f \quad (\because \Phi(z) = \Phi(z^+) \quad z = -z^+)$$



$$= [\mu(x) - f(x^+) - \xi] \cdot \Phi(z) + \frac{\sigma(x)}{\sqrt{2\pi}} \left[\exp(-\frac{z^2}{2}) - 0 \right]$$

$$= [\mu(x) - f(x^+) - \xi] \cdot \Phi(z) + \sigma(x) \cdot \phi(z)$$

$$\textcircled{2} \quad \sigma(x) = 0 \Rightarrow z = 0 \Rightarrow 0 = z^+ = f(x^+) + \xi - \mu(x) \Rightarrow z = 0 = \mu(x) - f(x^+) - \xi$$

$$EI(x) = \mathbb{E}_{z_f \sim N(0, 1)}[I(x)]$$

$$EI(x) = \int_{-\infty}^{+\infty} I(x) \phi(z_f) dz_f$$

$$= \int_0^{+\infty} [\bar{z}_f \sigma(x) + \mu(x) - f(x^+) - \xi] \phi(z_f) dz_f + 0$$

$$\begin{aligned}
&= \int_0^{+\infty} [u(x) - f(x^*) - \varepsilon] \phi(z_f) dz_f - 0 \cdot \int_0^{+\infty} z_f \cdot \phi(z_f) dz_f \\
&= [u(x) - f(x^*) - \varepsilon] \cdot \int_0^{+\infty} \phi(z_f) dz_f \\
&= [u(x) - f(x^*) - \varepsilon] \cdot \Phi(0) \\
&= 0
\end{aligned}$$

- b) After you run the code on `bayesopt_implemention_viewer.ipynb` for the higher dimensional case, comment on the benefit of including a hyperparameter (meaning the parameters of the covariance function θ of the surrogate model) optimisation approach on your Bayesian Optimisation algorithm. (2 Points)

After running the `bayesopt_implemention_viewer.ipynb`, we can observe that the surrogate model can simulate the original function better if we using the hyperparameter optimisation approach. As a result, the next point that it will select will have a higher chance that leads to the global optimum value. Therefore, we can use fewer iterations to do the optimisation compared with the Bayesian Optimisation without using the hyperparameters optimisation.

Another benefit is that it can save time as less iteration, less computation of time-consuming black block function f .

Name: Yixi Rao (NO GROUP)

UID: U682654

DATE: 07/05/2022