

A Formal Model - The statistical Learning Framework

Domain Set: An arbitrary set, X . This is the set of objects that we may wish to label. These data points are represented by a vector of features. This set is from some probability distribution D .

Label Set: An arbitrary set of all possible labels, Y .

As an example, in a two-element set we have $Y = \{0, 1\}$

Training data:

$S = ((x_1, y_1), \dots, (x_m, y_m))$, sample set $X \times Y$ used to train the learner.

Learner's Output: The learner will output a feature hypothesis (h) which is a function $h: X \rightarrow Y$, used to predict the label of new data points.

Target Function:

$$(f: X \rightarrow Y)$$

Target function maps the X to Y . Target function is unknown. Hypothesis, h actually approximates the target function.

Probability Distribution: Suppose $A \subseteq X$, the probability distribution D will assign a number $D(A)$, which determines how likely it is to observe a point $x \in A$.

Measure of Success :

The error of h is the probability that $h(x) \neq f(x)$ for any given x taken from the underlying distribution.

Error is given by

$$L_{D,f}(h) = P_{\text{UND}}[h(x) \neq f(x)] = D(\{x : h(x) \neq f(x)\})$$

$L_{D,f}(h)$ is also called Generalization Error, the risk or the true error of h .

Empirical Risk Minimization

The goal of the learning algorithm is to find h_s that minimizes the error w.r.t to the unknown $D \& f$.

Since $D \& f$ are unknown, the true error is not available to the learner. A useful notion of error that can be calculated by the learner is the training error. — the error the classifier incurs over the training sample

$$L_s(h) = \frac{|\{i \in [m] : h(x_i) \neq y_i\}|}{m}$$

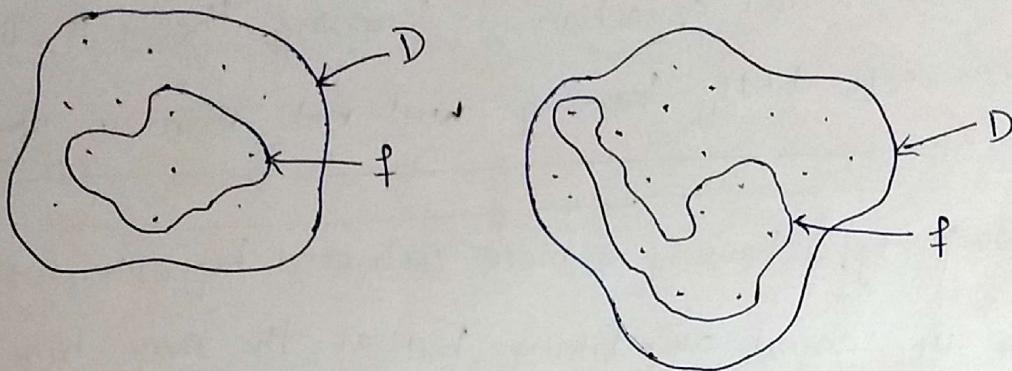
where $[m] = \{1, 2, \dots, m\}$

$L_s(h)$ is termed as the Empirical error | Empirical risk

The learning paradigm — coming up with a predictor h that minimizes $L_s(h)$ — is called Empirical Risk Minimization | EMR

Note: Since the training sample is the snapshot of the world that is available to the ^{learner}, it makes sense to search for a solution that works well on that data. This can result in over-fitting.

Side Note:



Overfitting model

$$h_S(x) = \begin{cases} y_i & \text{if } \exists i \in [m] \text{ s.t. } x_i = x \\ 0 & \text{otherwise} \end{cases}$$

Empirical Risk Minimization with Inductive Bias

A common solution to overfitting is to apply the ERM learning rule over a restricted search space. Formally, the learner should choose in advance a set of predictors. This set is called hypothesis set H . Each $h \in H$ is a function mapping from $X \rightarrow Y$. For a given class H , and a training sample, S , the $\underset{H}{\text{ERM}}_H$ learner uses the ERM rule to choose a predictor $h \in H$, with the lowest possible error over S . Formally

$$\text{ERM}_H(S) \in \operatorname{argmin}_{h \in H} L_S(h)$$

- ① Learning has lot to do with prior knowledge
- Ex: Bayes learners have prior knowledge
- ② Prior knowledge less data points - more sample data points

Where argmin stands for the set of hypotheses in H that achieve the minimum value of $L_S(h)$ over H . By restricting the learner to choosing a predictor from H , we bias it toward a particular set of predictors. Such restrictions are often called an inductive bias.

A fundamental question in learning theory is, over which hypothesis set ERM_H learning will not result in overfitting.

Intuitively, choosing a more restricted hypothesis class better protects us against overfitting but at the same time might cause us a stronger inductive bias.

Finite Hypothesis Set (then ERM_H will not overfit)

The hypothesis set contains infinite hypothesis. Then the simplest type of restriction is imposing an upper bound on its size (the number of predictors h in H)

$$h_S \in \underset{h \in H}{\text{argmin}} L_S(h)$$

Realizability assumption implies there exists $h^* \in H$ s.t. $L_{D,f}(h^*) = 0$, this implies $L_S(h^*) = 0$ for any sample drawn from D . $\&$ labelled by f . This means

$$h^* = f$$

The i.i.d assumption: (Independently & Identically Distributed).

Every x_i in S is freshly sampled according to the distribution D & then labelled according to the labelling function, f .

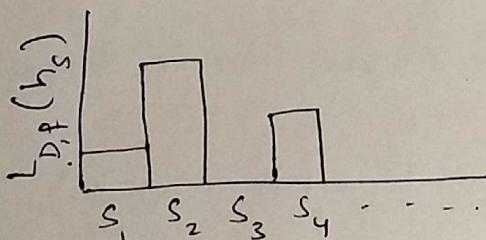
(3)

denote this assumption by $S \sim D^m$, where m is the size of $S \subseteq D^m$, the probability distribution over m tuples.

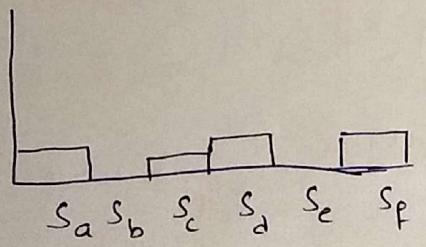
Intuitively, the training set S is a window through which the learner gets partial information about the distribution D & the labelling function, f . So, the larger the sample, the better it is.

Since $L_{D,f}(h_s)$ depends on the training set, S , and that training set is picked by a random process, there is randomness in the choice of the predictor h_s & consequently, in the risk $L_{D,f}(h_s)$. Formally, we say $L_{D,f}(h_s)$ is a random variable.

Random Variable: ω is a variable whose possible values are outcomes of a random phenomenon.



Here S_1, S_2, \dots all have the same sample size.



Here S_a, S_b, \dots all have same sample size but bigger than S_1, S_2, \dots

Let us denote probability of getting a non representative sample by S , and call $(1-S)$ the confidence parameter of our prediction.

On top of that, since we cannot guarantee perfect labels prediction, we introduce another parameter for the quality of prediction, the accuracy parameter, commonly denoted by ϵ .

We interpret the event $L_{D,f}(h_s) > \epsilon$ as failure of the learner, while if $L_{D,f}(h_s) \leq \epsilon$ we view the output of the algorithm as an approximately correct predictor.

Therefore we are interested in upper bounding the probability to sample of instances that will lead to failure of the learner.

Formally, let $S|_x = (x_1, x_2, \dots, x_m)$ be the instances of the training set. We would like to upper bound

$$D^m(\{S|_x : L_{D,f}(h_s) > \epsilon\}) \quad \left\{ \begin{array}{l} \text{find out the} \\ \text{optimal } m \end{array} \right\}$$

Let H_B be the set of "bad" hypotheses, that is,

$$H_B = \{ h \in H : L_{D,f}(h) > \epsilon \}$$

In addition, let

$$M = \{ S|_x : \exists h \in H_B, L_S(h) = 0 \}$$

be the set of misleading samples: Namely, for every

(4)

$S|_x \in M$, there is a "bad" hypothesis, $h \in H_B$, that looks like a "good" hypothesis on $S|_x$. Now, recall that we would like to bound the probability of the event $L_{D,f}(h_s) > \epsilon$. But since the realizability assumption implies that $L_s(h_s) = 0$, it follows that the event $L_{D,f}(h_s) > \epsilon$ can only happen if for some $h \in H_B$ we have $L_s(h) = 0$. In other words, this event will only happen if our sample is in the set of misleading samples, M . Formally, we have shown that

$$\{S|_x : L_{D,f}(h_s) > \epsilon\} \subseteq M$$

Note that we can rewrite M as

$$M = \bigcup_{h \in H_B} \{S|_x : L_s(h) = 0\}$$

Hence,

$$D^m(\{S|_x : L_{(D,f)}(h_s) > \epsilon\}) \leq D^m(M) = D^m\left(\bigcup_{h \in H_B} \{S|_x : L_s(h) = 0\}\right)$$

Next, we upper bound the right hand side of the preceding equation using the union bound - a basic

Property of probabilities.

Union Bound For any two sets A, B and a distribution D we have $D(A \cup B) \leq D(A) + D(B)$

Applying the union bound to the right hand side $y_{1 \in S}$,

$$D^m(\{S|_x : L_{D,f}(h_s) > \epsilon\}) \leq \sum_{h \in H_B} D^m(\{S|_x : L_s(h) = 0\})$$

$$\begin{aligned} D^m(\{S|_x : L_s(h) = 0\}) &= D^m(\{S|_x : \forall i, h(x_i) = f(x_i)\}) \\ &= \prod_{i=1}^m D(\{x_i : h(x_i) = f(x_i)\}) \end{aligned}$$

For each individual sampling of an element of the training set we have

$$D(\{x_i : h(x_i) = y_i\}) = 1 - L_{D,f}(h) \leq 1 - \epsilon,$$

Combining the previous equation with ϵ using the inequality $1 - \epsilon \leq e^{-\epsilon}$ we obtain that for every $h \in H_B$,

$$D^m(\{S|_x : L_s(h) = 0\}) \leq (1 - \epsilon)^m \leq e^{-\epsilon m}$$

$$\text{Finally } D^m(\{S|_x : L_{D,f}(h_s) > \epsilon\}) \leq |H_B| e^{-\epsilon m} \leq |H| e^{-\epsilon m}$$

Important

(5)

Let H be a finite hypothesis class. Let $\delta \in (0,1)$ & $\epsilon > 0$ and let m be an integer that satisfies

$$m \geq \frac{\log(|H|/\delta)}{\epsilon}$$

Then, for any labeling function, f , and for any distribution, D , for which the realizability assumption holds (that is, for some $h \in H$, $L_{D,f}(h) = 0$), with probability of at least $1 - \delta$ over the choice of an i.i.d sample S of size m , we have that for every ERM hypothesis, h_S , it holds that

$$L_{D,f}(h_S) \leq \epsilon.$$

The preceding corollary tells us that for a sufficiently large m , the ERM_H rule over a finite hypothesis class will be probably (with confidence $1 - \delta$) approximately (up to an error of ϵ) correct.

Finally, the Machine Learning Diagram

