



## Why

We want to discuss an inductor’s performance on (possibly) incomplete, but still consistent datasets.

## Definition

We define a probability measure over the inputs. Let  $(X, \mathcal{X}, \mu)$  be a probability space and  $(Y, \mathcal{Y})$  a measurable space. Let  $f : X \rightarrow Y$  measurable. We call the pair  $((X, \mathcal{X}, \mu), f)$  a *probabilistic data-generation model*.

The idea is that to each instance  $x \in X$ , we associate a label  $f(x) \in Y$ . We call the law  $\mu$  the *data-generating distribution* or *underlying distribution* and call  $f$  the *correct labeling function*. Many authors refer to a probabilistic data-generation model as the *statistical learning (theory) framework*.

## Representative datasets

Let  $(X, \mathcal{X}, \mu)$  and  $f : X \rightarrow Y$  be a probabilistic data-generation model. We put a measure on the set of datasets  $X^n$  so that such “nonrepresentative” datasets have low measure. We use the product measure  $\mu^n$ . We interpret this as a model of the inductor as acting on a training set of correctly labeled, independent and identically distributed inputs.

Let  $0 < \delta < 1$ . A set of datasets  $\mathcal{S} \subset X^n$  is  $1 - \delta$ -*representative* if  $\mu^n(\mathcal{S}) \geq 1 - \delta$ . We interpret this as follows: “the dataset we see is in  $\mathcal{S}$  with high probability.” Roughly

speaking, if  $\delta \ll 1$  then the set  $\mathcal{S}$  is “fairly likely.”

We think of  $\mathcal{S}$  as the set of datasets for which our learning algorithm “succeeds,” those datasets which are reasonable. In other words, the probability of a reasonable training set is  $1 - \delta$ . We call  $\delta$  the *confidence parameter*.

## Predictor errors

The *error of a predictor*  $h : X \rightarrow Y$  with respect to the probabilistic data-generation model is

$$\mu(\{x \in \mathcal{X} \mid h(x) \neq f(x)\}).$$

In other words, the error is the probability (w.r.t. the underlying distribution) that the predictor  $h$  mislabels a point. The error is measured with respect to the underlying distribution  $\mu$  and correct labeling function  $f$ .

The *accuracy* of a predictor is

$$\mu(\{x \in X \mid h(x) = f(x)\}).$$

It is the measure of the set of points for which the predictor matches the correct labeling function: In other words, it is one minus the error of the predictor.

## Inductor errors

Let  $S = (x_1, y_1), \dots, (x_n, y_n)$  a dataset of  $n$  paired records in  $X \times Y$ . Let  $\mathcal{M}_{X \rightarrow Y}$  be the set of measurable functions from  $X$  to  $Y$ . The *dataset error of an inductor*  $A : (X \times Y)^n \rightarrow \mathcal{M}_{X \rightarrow Y}$  on  $S$  is the error of the predictor  $A(S)$ .

Notice that if we have a probability measure on  $X^n$ , we can discuss the measure of a set of training sets. We may be interested in statements like: the measure of the set of training sets for which the error of the inductor is nonzero is small.

### **Other terminology**

A *hypothesis class* is a subset of the measurable functions from  $X \rightarrow Y$ . Other names for the error of a classifier include the *generalization error*, the *risk* or the *true error* or *loss*.

