

Intro to Learning Theory

Ruth Urner

1 Machine Learning and Learning Theory

Coming soon..

2 Formal Framework

2.1 Basic notions

In our formal model for machine learning, the instances to be classified are members of a set \mathcal{X} , the *domain set* or *feature space*. Instances are to be classified into a *label set* \mathcal{Y} . For now (and most of the class), we assume that the label set is binary, that is $\mathcal{Y} = \{0, 1\}$. For example, an instance $x \in \mathcal{X}$ could be an email and its label indicates whether the email is spam ($y = 1$) or not spam ($y = 0$). We often assume that the instances are represented as real-valued vectors, that is $\mathcal{X} \subseteq \mathbb{R}^d$ for some dimension d .

A *predictor* or *classifier* is a function $h : \mathcal{X} \rightarrow \mathcal{Y}$. A *learner* is a function that takes some training data and maps it to a predictor. We let the *training data* be denoted by a sequence $S = ((X_1, Y_1), \dots, (X_n, Y_n))$. Then, formally, a learner \mathcal{A} is a function

$$\begin{aligned}\mathcal{A} &: \bigcup_{i=1}^{\infty} (\mathcal{X} \times \mathcal{Y})^i \rightarrow \mathcal{Y}^{\mathcal{X}} \\ \mathcal{A} &: S \mapsto h,\end{aligned}$$

where $\mathcal{Y}^{\mathcal{X}}$ denotes the set of all functions from set \mathcal{X} to set \mathcal{Y} . For convenience, when the learner is clear from context, we use the notation h_n to denote the output of the learner on data of size n , that is $h_n = \mathcal{A}(S)$ for $|S| = n$.

The goal of learning is produce a predictor h that correctly classifies not only the training data, but also future instances that it has not seen yet. We thus need a mathematical description of how the environment produces instances. In particular, we would like to model that the environment (or nature) remains somehow stable, that the process that generated the training data is the same that will generate future data.

We model the data generation as a probability distribution P over $\mathcal{X} \times \mathcal{Y} = \mathcal{X} \times \{0, 1\}$. We further assume that the instances (X_i, Y_i) are *i.i.d.* (independently and identically distributed) according to P .

The performance of a classifier h on an instance (X, Y) is measured by a *loss function*. A loss function is a function

$$\ell : (\mathcal{Y}^{\mathcal{X}} \times \mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}.$$

The value $\ell(h, X, Y) \in \mathbb{R}$ indicates “how badly h predicts on example (X, Y) ”. We will, for now, work with the *binary loss* (or *0/1-loss*), defined as

$$\ell(h, X, Y) = \mathbf{1}[h(X) \neq Y],$$

where $\mathbf{1}[p]$ denotes the *indicator function* of predicate p , that is $\mathbf{1}[p] = 1$ if p is true and $\mathbf{1}[p] = 0$ if p is false. The binary loss is 1 if the prediction of h on example (X, Y) is wrong. If the prediction is correct, no loss is suffered and the binary loss assigns value 0.

We can now formally phrase the goal of learning, as aiming for a classifier that has low loss on expectation over the data generating distribution. That is, we would like to output a classifier that has low *expected loss*, or *risk*, defined as

$$L(h) = \mathbb{E}_{(X,Y) \sim P}[\ell(h, X, Y)] = \mathbb{E}_{(X,Y) \sim P}[\mathbf{1}[h(X) \neq Y]].$$

Since our loss function assumes only values in $\{0, 1\}$, the above expectation is equal to the probability of generating an example X on which h makes a wrong prediction. That is, we have

$$L(h) = \mathbb{E}_{(X,Y) \sim P} \mathbb{E}_{(X,Y) \sim P}[\mathbf{1}[h(X) \neq Y]] = \mathbb{P}_{(X,Y) \sim P}[h(X) \neq Y].$$

Note however, that the learner does not get to see the data generating distribution. It can thus not merely output a classifier of lowest expected loss. The learner needs to make its decisions based on the data S . Given a classifier h and data S , the learner can evaluate the *empirical risk* of h on S

$$L_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}[h(X_i) \neq Y_i].$$

2.2 On the relation of empirical and true risk

A natural strategy for the learner would be, to simply output a function that has small empirical risk. In favor of this approach, we now show that the empirical risk is an unbiased estimator of the true risk.

Claim 1. *For all functions $h : \mathcal{X} \rightarrow \{0, 1\}$ and for all sample sizes n we have*

$$\mathbb{E}_S L_n(h) = L(h)$$

Proof.

$$\begin{aligned}
\mathbb{E}_S L_n(h) &= \mathbb{E}_S \frac{1}{n} \sum_{i=1}^n \mathbf{1}[h(X_i) \neq Y_i] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_S \mathbf{1}[h(X_i) \neq Y_i] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{(X,Y)} \mathbf{1}[h(X) \neq Y] \\
&= \frac{1}{n} \sum_{i=1}^n L(h) \\
&= L(h)
\end{aligned}$$

where the second equality holds by linearity of expectation, and the third inequality holds since that expectation depends only on one (the i -th) example in S . \square

Thus, for any fixed function, the empirical risk gives us an unbiased estimate of the quantity that we are after, the true risk. Note that this holds even for small sample sizes. Moreover, by the *law of large numbers*, the above claim implies that, with large sample sizes, the empirical risk of a classifier converges to its true risk (in probability). As we see more and more data, the empirical risk of a function becomes a better and better estimate of its true risk.

This may lead us to believe that the simple learning strategy of just finding some function with low empirical risk should succeed at achieving low true risk as we see more and more data. However, the following phenomenon shows that this strategy can in fact go wrong arbitrarily badly.

Claim 2. *There exists a distribution P and a learner, such that for all n we have*

$$L_n(h_n) = 0 \text{ and } L(h_n) = 1$$

Proof. As the data generating distribution, consider the uniform distribution over $\mathbb{R} \times \{1\}$. That is, in any sample S , generated by this P , the examples are labeled with 1, that is $S = ((X_1, 1), \dots, (X_n, 1))$. We construct a “stubborn” learner \mathcal{A} . The stubborn learner outputs a function that agrees with the sample’s labels on points that were in the sample S , but keeps believing that the label is 0 everywhere else. Formally:

$$h_n(X) = \mathcal{A}(S)(X) = \begin{cases} 1 & \text{if } (X, 1) \in S \\ 0 & \text{otherwise} \end{cases}$$

Now we clearly have $L_n(h_n) = 0$ for all n . However, since S is finite, the set of instances X on which h_n predicts 1 has measure 0. Thus, with probability 1, h_n outputs the incorrect label 0. Thus $L(h) = 1$. \square

The difference between the situations in the above two claims is that, in the second case, the function h_n depends on the data. While, for every fixed function h (fixed before the data is seen), the empirical risk estimates converge to the true risk of this function,

this convergence is not uniform over all functions. Claim 2 shows that, at any given sample size, *there exist* functions, for which true and empirical risk are arbitrarily far apart.

Now, in machine learning, we do want the function that the learner outputs to be able to depend on the data. Furthermore, the learner only ever gets to see a finite amount of data. We have seen that, for any finite sample size, that is, on any finite amount of data, the empirical risk can be a very bad indicator of the true risk of a function.

Basic questions of learning theory thus are: How can we control the (true) risk of a function learned based on a finite amount of data? Can we identify situations where we can relate the true and empirical risk?

2.3 Fixing a hypothesis space

We have seen that, if we want our learned function h_n to depend on the data, we have to change the rules for the learner. In Claim 2, we let the learner output any function it wanted. This resulted in the learner adapting itself very well to the data it has seen in the sample, achieving 0 empirical risk, while not making any progress towards predicting well on unseen examples.

The construction of Claim 2 is an extreme version of a phenomenon called *overfitting*. In informal terms, if a learning method has too much freedom with regards to the functions it can output, it may “overadapt” to the training data, rather than extracting structure that will also apply to the unseen examples. Overfitting is frequently encountered phenomenon in practice, that one has to guard against. To prevent the learner from overfitting, we need to *restrict the class of predictors*.

A *hypothesis class* \mathcal{H} is a set of predictors $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$. Instead of allowing the learner to output any function, we will now consider learners that output functions from \mathcal{H} . We will see that, in many cases, fixing the hypothesis class before we see the data, will let us regain control over the relation between empirical and true risk.

However, fixing the Hypothesis class also means that there may not be any good function in the class. We will thus rephrase the goal of learning, to only require the learner to come up with a function that is (approximately) *as good as the best function in the class \mathcal{H}* .

Thus, our new goal is to show that

$$L(h_n) \leq \inf_{h \in \mathcal{H}} L(h) + f(n),$$

where f is decreasing function of sample size n . That is, as we see more and more data, we would like that the true risk of the output of the learner approaches the best risk possible with the class \mathcal{H} . Or equivalently, we would like to show that

$$L(h_n) - \inf_{h \in \mathcal{H}} L(h) \leq f(n).$$

2.4 Learnability of finite classes

We now show that the above goal is achievable for finite classes $\mathcal{H} = \{h_1, \dots, h_N\}$. We will analyze the learner ERM (Empirical Risk Minimization), which outputs a function from \mathcal{H} that has minimal empirical risk.

$$\text{ERM} : S \mapsto \hat{h}_n \in \operatorname{argmin}_i L_n(h_i)$$

There may be several functions in \mathcal{H} that have lowest empirical risk on a data set S . But since every function $h \in H$ has some empirical risk on data S , and the empirical risk can only assume finitely many values (namely multiples of $\frac{1}{n} = \frac{1}{|S|}$), the argmin is a nonempty subset of \mathcal{H} . A learner is an ERM learner, if it always outputs some function from this subset.

For now, we will further make a simplifying assumption on the data generating distribution P . We will assume that P is *realizable with respect to the class* \mathcal{H} . A distribution is realizable with respect to a hypothesis class \mathcal{H} if there is an $h^* \in \mathcal{H}$ with $L(h^*) = 0$.

Theorem 1. *Let $\mathcal{H} = \{h_1, \dots, h_N\}$ and $\delta \in (0, 1]$. Under the realizability assumption, we have with probability at least $(1 - \delta)$ over the generation of the sample S*

$$L(\hat{h}_n) \leq \frac{\log N + \log(1/\delta)}{n}.$$

Proof. Note that $L_n(h^*) = 0$ for all possible samples S . Thus, for any $\epsilon > 0$, ERM only outputs a function error larger than ϵ , if $L_n(\hat{h}_n) = 0$ while $L(\hat{h}_n) \geq \epsilon$.

For every $h \in \mathcal{H}$ with $L(h) > \epsilon$, we have

$$\mathbb{P}_S \{L_n(h) = 0\} \leq (1 - \epsilon)^n \leq e^{-\epsilon n}$$

(Recall that, for all $x \in \mathbb{R}$, we have $(1 + x) \leq e^x$.)

Let \mathcal{H}_ϵ denote the set of functions h in \mathcal{H} with $L(h) > \epsilon$. We get, using the union bound,

$$\begin{aligned} \mathbb{P}_S \left\{ L(\hat{h}_n) \geq \epsilon \right\} &\leq \mathbb{P}_S \{ \exists h \in \mathcal{H}_\epsilon : L_n(h) = 0 \} \\ &= \mathbb{P}_S \{ \vee_{h \in \mathcal{H}_\epsilon} L_n(h) = 0 \} \\ &\leq |\mathcal{H}_\epsilon|(1 - \epsilon)^n \\ &\leq |\mathcal{H}|(1 - \epsilon)^n \\ &\leq |\mathcal{H}|e^{-\epsilon n} = Ne^{-\epsilon n} \end{aligned}$$

Now we set $\epsilon = \frac{\log \frac{1}{\delta} + \log N}{n}$.

Plugging in this value for ϵ , we have shown

$$\mathbb{P}_S \left\{ L(\hat{h}_n) \geq \frac{\log N + \log(1/\delta)}{n} \right\} \leq \delta$$

which is equivalent to the statement of the theorem. \square

Thus, under the realizability we have $L(h_n) - \min_{h \in \mathcal{H}} L(h) \leq f(n)$ for $f(n) = \frac{\log N + \log(1/\delta)}{n}$.