

Chapter 17

Bias-variance tradeoff I

17.1 Hypothesis space \mathcal{H}

In this book we choose the hypothesis space \mathcal{H} as a compact subset of the space of continuous functions $\mathbb{C}(X)$ equipped with the standard norm

$$\|f\|_{\mathbb{C}(X)} := \|f\|_{\infty} := \sup_X |f(\mathbf{x})|, \quad \forall f \in \mathbb{C}(X).$$

While this choice does not cover all the practical problems of interest, it provides rich mathematical structures to study machine learning methods.

17.2 Empirical target function

Recall from Chapter 16 that the best target function is the regression function h^* which is, unfortunately, not computable. In general, h^* does not reside in \mathcal{H} , and thus the best that we hope for is to find a *target function* \hat{h} in \mathcal{H} that is closest, e.g. in the least squares sense, to h^* :

$$\hat{h} := \arg \min_{\mathcal{H}} \int_X (f(\mathbf{x}) - h^*(\mathbf{x}))^2 d\pi(\mathbf{x}) \quad (17.1)$$

Exercise 17.1. Show that \hat{h} is also closest to \mathbf{y} , that is, \hat{h} is a minimizer of the following least squares problem

$$\min_{\mathcal{H}} \int_Z (f(\mathbf{x}) - y)^2 d\pi(\mathbf{x}, \mathbf{y})$$

•

Exercise 17.1 implies

$$\hat{h} := \arg \min_{\mathcal{H}} \mathcal{R}(f). \quad (17.2)$$

Compared to (17.1), for which we have no access to h^* , we partially know y in (17.2) through the training set $S := \{(\mathbf{x}^i, y^i)\}_{i=1}^N$. Therefore, instead of seeking \hat{h} , which is not possible, we resort to a minimizer, the *empirical target function*, \hat{h}_N of the empirical risk minimization problem defined as

$$\hat{h}_N := \arg \min_{\mathcal{H}} \mathcal{R}_N(f) = \frac{1}{N} \sum_i^N (f(\mathbf{x}^i) - y^i)^2. \quad (17.3)$$

The question that needs to be addressed is if the target and the empirical target functions exist. To that end, by Theorem 17.1, we need to show that the risk function and the empirical risk function are continuous.

Assumption 17.1 (*M-Boundedness of the misfit on \mathcal{H}*). For any $h \in \mathcal{H}$ and almost everywhere (a.e. or aka a.s.) in X , there holds

$$|h(\mathbf{x}) - y| \leq M. \quad (17.4)$$

Proposition 17.1. *Suppose \mathcal{H} is bounded in the sense of Assumption 17.1. Then $\mathcal{R}, \mathcal{R}_N : \mathcal{H} \rightarrow \mathbb{R}$ are Lipschitz continuous.*

Proof. The proof is straightforward using (17.4). Indeed, we have

$$|\mathcal{R}(h_1) - \mathcal{R}(h_2)| = \left| \int_Z (h_1 + h_2 - 2y)(h_1 - h_2) d\pi(\mathbf{x}, \mathbf{y}) \right| \leq 2M \|h_1 - h_2\|_\infty.$$

Similarly,

$$\begin{aligned} |\mathcal{R}_N(h_1) - \mathcal{R}_N(h_2)| &= \left| \frac{1}{N} \sum_i^N (h_1(\mathbf{x}^i) + h_2(\mathbf{x}^i) - 2y^i)(h_1(\mathbf{x}^i) - h_2(\mathbf{x}^i)) \right| \\ &\leq \frac{1}{N} \sum_i^N |h_1(\mathbf{x}^i) + h_2(\mathbf{x}^i) - 2y^i| \|h_1 - h_2\|_\infty \leq 2M \|h_1 - h_2\|_\infty. \end{aligned}$$

□

Let us define the *sampling error* in computing the risk for any h as

$$\mathcal{E}(h) := \mathcal{R}(h) - \mathcal{R}_N(h). \quad (17.5)$$

As will be shown in the next section this sampling error plays the key role in estimating the sample/estimation error.

Exercise 17.2 (Continuity of Sampling Error). Show that $\mathcal{E} : \mathcal{H} \rightarrow \mathbb{R}$ is Lipschitz continuous. •

Corollary 17.1 (Existence of \hat{h} and \hat{h}_N). *Suppose \mathcal{H} is bounded in the sense of Assumption 17.1. Then \hat{h} and \hat{h}_N exist.*

Proof. The assertion is obvious due to Proposition 17.1 and Theorem 17.1. □

17.3 Bias-Variance Tradeoff

Suppose the empirical target \hat{h}_N is already computed/estimated (to be discussed in details later), the question of interest is how to estimate the actual risk $\mathcal{R}(\hat{h}_N)$. This section decomposes this actual risk into two parts: the *sample or estimation* error and the approximation error. The compromise between these two errors is known as the *bias-variance tradeoff*. We begin with the following decomposition

$$\mathcal{R}(\hat{h}_N) = \underbrace{\mathcal{R}(\hat{h}_N) - \mathcal{R}(\hat{h})}_{\mathcal{S}(\hat{h}_N):=} + \underbrace{\mathcal{R}(\hat{h})}_{\mathcal{B}(\hat{h}):=}.$$

Let us first consider $\mathcal{S}(\hat{h}_N)$. We have

$$\begin{aligned} \mathcal{S}(\hat{h}_N) &= \underbrace{\mathcal{R}_N(\hat{h}_N) - \mathcal{R}_N(\hat{h})}_{\leq 0} + \mathcal{R}(\hat{h}_N) - \mathcal{R}_N(\hat{h}_N) + \mathcal{R}_N(\hat{h}) - \mathcal{R}(\hat{h}) \\ &\leq |\mathcal{R}(\hat{h}_N) - \mathcal{R}_N(\hat{h}_N)| + |\mathcal{R}_N(\hat{h}) - \mathcal{R}(\hat{h})|, \end{aligned} \quad (17.6)$$

where the negativeness of the first difference is due to the fact that \hat{h}_N is a minimizer of \mathcal{R}_N . Note that the last two terms are the sampling errors incurred by approximating $\mathcal{R}(\hat{h}_N)$ and $\mathcal{R}(\hat{h})$ using the training set S . Since $\mathcal{S}(\hat{h}_N)$ is bounded by the sampling errors, it is conventionally named as the *sample/estimation* error or the *variance*.

Now due to Exercise 16.1 we can rewrite $\mathcal{B}(\hat{h})$ as

$$\mathcal{B}(\hat{h}) = \int_X (\hat{h}(\mathbf{x}) - h^*(\mathbf{x}))^2 d\pi(\mathbf{x}) + \sigma_\pi^2, \quad (17.7)$$

which is independent of the training set. It actually depends only on the distance between \hat{h} and the regression function h^* . That is, it is completely determined by the approximation capability of the hypothesis space \mathcal{H} . Thus, $\mathcal{B}(\hat{h})$ is often called the *bias*.

We are in the position to discuss *bias-variance tradeoff*. For a fixed \mathcal{H} , the sample error clearly decreases as the sample size N increases. Now if we fix the sample size N , enlarging the hypothesis space \mathcal{H} clearly reduces the bias but increase the sample error in general (will be shown in the next chapters). A popular tradeoff is to increase the dimension of the hypothesis space \mathcal{H} as the sample size N increases. The question is how fast we should enlarge \mathcal{H} . To answer this question in a sensible and rigorous way, we need to first estimate both sample error and the bias, and then balance out these errors. This is the task of the following chapters.

17.4 Appendix

Theorem 17.1. *Let K be a compact set and $L : K \rightarrow \mathbb{R}$ be a continuous function. Then L is bounded. In particular, there exist $f, g \in K$ at which L attains its minimum and maximum.*