# Chapter 16
# Statistical machine learning

## 16.1 What is machine learning in this book?

Given a set of data $S := \left\{ \left( \mathbf{x}^i, y^i \right) \right\}_{i=1}^N$, where $y^i = f^* \left( \mathbf{x}^i \right)$ for some unknown function $f^*$, *in this book machine learning is the task of learning the function (or the map) $f^*$ from the incomplete information described by the training set* $\left\{ \left( \mathbf{x}^i, y^i \right) \right\}_{i=1}^N$. Note that we have assume that $y(\mathbf{x})$ is a scalar-valued function, and the extension to vector-valued function is straightforward.

Let $X$ be a closed and bounded (and hence compact) subset of $\mathbb{R}^k$, and $y \in \mathbb{R}$. Let $h$ be any approximation of $f^*$ and we are interested in measuring the error that we commit in approximating $f^*$ using $h$, the error is also known as the "risk" in the machine learning literature. To measure the risk we let $\pi$ be a (Borel) probability measure on the on product space $Z := X \times y$, $\pi(y|\mathbf{x})$ be the conditional probability measure of $y$ given $\mathbf{x}$, and $\pi(\mathbf{x}) := \int_y d\pi(\mathbf{x}, y)$ be the marginal probability measure on $X$. By Bayes' theorem, we know that $\pi(\mathbf{x}, y) = \pi(y|\mathbf{x}) \times \pi(\mathbf{x})$. One of the popular risk functions is the *least squares error* with respect to the product measure $\pi$:

$$\mathscr{R}(h) := \int_Z \left( h(\mathbf{x}) - y \right)^2 d\pi(\mathbf{x}, y).$$

Cleary, the best $h^*(\mathbf{x}) \in \mathbb{L}^2(X, \pi(\mathbf{x}))$, i.e. $\int_X \|h^*\|_{\mathbb{R}^d}^2 \, d\pi(\mathbf{x}) < \infty$, is the one which minimizes the risk. Chapter 2 shows that the first variation of the risk $\mathscr{R}$ at $h^*$ in any direction $\tilde{h}(\mathbf{x}) \in \mathbb{L}^2(X, \pi(\mathbf{x}))$ must vanish, i.e.,

$$\mathscr{D}\mathscr{R}\left( h^*; \tilde{h} \right) = 2 \int_Z \left( h^*(\mathbf{x}) - y \right) \tilde{h}(\mathbf{x}) \, d\pi(\mathbf{x}, y) = 0, \quad \forall \tilde{h}(\mathbf{x}) \in \mathbb{L}^2(X, \pi(\mathbf{x})),$$

which, after using the Bayes' theorem, becomes

$$\int_X \left( h^*(\mathbf{x}) - \int_y y \, d\pi(y|\mathbf{x}) \right) \tilde{h}(\mathbf{x}) \, d\pi(\mathbf{x}) = 0, \quad \forall \tilde{h}(\mathbf{x}) \in \mathbb{L}^2(X, \pi(\mathbf{x})),$$

which, by Riesz representation theorem, implies that

$$h^*(\mathbf{x}) = \int_y y\, d\pi(y|\mathbf{x}).$$

(16.1)

Note that $h^*(\mathbf{x})$ is known as the *regresstion function of* $\pi(\mathbf{x}, y)$.

**Exercise 16.1.** Show that the risk for any function $h$ is given as

$$\mathscr{R}(h) = \int_X (h(\mathbf{x}) - h^*(\mathbf{x}))^2\, d\pi(\mathbf{x}) + \sigma_\pi^2,$$

where we have defined

$$\sigma_\pi^2 := \int_Z (h^*(\mathbf{x}) - y)^2\, d\pi(\mathbf{x}, y)$$

•

Note that since $\sigma_\pi^2$ is independent of $h$, minimizing the risk is the same as minimizing distance between $h$ and $h^*$. In particular, $h^*$ is the minimizer of the risk.

## 16.2 Empirical Risk Minimization (ERM)

**Assumption 16.1 (i.i.d. assumption of the training set $S$).** The training set $S = \left\{(\mathbf{x}^i, y^i)\right\}_{i=1}^N$ are i.i.d. draws from $\pi(\mathbf{x}, y)$ on the product space $Z = X \times y$. Note that $S = \left\{(\mathbf{x}^i, y^i)\right\}_{i=1}^N$ can be equivalently considered as a random draw from on $Z^N$ with the $N$-fold product measure $\pi(\mathbf{x}, y) \times \pi(\mathbf{x}, y) \times \ldots \times \pi(\mathbf{x}, y)$.

Unfortunately the regression function (16.1) is not computable, and hence the risk, since it would require the joint distribution $\pi(\mathbf{x}, y)$, which is unknown. What available about $\pi(\mathbf{x}, y)$ is the limited information given by the training set $S$. *We assuming that the training set is i.i.d. in the sense of Assumption 16.1*, then by the law of large numbers (see Chapter **??**) we have

$$\mathscr{R}_N(h) := \frac{1}{N} \sum_{i=1}^N \left(h(\mathbf{x}^i) - y^i\right)^2 \overset{a.s.}{\to} \mathscr{R}(h),$$

where $\mathscr{R}_N(h)$ *is known as the empirical risk*. Thus we have to resort to minimizing the empirical risk, i.e.,

$$h_N(\mathbf{x}) \in \arg\min_h \mathscr{R}_N(h).$$

(16.2)

This is the well-known *empirical risk minimization* (ERM) problem, and here $h_N$ is a solution to the ERM problem.

Note that since the training set contain random realizations of $\pi(\mathbf{x}, y)$, $h_N(\mathbf{x})$ is a random function.

## 16.3 Overfitting and No-Free-Lunch theorem

It turns out that there we can always find a *learner*, i.e. a learning algorithm, $h_N(\mathbf{x})$ that makes that empirical risk vanish, i.e. $\mathscr{R}_N(h_N) = 0$ while having (significant) actual risk $\mathscr{R}(h_N) > 0$. This is known as overfitting: the learner does not *generalize* beyond the data, that is, the training error vanishes but the generalized error can be significant. Since the actual risk is our primary object of interest, overfitting is undersirable. You may hope to find a learner that avoid overfitting. Unfortunately this is impossible as pointed out by the following No-Free-Lunch theorem.

**Theorem 16.1 (No-Free-Lunch).** *For any $N \geq 1$, and any learner $h_N(\mathbf{x})$ built from the training data $S = \left\{ \left( \mathbf{x}^i, y^i \right) \right\}_{i=1}^{N}$, where $y \in \{0, 1\}$, there exists a joint distribution $\pi(\mathbf{x}, y)$ such that*

- $\mathscr{R}_N(h_N) = 0$ *and*
- $\mathbb{E}_{\pi}[\mathscr{R}(h_N)] \geq 1/8$.

*Proof.* See [26]. □

In other words, there is no universal learner—no learner can succeed in all learning tasks—unless further hypothesis or prior structure is present at the beginning of the learning process. This is, in the language of inverse problem, the ill-posedness nature of the learning problem in which we need to make inference on the high (possibly infinite) dimensional learner with limited information from data. We have seen in Chapter 1 that this problem can be overcome by using regularization, and restricting the learner having a prior structure is in fact a regularization technique.

In this book, the prior structure takes the form of a class of functions, namely, the *hypothesis space* in which we find a learner $h_N(\mathbf{x})$ that approximates the regression function (16.1) as well as it can while avoiding the overfitting problem. We shall also determine the compromise between the sample size $N$ and the hypothesis space in order to minimize the *generalized error*, i.e. the actual risk.