

EXERCISE 1 (BAYES CLASSIFIER) Prove that the optimal binary Bayes Classifier for the 0 – 1 loss is given by

$$f^*(X) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X) \geq \mathbb{P}(Y = -1|X) \\ -1 & \text{otherwise} \end{cases}$$

Solution.

$$\begin{aligned} \mathbb{E}[\ell^{0-1}(Y, f(X))] &= \mathbb{E}[\ell^{0-1}(Y, f(X))|Y = 1, X] \mathbb{P}(Y = 1|X) + \mathbb{E}[\ell^{0-1}(Y, f(X))|Y = -1, X] \mathbb{P}(Y = -1|X) \\ &= \mathbf{1}_{f(X)=1} \mathbb{P}(Y = 1|X) + \mathbf{1}_{f(X)=-1} \mathbb{P}(Y = -1|X) \end{aligned}$$

which is indeed minimized by f^* . □

EXERCISE 2 (TRAINING ERROR OPTIMISM) Assume that $((x_i, y_i))_i$ is a sequence of i.i.d. random vectors. Consider a parametric model f_β fit by least squares to a set of training data $(x_1, y_1), \dots, (x_N, y_N)$ drawn at random from a population. Let $\hat{\beta} \in \mathbb{R}^d$ be the least squares estimate, i.e.

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N (y_i - f_\beta(x_i))^2$$

Suppose we have some test data $(\tilde{x}_1, \tilde{y}_1), \dots, (\tilde{x}_M, \tilde{y}_M)$ drawn at random from the same population $((x_i, y_i))_i$ as the training data.

1. If we define

$$R_{train}(\beta) = \frac{1}{N} \sum_{i=1}^N (y_i - f_\beta(x_i))^2$$

and

$$R_{test}(\beta) = \frac{1}{M} \sum_{i=1}^M (\tilde{y}_i - f_\beta(\tilde{x}_i))^2$$

Prove that

$$\mathbb{E}[R_{train}(\hat{\beta})] \leq \mathbb{E}[R_{test}(\hat{\beta})]$$

where the expectations are over all that is random in each expression.

Solution.

Let $\beta^* = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \mathbb{E}[R_{train}(\beta)] = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \mathbb{E}[(y_0 - f_\beta(x_0))^2]$. Note that β^* is a vector of \mathbb{R}^d that is not random.

By construction, $R_{train}(\hat{\beta}) \leq R_{train}(\beta^*)$ and thus

$$\mathbb{E}[R_{train}(\hat{\beta})] \leq \mathbb{E}[R_{train}(\beta^*)]. \quad (1)$$

Note that there is no contradiction with $\beta^* = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \mathbb{E}[R_{train}(\beta)]$ because $\hat{\beta}$ is a random vector that depends on and could be different for each $(x_1, y_1), \dots, (x_N, y_N)$.

Denote $g(\beta) = \mathbb{E}[R_{test}(\beta)]$ for $\beta \in \mathbb{R}^d$. Since we also have $\mathbb{E}[R_{test}(\beta)] = \mathbb{E}[(y_0 - f_\beta(x_0))^2] = \mathbb{E}[R_{train}(\beta)]$ for any $\beta \in \mathbb{R}^d$, we deduce that $\beta^* = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \mathbb{E}[R_{test}(\beta)] = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} g(\beta)$. But $\hat{\beta} \in \sigma((x_1, y_1), \dots, (x_N, y_N))$ and is independant from $(\tilde{x}_1, \tilde{y}_1), \dots, (\tilde{x}_M, \tilde{y}_M)$.

Therefore, for any $\beta \in \mathbb{R}^d$, we also have $g(\beta) = \mathbb{E}[R_{test}(\beta) | \sigma((x_1, y_1), \dots, (x_N, y_N))]$. Hence $g(\beta^*) \leq \mathbb{E}[R_{test}(\beta) | \sigma((x_1, y_1), \dots, (x_N, y_N))]$ for any $\beta \in \mathbb{R}^d$.

Since $g(\hat{\beta}) = \mathbb{E}[R_{test}(\hat{\beta}) | \sigma((x_1, y_1), \dots, (x_N, y_N))]$, we have $g(\beta^*) \leq g(\hat{\beta})$. Finally this induces $g(\beta^*) = \mathbb{E}[g(\beta^*)] \leq \mathbb{E}[g(\hat{\beta})]$ and this leads to

$$\mathbb{E}[R_{train}(\beta^*)] = g(\beta^*) = \mathbb{E}[R_{test}(\beta^*)] \leq \mathbb{E}[R_{test}(\hat{\beta})] = \mathbb{E}[g(\hat{\beta})]. \quad (2)$$

As a consequence, we obtain

$$\mathbb{E}[R_{train}(\hat{\beta})] \leq \mathbb{E}[R_{train}(\beta^*)] = \mathbb{E}[R_{test}(\beta^*)] \leq \mathbb{E}[R_{test}(\hat{\beta})].$$

□

2. Can we replace $R_{test}(\beta)$ by the risk $R(\beta) = \mathbb{E}[(\tilde{y} - f_\beta(\tilde{x}))^2]$ where (\tilde{x}, \tilde{y}) follows the population law?

Solution.

From the assumption, we establish for any $\beta \in \mathbb{R}^d$, $g(\beta) = R(\beta) = \mathbb{E}[R_{test}(\beta)] = \mathbb{E}[(\tilde{y} - f_\beta(\tilde{x}))^2] = \mathbb{E}[(y_1 - f_\beta(x_1))^2] = \mathbb{E}[R_{train}(\beta)]$. \square

3. Let β^* be the minimizer of the risk $R(\beta)$, prove that

$$\mathbb{E}[R_{train}(\hat{\beta})] \leq R(\beta^*) \leq \mathbb{E}[R(\hat{\beta})]$$

Solution.

From (2), $g(\beta^*) = R(\beta^*) \leq \mathbb{E}[g(\hat{\beta})]$ and (1), $\mathbb{E}[R_{train}(\hat{\beta})] \leq R(\beta^*) = g(\beta^*)$, and this achieves the proof. \square

EXERCISE 3 (TESTING ERROR) Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space. Assume that (X, Y) is a couple of random variables defined on $(\Omega, \mathcal{F}, \mathbb{P})$ and taking values in $\mathcal{X} \times \{-1, 1\}$ where \mathcal{X} is a given state space. One aim of supervised classification is to define a function $h : \mathcal{X} \rightarrow \{-1, 1\}$, called *classifier*, such that $h(X)$ is the best prediction of Y in a given context. For instance, the probability of misclassification of h is

$$L_{\text{miss}}(h) = \mathbb{P}(Y \neq h(X)) .$$

Note that $\mathbb{E}[Y|X]$ is a random variable measurable with respect to the σ -algebra $\sigma(X)$. Therefore, there exists a function $\eta : \mathcal{X} \rightarrow [-1, 1]$ so that $\mathbb{E}[Y|X] = \eta(X)$ almost surely.

In Exercise 1, we have shown that h_* , defined for all $x \in \mathcal{X}$, by

$$h_*(x) = \begin{cases} 1 & \text{if } \eta(x) > 0, \\ -1 & \text{otherwise,} \end{cases}$$

is such that

$$h_* = \underset{h: \mathcal{X} \rightarrow \{-1, 1\}}{\operatorname{argmin}} L_{\text{miss}}(h) .$$

1. In practice, the minimization of L_{miss} holds on a specific set \mathcal{H} of classifiers (often called the *dictionary*), which may possibly not contain the Bayes classifier. Moreover, since in most cases, the classification risk L_{miss} cannot be computed nor minimized, it is instead estimated by the empirical classification risk defined as

$$\hat{L}_{\text{miss}}^n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{Y_i \neq h(X_i)} ,$$

where $(X_i, Y_i)_{1 \leq i \leq n}$ are independent observations with the same distribution as (X, Y) . The classification problem then boils down to solving

$$\hat{h}_{\mathcal{H}}^n \in \underset{h \in \mathcal{H}}{\operatorname{argmin}} \hat{L}_{\text{miss}}^n(h) .$$

Prove that for all set \mathcal{H} of classifiers and all $n \geq 1$,

$$L_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) - \inf_{h \in \mathcal{H}} L_{\text{miss}}(h) \leq 2 \sup_{h \in \mathcal{H}} \left| \hat{L}_{\text{miss}}^n(h) - L_{\text{miss}}(h) \right| .$$

Solution.

By definition of $\hat{h}_{\mathcal{H}}^n$, for any $h \in \mathcal{H}$,

$$\begin{aligned} L_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) - \inf_{h \in \mathcal{H}} L_{\text{miss}}(h) &= L_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) - \hat{L}_{\text{miss}}^n(\hat{h}_{\mathcal{H}}^n) + \hat{L}_{\text{miss}}^n(\hat{h}_{\mathcal{H}}^n) - \inf_{h \in \mathcal{H}} L_{\text{miss}}(h) , \\ &\leq L_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) - \hat{L}_{\text{miss}}^n(\hat{h}_{\mathcal{H}}^n) + \hat{L}_{\text{miss}}^n(h) - \inf_{h \in \mathcal{H}} L_{\text{miss}}(h) . \end{aligned}$$

For all $\varepsilon > 0$ there exists $h_\varepsilon \in \mathcal{H}$ such that $L_{\text{miss}}(h_\varepsilon) < \inf_{h \in \mathcal{H}} L_{\text{miss}}(h) + \varepsilon$ so that

$$\begin{aligned} L_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) - \inf_{h \in \mathcal{H}} L_{\text{miss}}(h) &\leq L_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) - \hat{L}_{\text{miss}}^n(\hat{h}_{\mathcal{H}}^n) + \hat{L}_{\text{miss}}^n(h_\varepsilon) - L_{\text{miss}}(h_\varepsilon) + \varepsilon , \\ &\leq 2 \sup_{h \in \mathcal{H}} \left| \hat{L}_{\text{miss}}^n(h) - L_{\text{miss}}(h) \right| + \varepsilon , \end{aligned}$$

which concludes the proof. \square

2. Using Hoeffding's inequality, stating that if $X_i \in [a_i, b_i]$ then for all $t > 0$,

$$\mathbb{P}\left(\left|\sum_{i=1}^n X_i - \sum_{i=1}^n \mathbb{E}[X_i]\right| > t\right) \leq 2\exp\left(\frac{-2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right), \quad (3)$$

prove that when $\mathcal{H} = \{h_1, \dots, h_M\}$ for a given $M \geq 1$, then, for all $\delta > 0$,

$$\mathbb{P}\left(L_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) \leq \min_{1 \leq j \leq M} L_{\text{miss}}(h_j) + \sqrt{\frac{2}{n} \log\left(\frac{2M}{\delta}\right)}\right) \geq 1 - \delta.$$

Solution.

By the previous question, for all $u > 0$,

$$\mathbb{P}\left(L_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) > \min_{1 \leq j \leq M} L_{\text{miss}}(h_j) + u\right) \leq \mathbb{P}\left(\sup_{h \in \mathcal{H}} |\hat{L}_{\text{miss}}^n(h) - L_{\text{miss}}(h)| > \frac{u}{2}\right) \leq \sum_{j=1}^M \mathbb{P}\left(|\hat{L}_{\text{miss}}^n(h_j) - L_{\text{miss}}(h_j)| > \frac{u}{2}\right).$$

By Hoeffding's inequality,

$$\mathbb{P}\left(L_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) > \min_{1 \leq j \leq M} L_{\text{miss}}(h_j) + u\right) \leq 2Me^{-nu^2/2},$$

which concludes the proof by choosing

$$u = \sqrt{\frac{2}{n} \log\left(\frac{2M}{\delta}\right)}.$$

□