

## Exercise 1.3

(a)

The misclassification loss can be reformulated as:

$$R(f) = \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{y|\mathbf{x}} \left[ \mathbb{1}_{\{f(\mathbf{x}) \neq y\}} \right] \right].$$

Note that the inner expectation is:

$$\mathbb{E}_{y|\mathbf{x}} \left[ \mathbb{1}_{\{f(\mathbf{x}) \neq y\}} \right] = 1 - P(y = f(\mathbf{x})|\mathbf{x}).$$

Minimizing the misclassification loss is equivalent to maximizing  $P(y|\mathbf{x})$  for given  $\mathbf{x}$  w.r.t.  $y$  (MAP of  $y$ ):

$$\hat{y}_{\text{MAP}} = \arg \max_{y \in \{0,1\}} P(y|\mathbf{x}).$$

(b)

We have  $R(\mathbf{f}) = \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{y|\mathbf{x}} [\|\mathbf{y} - \mathbf{f}(\mathbf{x})\|_2^2] \right]$  for the Linear Regression problem. Our task is to find the  $\mathbf{f}(\mathbf{x})$  that minimizes  $R(\mathbf{f})$ .

Let  $\frac{\partial}{\partial \mathbf{f}} \mathbb{E}_{y|\mathbf{x}} [\|\mathbf{y} - \mathbf{f}(\mathbf{x})\|_2^2] = 0$ , we can solve for that optimal regression function:

$$\mathbf{f}^*(\mathbf{x}) = \mathbb{E}[\mathbf{y}|\mathbf{x}].$$

## Exercise 1.4

(a)

Substitute  $X_i$  by  $\mathcal{L}(\mathbf{y}_i, \mathbf{f}(\mathbf{x}_i))$ , note that  $R_{\text{emp}}(\mathbf{f} | \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N X_i$  and  $R(\mathbf{f}) = \mathbb{E}[X_i]$ .

According to Hoeffding's inequality:

$$\Pr \left( \left| \frac{1}{N} \sum_{i=1}^N X_i - \mathbb{E}[X_i] \right| \geq \varepsilon \right) \leq 2 \exp \left( -\frac{2N\varepsilon^2}{M^2} \right).$$

Thus,

$$\Pr \left( \left| \frac{1}{N} \sum_{i=1}^N X_i - \mathbb{E}[X_i] \right| \leq \varepsilon \right) \geq 1 - 2 \exp \left( -\frac{2N\varepsilon^2}{M^2} \right).$$

Let  $\varepsilon = \sqrt{\frac{M^2 \ln(2/\delta)}{2N}}$ , we can show that:

$$\Pr \left( \left| \frac{1}{N} \sum_{i=1}^N X_i - \mathbb{E}[X_i] \right| \leq \sqrt{\frac{M^2 \ln(2/\delta)}{2N}} \right) \geq 1 - \delta.$$

(b)

For finite number of bounded training samples, it is always possible for ERM to approximate the correct expectation value, if the size of the training set is big enough. However, the decreasing speed of the absolute value maybe relatively slow (of order  $O(1/\sqrt{n})$ ).