# Homework #1

CSE 546: Machine Learning

Ray Chen

Collaborator: Kevin Shao

October 20, 2022

**B1:**

- **Part a:** Proof
  For Trian error:

$$\mathbb{E}_{\text{train}}\big[\hat{\epsilon}_{\text{train}}(\hat{f}_{\text{train}})\big] = \mathbb{E}\big[\frac{1}{N_{train}} \sum_{(x,y)\in S_{trian}} (f(x)-y)^2\big]$$

$$= \frac{1}{N_{train}} \sum_{(x,y)\in S_{trian}} \mathbb{E}(f(x)-y)^2$$

$$= N_{train}\frac{1}{N_{train}} \sum_{(x,y)\sim D} \mathbb{E}(f(x)-y)^2$$

$$= \epsilon(f)$$

Similarly, for test error

$$\mathbb{E}_{\text{test}}\big[\hat{\epsilon}_{\text{test}}(\hat{f}_{\text{test}})\big] = \mathbb{E}\big[\frac{1}{N_{test}} \sum_{(x,y)\in S_{test}} (\hat{f}(x)-y)^2\big]$$

$$= \frac{1}{N_{test}} \sum_{(x,y)\in S_{test}} \mathbb{E}(\hat{f}(x)-y)^2$$

$$= N_{test}\frac{1}{N_{test}} \sum_{(x,y)\sim D} \mathbb{E}(\hat{f}(x)-y)^2$$

$$= \epsilon(\hat{f})$$

- **Part b:** Brief Explanation (3-5 sentences)
  No, it will not be true with regard to the training loss. Also, $\mathbb{E}_{\text{train}}\big[\hat{\epsilon}_{\text{train}}(\hat{f}_{\text{train}})\big] \neq \epsilon(\hat{f})$. Because $x_i$ and $y_i$ are not independent to the $\hat{f}$.

- **Part c:** Proof

$$\mathbb{E}_{\text{train}}\big[\hat{\epsilon}_{\text{train}}(\hat{f}_{\text{train}})\big] = \sum_{f\in\mathcal{F}} \mathbb{E}_{\text{train}}\big[\hat{\epsilon}_{\text{train}}(\hat{f}_{\text{train}})\big]\mathbb{P}_{\text{train}}(\hat{f}_{\text{train}} = f)$$

$$= \sum_{f\in\mathcal{F}} \mathbb{E}_{\text{test}}\big[\hat{\epsilon}_{\text{test}}(f)\big]\mathbb{P}_{\text{train}}(\hat{f}_{\text{train}} = f)$$

$$= \mathbb{E}_{\text{train,test}}\big[\hat{\epsilon}_{\text{test}}(\hat{f}_{\text{train}})\big]$$

**B2:**

- **Part a:** 1-2 sentences Intuitively, for a small m, it should have low bias and high variance; for a large m, it should have high bias and low variance. According to formula of bias-variance decomposition at $x_i$:

$$\mathbb{E}\left[(\widehat{f}_m(x_i) - f(x_i))^2\right] = \underbrace{(\mathbb{E}[\widehat{f}_m(x_i)] - f(x_i))^2}_{\text{Bias}^2(x_i)} + \underbrace{\mathbb{E}\left[(\widehat{f}_m(x_i) - \mathbb{E}[\widehat{f}_m(x_i)])^2\right]}_{\text{Variance}(x_i)}$$

- **Part b:** Proof

$$\frac{1}{n}\sum_{i=1}^{n}\left(\mathbb{E}[\hat{f}_m(x_i)] - f(x_i)\right)^2 = \frac{1}{n}\sum_{i=1}^{n}\left(\mathbb{E}\left[\sum_{j=1}^{n/m} c_j \mathbf{1}\right] - f(x_i)\right)^2$$

$$= \frac{1}{n}\sum_{i=1}^{n}\left(\mathbb{E}\left[\sum_{j=1}^{n/m}\left(\frac{1}{m}\sum_{k=(j-1)m+1}^{jm} y_k\right)\mathbf{1}\right] - f(x_i)\right)^2$$

$$= \frac{1}{n}\sum_{i=1}^{n}\left(\sum_{j=1}^{n/m}\left(\frac{1}{m}\sum_{k=(j-1)m+1}^{jm}\mathbb{E}\left[f(x_k) + \epsilon_k\right]\right)\mathbf{1} - f(x_i)\right)^2$$

$$= \frac{1}{n}\sum_{i=1}^{n}\left(\sum_{j=1}^{n/m}\bar{f}^{(j)}\mathbf{1} - f(x_i)\right)^2$$

$$= \frac{1}{n}\sum_{j=1}^{n/m}\sum_{i=(j-1)m+1}^{jm}\left(\bar{f}^{(j)} - f(x_i)\right)^2$$

- **Part c:** Proof

$$\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left(\hat{f}_m(x_i) - \mathbb{E}[\hat{f}_m(x_i)]\right)^2\right] = \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}\left[\left(\sum_{j=1}^{n/m}c_j\mathbf{1}\{x_i \in ((j-1)m, jm]\} - \mathbb{E}\left[\sum_{j=1}^{n/m}c_j\mathbf{1}\{x_i \in ((j-1)m, jm]\}\right]\right)^2\right]$$

$$= \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}\left[\left(\sum_{j=1}^{n/m}(c_j - \mathbb{E}[c_j])\mathbf{1}\{x_i \in ((j-1)m, jm]\}\right)^2\right]$$

$$= \frac{1}{n}\sum_{j=1}^{n/m}\mathbb{E}\left[\sum_{i=(j-1)m+1}^{jm}(c_j - \mathbb{E}[c_j])^2\right]$$

$$= \frac{1}{n}\sum_{j=1}^{n/m}m\mathbb{E}\left[(c_j - \bar{f}^{(j)})^2\right]$$

$$= \frac{1}{n}\sum_{j=1}^{n/m}m\mathbb{E}\left[\left(\frac{1}{m}\sum_{i=(j-1)m+1}^{jm}y_i - f(x_i)\right)^2\right]$$

$$= \frac{1}{n}\sum_{j=1}^{n/m}m\mathbb{E}\left[\left(\sum_{i=(j-1)m+1}^{jm}\frac{\epsilon_i}{m}\right)^2\right]$$

$$= \frac{1}{n}\sum_{j=1}^{n/m}\frac{1}{m}\sum_{i=(j-1)m+1}^{jm}\sigma^2$$

$$= \frac{\sigma^2}{m}$$

- **Part d:** Derivation of minimal error with respect to $m$. 1-2 sentences about scaling of $m$ with parameters.

$$\frac{1}{n}\sum_{j=1}^{n/m}\sum_{i=(j-1)m+1}^{jm}\left(\bar{f}^{(j)} - f(x_i)\right)^2 \leq \frac{L^2}{n^2}(argmin(x_i) - i)^2$$

$$\leq \frac{1}{n}\sum_{j=1}^{n/m}\sum_{i=(j-1)m+1}^{jm}\left(\frac{L}{n}m\right)^2$$

$$= \mathcal{O}\left(\frac{L^2 m^2}{n^2}\right)$$

Total error is $\mathcal{O}\left(\frac{L^2 m^2}{n^2} + \frac{\sigma^2}{m}\right)$. If we want to minimize with respect to m, we can set its derivative to zero:

$$\frac{d}{dm}\left(\frac{L^2 m^2}{n^2} + \frac{\sigma^2}{m}\right) = 0$$

$$\frac{2L^2 m}{n^2} - \frac{\sigma^2}{m^2} = 0$$

Then we have:

$$m = \left(\frac{\sigma^2 n^2}{2L^2}\right)^{1/3}$$

$$= \mathcal{O}\left(\left(\frac{\sigma n}{L}\right)^{2/3}\right) = \mathcal{O}\left(\left(\frac{n}{L}\right)^{2/3}\right)$$

Plug $m$ back in:

$$\frac{L^2\left(\frac{n^2\sigma^2}{2L^2}\right)^{2/3}}{n^2} + \sigma^2\left(\frac{n^2\sigma^2}{2L^2}\right)^{-1/3} = \frac{3}{4^{\frac{1}{3}}}\sigma^{\frac{4}{3}}\left(\frac{L}{n}\right)^{\frac{2}{3}}$$

$$= \mathcal{O}\left(\left(\frac{L}{n}\right)^{\frac{2}{3}}\right)$$

It just as intuition, minimized total error increased with the decrease of number of samples.