

Statistical Learning

In order to analyze a learning algorithm, we must define a mathematical model of how examples (\mathbf{x}, y) are generated. In the statistical learning framework, we assume that every example (\mathbf{x}, y) is obtained through an independent draw from a fixed but unknown probability distribution on $\mathcal{X} \times \mathcal{Y}$. We often write (\mathbf{X}, Y) to highlight that \mathbf{x} and y are random variables. The assumption that not all data points \mathbf{x} are equally likely is quite natural (for example, when data points are images, only a small fraction of all possible pixel configurations correspond to real-world images). Similarly, as we previously argued labels are typically noisy, and thus a probability distribution is a good mathematical abstraction to explain the variety of labels which a same data point may be associated with.

As we just explained, in statistical learning every example (\mathbf{x}, y) is the realization of an independent random draw from the same joint probability distribution \mathcal{D} . This implies that every dataset (e.g., a training set or a test set) is a **random sample**, of the kind we study in statistics. Note that the independence assumption is actually violated in many practical domains. Consider for example the problem of categorizing news stories. The newsfeed is clearly far from being an independent process, as the evolution of news reflects developing and related stories. Although not very realistic, the independence assumption is nevertheless convenient from the viewpoint of the analytical tractability of the problem.

In statistical learning, a problem is fully specified by a pair (\mathcal{D}, ℓ) , where \mathcal{D} is the data distribution and ℓ is a loss function. The performance of a predictor $h : \mathcal{X} \rightarrow \mathcal{Y}$ with respect to (\mathcal{D}, ℓ) is evaluated via the **statistical risk**, defined by

$$\ell_{\mathcal{D}}(h) = \mathbb{E}[\ell(Y, h(\mathbf{X}))]$$

This is the expected value of the loss function on a random example (\mathbf{X}, Y) drawn from \mathcal{D} . The best possible predictor $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ given \mathcal{D} is known as **Bayes optimal predictor**, and is defined by

$$f^*(\mathbf{x}) = \operatorname{argmin}_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(Y, \hat{y}) \mid \mathbf{X} = \mathbf{x}]$$

Note that $f^*(x)$ is the prediction minimizing the conditional risk, which is the expected loss of the prediction with respect to the distribution of the label Y conditioned on \mathbf{x} . By definition of f^* , we have that

$$\mathbb{E}[\ell(Y, f^*(x)) \mid \mathbf{X} = \mathbf{x}] \leq \mathbb{E}[\ell(Y, h(x)) \mid \mathbf{X} = \mathbf{x}]$$

for every predictor $h : \mathcal{X} \rightarrow \mathcal{Y}$. Because the above inequality holds for every $\mathbf{x} \in \mathcal{X}$, it also holds in expectation with respect to the random draw of \mathbf{X} . But since the average of the conditional risk is the risk, we have that $\ell_{\mathcal{D}}(f^*) \leq \ell_{\mathcal{D}}(h)$ for every predictor h . The risk $\ell_{\mathcal{D}}(f^*)$ of the Bayes optimal predictor is called **Bayes error**. Typically, the Bayes error is larger than zero because labels are stochastic.

We now compute the Bayes optimal predictor for the quadratic loss function $\ell(y, \hat{y}) = (y - \hat{y})^2$ when $\mathcal{Y} \equiv \mathbb{R}$,

$$\begin{aligned}
f^*(\mathbf{x}) &= \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \mathbb{E}[(Y - \hat{y})^2 \mid \mathbf{X} = \mathbf{x}] \\
&= \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \left(\mathbb{E}[Y^2 \mid \mathbf{X} = \mathbf{x}] + \hat{y}^2 - 2\hat{y} \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] \right) \\
&= \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \left(\hat{y}^2 - 2\hat{y} \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] \right) \quad (\text{ignoring the term that does not depend on } \hat{y}) \\
&= \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] \quad (\text{minimizing the function } F(\hat{y}) = \hat{y}^2 - 2\hat{y} \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}])
\end{aligned}$$

Thus, the Bayes optimal prediction for the quadratic loss function is the expected value of the label conditioned on the instance.

Substituting in the conditional risk formula $\mathbb{E}[(Y - f^*(\mathbf{X}))^2 \mid \mathbf{X} = \mathbf{x}]$ the Bayes optimal predictor $f^*(\mathbf{x}) = \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}]$ we obtain

$$\mathbb{E}[(Y - f^*(\mathbf{X}))^2 \mid \mathbf{X} = \mathbf{x}] = \mathbb{E}[(Y - \mathbb{E}[Y \mid \mathbf{x}])^2 \mid \mathbf{X} = \mathbf{x}] = \operatorname{Var}[Y \mid \mathbf{X} = \mathbf{x}].$$

In words, the conditional risk of the Bayes optimal predictor for the quadratic loss is the variance of the label conditioned on the instance.

We now focus on binary classification, where $\mathcal{Y} = \{-1, +1\}$. In this case we denote the joint distribution of (\mathbf{X}, Y) with the pair (D_X, η) , where D_X is the marginal of \mathcal{D} on \mathcal{X} and η is the distribution of $Y \in \{-1, +1\}$ conditioned on \mathcal{X} . We view the distribution η as a function $\eta : \mathcal{X} \rightarrow [0, 1]$, where $\eta(\mathbf{x}) = \mathbb{P}(Y = +1 \mid \mathbf{X} = \mathbf{x})$ is the probability that \mathbf{x} occurs with label +1.

Let $\mathbb{I}\{A\} \in \{0, 1\}$ be the indicator function of an event A ; that is, $\mathbb{I}\{A\} = 1$ if and only if A occurs. The statistical risk with respect to the zero-one loss $\ell(y, \hat{y}) = \mathbb{I}\{\hat{y} \neq y\}$ is therefore defined by

$$\ell_{\mathcal{D}}(h) = \mathbb{E}[\ell(Y, h(\mathbf{X}))] = \mathbb{E}[\mathbb{I}\{h(\mathbf{X}) \neq Y\}] = \mathbb{P}(h(\mathbf{X}) \neq Y).$$

In words, the risk of h is the probability that h misclassifies a random \mathbf{X} drawn from the distribution D_X and having a random label Y drawn from the distribution $\{1 - \eta(\mathbf{X}), \eta(\mathbf{X})\}$ on $\{-1, +1\}$.

The Bayes optimal predictor $f^* : \mathcal{X} \rightarrow \{-1, +1\}$ for binary classification is derived as follows

$$\begin{aligned}
f^*(\mathbf{x}) &= \operatorname{argmin}_{\hat{y} \in \{-1, +1\}} \mathbb{E}[\ell(Y, \hat{y}) \mid \mathbf{X} = \mathbf{x}] \\
&= \operatorname{argmin}_{\hat{y} \in \{-1, +1\}} \mathbb{E}[\mathbb{I}\{Y = +1\} \mathbb{I}\{\hat{y} = -1\} + \mathbb{I}\{Y = -1\} \mathbb{I}\{\hat{y} = +1\} \mid \mathbf{X} = \mathbf{x}] \\
&= \operatorname{argmin}_{\hat{y} \in \{-1, +1\}} \left(\mathbb{P}(Y = +1 \mid \mathbf{X} = \mathbf{x}) \mathbb{I}\{\hat{y} = -1\} + \mathbb{P}(Y = -1 \mid \mathbf{X} = \mathbf{x}) \mathbb{I}\{\hat{y} = +1\} \right) \\
&= \operatorname{argmin}_{\hat{y} \in \{-1, +1\}} \left(\eta(\mathbf{x}) \mathbb{I}\{\hat{y} = -1\} + (1 - \eta(\mathbf{x})) \mathbb{I}\{\hat{y} = +1\} \right) \\
&= \begin{cases} -1 & \text{se } \eta(\mathbf{x}) < 1/2, \\ +1 & \text{se } \eta(\mathbf{x}) \geq 1/2. \end{cases}
\end{aligned}$$

Hence, the Bayes optimal classifier predicts the label whose probability is the highest when conditioned on the instance. Finally, it is easy to verify that the Bayes error in this case is $\ell_{\mathcal{D}}(f^*) = \mathbb{E}[\min\{\eta(\mathbf{X}), 1 - \eta(\mathbf{X})\}]$.

Bounding the risk. Next, we study the problem of bounding the risk of a predictor. From now on, we assume $\ell(y, \hat{y}) \in [0, 1]$. However, keep in mind that our analysis continues to hold also when $\ell(y, \hat{y}) \in [0, M]$ for any $M > 0$.

It should be clear that, given an arbitrary predictor h , we cannot directly compute its risk $\ell_{\mathcal{D}}(h)$ with respect to \mathcal{D} because \mathcal{D} is typically unknown (if we knew \mathcal{D} , we could directly construct the Bayes optimal predictor). We thus consider the problem of estimating the risk of a given predictor h . In order to compute this estimate, we use a so-called **test set** $S' = \{(\mathbf{x}'_1, y'_1), \dots, (\mathbf{x}'_n, y'_n)\}$ of examples. Then, we estimate $\ell_{\mathcal{D}}(h)$ with the **test error**, which is the average loss of h on the test set,

$$\widehat{\ell}_{S'}(h) = \frac{1}{n} \sum_{t=1}^n \ell(y'_t, h(\mathbf{x}'_t)) .$$

Under the assumption that the test set is generated through independent draws from \mathcal{D} , the test error corresponds to the sample mean of the risk. Indeed, for each $t = 1, \dots, n$ the example (\mathbf{X}'_t, Y'_t) is an independent draw from \mathcal{D} . Therefore,

$$\mathbb{E} \left[\ell(Y'_t, h(\mathbf{X}'_t)) \right] = \ell_{\mathcal{D}}(h) .$$

In order to compute a confidence interval for the risk using the test error, we can use the following result about the law of large numbers.

Lemma 1 (Chernoff-Hoeffding). *Let Z_1, \dots, Z_n be independent and identically distributed random variables with expectation μ and such that $0 \leq Z_t \leq 1$ for each $t = 1, \dots, n$. Then, for any given $\varepsilon > 0$,*

$$\mathbb{P} \left(\frac{1}{n} \sum_{t=1}^n Z_t > \mu + \varepsilon \right) \leq e^{-2\varepsilon^2 n} \quad \text{and} \quad \mathbb{P} \left(\frac{1}{n} \sum_{t=1}^n Z_t < \mu - \varepsilon \right) \leq e^{-2\varepsilon^2 n} .$$

Using the Chernoff-Hoeffding bound with $Z_t = \ell(y_t, h(x_t)) \in [0, 1]$ we can compute a confidence interval for the risk as follows

$$\mathbb{P} \left(\left| \ell_{\mathcal{D}}(h) - \widehat{\ell}(h) \right| > \varepsilon \right) = \mathbb{P} \left(\ell_{\mathcal{D}}(h) - \widehat{\ell}(h) > \varepsilon \right) + \mathbb{P} \left(\widehat{\ell}(h) - \ell_{\mathcal{D}}(h) > \varepsilon \right) \leq 2e^{-2\varepsilon^2 n} \quad (1)$$

where the probability is computed with respect to random draw of the test set. This inequality shows that the probability that a test set gives a test error $\widehat{\ell}_{S'}(h)$ differing from the true risk $\ell_{\mathcal{D}}(h)$ for more than ε quickly decreases with the size n of the test set.

More specifically, by setting to δ the right-hand side of (1), for any given $\delta \in (0, 1)$, we get that

$$\left| \ell_{\mathcal{D}}(h) - \widehat{\ell}_{S'}(h) \right| \leq \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}$$

holds with probability at least $1 - \delta$ with respect to the random draw of the test set.

The inequality (1) is telling us how to use a test set to estimate the risk of a classifier. More precisely, the inequality shows that the test set, which is how we measure in practice the performance of a classifier on unseen data, is close to the statistical risk with high probability.

In order to understand the occurrence of overfitting, we now study a specific learning algorithm in the statistical framework we just defined. As usual, we assume the algorithm obtains a training set $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$, in input, where $(\mathbf{x}_t, y_t) \in \mathcal{X} \times \mathcal{Y}$. The algorithm generates a predictor $h : \mathcal{X} \rightarrow \mathcal{Y}$ belonging to a given model space \mathcal{H} .

If the algorithm is forced to pick a predictor from \mathcal{H} , then the best thing the algorithm can do is to choose the best possible predictor $h^* \in \mathcal{H}$, that is the one satisfying

$$\ell_{\mathcal{D}}(h^*) = \min_{h \in \mathcal{H}} \ell_{\mathcal{D}}(h) .$$

Thanks to the law of large numbers, we know that the training error $\widehat{\ell}_S(h^*)$ is close to $\ell_{\mathcal{D}}(h^*)$ with high probability with respect to the random draw of the training set on which $\widehat{\ell}_S(h^*)$ is computed.

Consider the algorithm choosing $\widehat{h} \in \mathcal{H}$ minimizing the training error,

$$\widehat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \widehat{\ell}_S(h) .$$

This algorithm, which is very natural in many applications, is known as Empirical Risk Minimizer (ERM).

Unfortunately, we cannot directly apply the Chernoff-Hoeffding bound to \widehat{h} to show that $\ell_{\mathcal{D}}(\widehat{h})$ is close to $\widehat{\ell}_S(\widehat{h})$. The reason is that \widehat{h} is a function of the training set, and thus a random variable. Chernoff-Hoeffding ensures that $\widehat{\ell}_S(h)$ is close to $\ell_{\mathcal{D}}(h)$ for any fixed h , whereas \widehat{h} is not fixed as it depends on the sample.

In order to analyze the risk of \widehat{h} , we then proceed as follows:

$$\begin{aligned} \ell_{\mathcal{D}}(\widehat{h}) &= \ell_{\mathcal{D}}(\widehat{h}) - \ell_{\mathcal{D}}(h^*) && \text{variance error} \\ &+ \ell_{\mathcal{D}}(h^*) - \ell_{\mathcal{D}}(f^*) && \text{bias error} \\ &+ \ell_{\mathcal{D}}(f^*) && \text{Bayes error} \end{aligned}$$

where f^* is the Bayes optimal predictor for \mathcal{D} . The Bayes error is not controllable, because it only depends on \mathcal{D} and the loss ℓ . The bias error arises because \mathcal{H} does not necessarily contain the Bayes optimal predictor. The variance error arises because the risk of \widehat{h} is generally different from the training error of \widehat{h} (which is the smallest in \mathcal{H} by design of the algorithm).

Next, we bound the variance error. For every given training set S , we have that

$$\begin{aligned} \ell_{\mathcal{D}}(\widehat{h}) - \ell_{\mathcal{D}}(h^*) &= \ell_{\mathcal{D}}(\widehat{h}) - \widehat{\ell}_S(\widehat{h}) + \widehat{\ell}_S(\widehat{h}) - \ell_{\mathcal{D}}(h^*) \\ &\leq \ell_{\mathcal{D}}(\widehat{h}) - \widehat{\ell}_S(\widehat{h}) + \widehat{\ell}_S(h^*) - \ell_{\mathcal{D}}(h^*) \\ &\leq |\ell_{\mathcal{D}}(\widehat{h}) - \widehat{\ell}_S(\widehat{h})| + |\widehat{\ell}_S(h^*) - \ell_{\mathcal{D}}(h^*)| \\ &\leq 2 \max_{h \in \mathcal{H}} |\widehat{\ell}_S(h) - \ell_{\mathcal{D}}(h)| \end{aligned}$$

where we used the assumption that \widehat{h} minimizes $\widehat{\ell}_S(h)$ among all $h \in \mathcal{H}$. Therefore, for all $\varepsilon > 0$,

$$\ell_{\mathcal{D}}(\widehat{h}) - \ell_{\mathcal{D}}(h^*) > \varepsilon \quad \Rightarrow \quad \max_{h \in \mathcal{H}} |\widehat{\ell}_S(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2} \quad \Rightarrow \quad \exists h \in \mathcal{H} : |\widehat{\ell}_S(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2} .$$

Since the above chain of implications holds for any realization of the training set, we can write

$$\mathbb{P}\left(\ell_{\mathcal{D}}(\hat{h}) - \ell_{\mathcal{D}}(h^*) > \varepsilon\right) \leq \mathbb{P}\left(\exists h \in \mathcal{H} : |\hat{\ell}(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2}\right).$$

We now study the case $|\mathcal{H}| < \infty$, that is when the model space contains a finite number of predictors. Note that the event

$$\exists h \in \mathcal{H} : |\hat{\ell}(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2}$$

is the union over $h \in \mathcal{H}$ of the (not necessarily disjoint) events

$$|\hat{\ell}(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2}$$

We may then use the union bound

$$\mathbb{P}(A_1 \cup \dots \cup A_n) \leq \sum_{i=1}^n \mathbb{P}(A_i)$$

which holds for any collection A_1, \dots, A_n of events. By doing so we get

$$\begin{aligned} \mathbb{P}\left(\exists h \in \mathcal{H} : |\hat{\ell}(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2}\right) &= \mathbb{P}\left(\bigcup_{h \in \mathcal{H}} \left(|\hat{\ell}(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2}\right)\right) \\ &\leq \sum_{h \in \mathcal{H}} \mathbb{P}\left(|\hat{\ell}(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2}\right) \\ &\leq |\mathcal{H}| \max_{h \in \mathcal{H}} \mathbb{P}\left(|\hat{\ell}(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2}\right) \\ &\leq |\mathcal{H}| 2e^{-m\varepsilon^2/2} \end{aligned} \tag{2}$$

where in the last step we used the Chernoff-Hoeffding bound.

In conclusion, we have that

$$\mathbb{P}\left(\ell_{\mathcal{D}}(\hat{h}) - \ell_{\mathcal{D}}(h^*) > \varepsilon\right) \leq \mathbb{P}\left(\exists h \in \mathcal{H} : |\hat{\ell}(h) - \ell_{\mathcal{D}}(h)| > \frac{\varepsilon}{2}\right) \leq 2|\mathcal{H}|e^{-m\varepsilon^2/2}. \tag{3}$$

Setting the left-hand side of (3) equal to δ and solving for ε we obtain that

$$\ell_{\mathcal{D}}(\hat{h}) \leq \ell_{\mathcal{D}}(h^*) + \sqrt{\frac{2}{m} \ln \frac{2|\mathcal{H}|}{\delta}}$$

holds with probability at least $1 - \delta$ with respect to the random draw of a training set of cardinality m .

Note that the risk of \hat{h} is bounded through a sum of two terms: $\ell_{\mathcal{D}}(h^*)$ and $\sqrt{\frac{2}{m} \ln \frac{2|\mathcal{H}|}{\delta}}$. Lacking further information on \mathcal{D} , and for a given cardinality m of the training set, in order to decrease $\sqrt{\frac{2}{m} \ln \frac{2|\mathcal{H}|}{\delta}}$, which is our bound on the variance error, we must decrease $|\mathcal{H}|$. But decreasing $|\mathcal{H}|$ might cause an increase of $\ell_{\mathcal{D}}(h^*)$, which produces a corresponding increase of the bias error. In light of this statistical analysis, we may say that overfitting in the ERM algorithm is caused by an unbalance between the variance error and the bias error. In particular, overfitting is when the

variance error is much larger than the bias error, and underfitting is when the bias error is much larger than the variance error.

In the proof of the bound on the variance error, we have also shown in (2) that

$$\forall h \in \mathcal{H} \quad |\widehat{\ell}(h) - \ell_{\mathcal{D}}(h)| \leq \sqrt{\frac{1}{2m} \ln \frac{2|\mathcal{H}|}{\delta}}$$

with probability at least $1 - \delta$ with respect to the random draw of the training set. This implies that when the cardinality of the training set is sufficiently large with respect to $\ln |\mathcal{H}|$, then the training error $\widehat{\ell}_S(h)$ becomes a good estimate for the statistical risk $\ell_{\mathcal{D}}(h)$ *simultaneously* for all predictors $h \in \mathcal{H}$. In this regime, namely when the law of large numbers holds uniformly with respect to the choice of $h \in \mathcal{H}$, it is clear that any algorithm ranking the predictors in \mathcal{H} according to their training error is protected from overfitting.