
INTRODUCTION TO MACHINE LEARNING

1 Bayes classifier

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, which means that the focus is set on a two-class classification problem. 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 risk of misclassification of h is

$$R_{\text{miss}}(h) = \mathbb{E} [\mathbb{1}_{Y \neq h(X)}] = \mathbb{P}(Y \neq h(X)) .$$

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.

Lemma 1. *The classifier 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

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

Proof. For all $u, v \in \{-1, 1\}$, $\mathbb{1}_{u \neq v} = \mathbb{1}_{uv = -1} = (1 - uv)/2$. Since Y and $h(X)$ take values in $\{-1, 1\}$, this implies

$$R_{\text{miss}}(h) = \mathbb{P}(Y \neq h(X)) = (1 - \mathbb{E}[Yh(X)]) / 2 .$$

Now, using successively the tower property, the equality $|u| = u \times \operatorname{sign}(u)$, and the tower property again,

$$\mathbb{E}[Yh(X)] = \mathbb{E}[\mathbb{E}[Y|X]h(X)] \leq \mathbb{E}[|\mathbb{E}[Y|X]| \underbrace{|h(X)|}_{=1}] = \mathbb{E}[\mathbb{E}[Y|X] \underbrace{\operatorname{sign}(\mathbb{E}[Y|X])}_{h_*(X)}] = \mathbb{E}[Yh_*(X)] .$$

This yields $R_{\text{miss}}(h) \geq R_{\text{miss}}(h_*)$, which concludes the proof. \square

2 Empirical risk minimization

First, we do not assume that the joint law of (X, Y) belongs to any parametric or semiparametric family of models. Instead, we make some restrictions on the set of classifiers on which the optimisation occurs.

More precisely, we consider that the optimization of classifiers 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 R_{miss} cannot be computed nor minimized, it is instead estimated by the empirical classification risk defined as

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

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

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

In this context several practical and theoretical challenges arise from the minimization of the empirical classification risk. The choice of \mathcal{H} is pivotal in designing an efficient classification procedure. Note that choosing \mathcal{H} as all possible classifiers is meaningless, in this case, $\hat{h}_{\mathcal{H}}^n$ is such that $\hat{h}_{\mathcal{H}}^n(X_i) = Y_i$ for all $1 \leq i \leq n$ and $\hat{h}_{\mathcal{H}}^n(x)$ is any element of $\{-1, 1\}$ for all $x \notin \{X_1, \dots, X_n\}$. Although $\hat{R}_{\text{miss}}^n(h_{\mathcal{H}}^n) = 0$, is likely to be a poor approximation of $R_{\text{miss}}(h_{\mathcal{H}}^n)$. To understand this, the excess misclassification risk may be decomposed as follows

$$R_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) - R_{\text{miss}}(h_{\star}) = R_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) - \min_{h \in \mathcal{H}} R_{\text{miss}}(h) + \min_{h \in \mathcal{H}} R_{\text{miss}}(h) - R_{\text{miss}}(h_{\star}) \geq 0.$$

The first term of the decomposition $R_{\text{miss}}(\hat{h}_{\mathcal{H}}^n) - \min_{h \in \mathcal{H}} R_{\text{miss}}(h)$ is a **stochastic error** which is likely to grow when the size of \mathcal{H} grows while $\min_{h \in \mathcal{H}} R_{\text{miss}}(h) - R_{\text{miss}}(h_{\star})$ is **deterministic** and likely to decrease as the size of \mathcal{H} grows.

Lemma 2. For all set \mathcal{H} of classifiers and all $n \geq 1$,

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

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

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

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

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

which concludes the proof. \square