

Title: Learnable function class

URL: [https://en.wikipedia.org/wiki/Learnable\\_function\\_class](https://en.wikipedia.org/wiki/Learnable_function_class)

PageID: 48827727

Categories: Category:Machine learning

Source: Wikipedia (CC BY-SA 4.0).

-----

In statistical learning theory , a learnable function class is a set of functions for which an algorithm can be devised to asymptotically minimize the expected risk , uniformly over all probability distributions. The concept of learnable classes are closely related to regularization in machine learning , and provides large sample justifications for certain learning algorithms.

Definition

Background

Let  $\Omega = X \times Y = \{ (x, y) \}$  be the sample space, where  $y$  are the labels and  $x$  are the covariates (predictors).  $F = \{ f : X \rightarrow Y \}$  is a collection of mappings (functions) under consideration to link  $x$  to  $y$ .  $L : Y \times Y \rightarrow \mathbb{R}$  is a pre-given loss function (usually non-negative). Given a probability distribution  $P(x, y)$  on  $\Omega$ , define the expected risk  $I_P(f)$  to be:

The general goal in statistical learning is to find the function in  $F$  that minimizes the expected risk. That is, to find solutions to the following problem: [ 1 ]

But in practice the distribution  $P$  is unknown, and any learning task can only be based on finite samples. Thus we seek instead to find an algorithm that asymptotically minimizes the empirical risk, i.e., to find a sequence of functions  $\{ \hat{f}_n \}_{n=1}^{\infty}$  that satisfies

One usual algorithm to find such a sequence is through empirical risk minimization .

Learnable function class

We can make the condition given in the above equation stronger by requiring that the convergence is uniform for all probability distributions. That is:

The intuition behind the more strict requirement is as such: the rate at which sequence  $\{ \hat{f}_n \}$  converges to the minimizer of the expected risk can be very different for different  $P(x, y)$ . Because in real world the true distribution  $P$  is always unknown, we would want to select a sequence that performs well under all cases.

However, by the no free lunch theorem , such a sequence that satisfies ( 1 ) does not exist if  $F$  is too complex. This means we need to be careful and not allow too "many" functions in  $F$  if we want ( 1 ) to be a meaningful requirement. Specifically, function classes that ensure the existence of a sequence  $\{ \hat{f}_n \}$  that satisfies ( 1 ) are known as learnable classes . [ 1 ]

It is worth noting that at least for supervised classification and regression problems, if a function class is learnable, then the empirical risk minimization automatically satisfies ( 1 ). [ 2 ] Thus in these settings not only do we know that the problem posed by ( 1 ) is solvable, we also immediately have an algorithm that gives the solution.

Interpretations

If the true relationship between  $y$  and  $x$  is  $y = f^*(x)$ , then by selecting the appropriate loss function,  $f^*$  can always be expressed as the minimizer of the expected loss across all possible functions. That is,

Here we let  $F^*$  be the collection of all possible functions mapping  $X$  onto  $Y$ .  $f^*$  can be interpreted as the actual data generating mechanism. However, the no free lunch theorem tells us that in practice, with finite samples we cannot hope to search for the expected risk minimizer over  $F^*$ . Thus we often consider a subset of  $F^*$ ,  $F$ , to carry out searches on. By doing so, we risk that  $f^*$  might not be an element of  $F$ . This tradeoff can be mathematically expressed as

In the above decomposition, part (b) does not depend on the data and is non-stochastic. It describes how far away our assumptions ( $F$ ) are from the truth ( $F^*$ ). (b) will be strictly greater than 0 if we make assumptions that are too strong ( $F$  too small). On the other hand, failing to put enough restrictions on  $F$  will cause it to be not learnable, and part (a) will not stochastically converge to 0. This is the well-known overfitting problem in statistics and machine learning literature.

Example: Tikhonov regularization

A good example where learnable classes are used is the so-called Tikhonov regularization in reproducing kernel Hilbert space (RKHS). Specifically, let  $F^*$  be an RKHS, and  $\|\cdot\|_2$  be the norm on  $F^*$  given by its inner product. It is shown in [3] that  $F = \{f : \|f\|_2 \leq \gamma\}$  is a learnable class for any finite, positive  $\gamma$ . The empirical minimization algorithm to the dual form of this problem is

This was first introduced by Tikhonov [4] to solve ill-posed problems. Many statistical learning algorithms can be expressed in such a form (for example, the well-known ridge regression).

The tradeoff between (a) and (b) in (2) is geometrically more intuitive with Tikhonov regularization in RKHS. We can consider a sequence of  $F_\gamma$ , which are essentially balls in  $F^*$  with centers at 0. As  $\gamma$  gets larger,  $F_\gamma$  gets closer to the entire space, and (b) is likely to become smaller. However we will also suffer smaller convergence rates in (a). The way to choose an optimal  $\gamma$  in finite sample settings is usually through cross-validation.

Relationship to empirical process theory

Part (a) in (2) is closely linked to empirical process theory in statistics, where the empirical risk  $\sum_{i=1}^n L(y_i, f(x_i))$ ,  $f \in F$  are known as empirical processes. [5] In this field, the function class  $F$  that satisfies the stochastic convergence

are known as uniform Glivenko–Cantelli classes. It has been shown that under certain regularity conditions, learnable classes and uniformly Glivenko-Cantelli classes are equivalent. [1] Interplay between (a) and (b) in statistics literature is often known as the bias-variance tradeoff.

However, note that in [2] the authors gave an example of stochastic convex optimization for General Setting of Learning where learnability is not equivalent with uniform convergence.

References