

Introduction to Machine Learning -

Machine learning in a general sense refers to the culmination of useful data obtained from given data to be able to process previously unseen cases of the same category. A good learner should be able to generalize data, and this is referred to as **inductive reasoning**.

Great care must be taken to ensure that only meaningful data is being extracted, and not the properties which are favourable by chance, like in the **pigeon superstition** example.

Having prior knowledge regarding the property in question greatly helps to deduce whether the information is useful, or not. The bias obtained due to prior knowledge is called **inductive bias**. However, having a very rigid prior makes the learning to be less flexible.

* Supervised Learning :-

The training data provided has additional information than the testing data, such as labels indicating what the given picture represents. This additional data "supervises" the learning process, hence the name.

* Unsupervised Learning :-

No difference between training and testing data. The algorithm needs to find classifications by itself. Clustering of data sets is a typical example.

* Reinforcement Learning

Training examples contain more information than testing examples. This involves the agent interacting with the environment in a Trial-and-error manner. A good example for this would be Q-Learning.

- Other Classifications Include -

- Active and passive
 - Active Learner refers to when the agent actively seeks for feedback and user input.
Passive Learners usually do neither
- Helpfulness of teacher
 - Humans trying to understand nature can be labelled as nature being a passive teacher, as nature simply doesn't care what humans do.
 - A security software learning by dealing with hackers would model the hackers adversarial teachers, and they act as "worst-case" training scenario.
- Online vs Batch Learning
 - Batch Learning usually has large amounts of data for the agent to learn from before making the final call, whereas online learning requires the agent to give answers on the go, making mistakes as it learns.

Formal model - Nomenclature

* Domain Set :-

- A set X which contains all the objects we wish to label. Each element in X is represented by a vector of its features.

* Label Set :-

- A set Y which contains all possible labels that every element in X may be awarded.

* Training data :-

- A finite sequence of pairs, where each pair $\in X \times Y$ (looking at supervised learning) which the agent has access to.

* Learner's Output :-

- A function $h: X \rightarrow Y$ obtained from analyzing the training data.

The function is also called as classifier, predictor, Hypothesis or a prediction rule.

- Let the training data be given by S , the predictor obtained is represented as $A(S)$ where 'A' is the predictor.

* Data generation model :-

- We first assume the instances in X are generated via some pdf D , and that there is a "correct" labelling function $f: X \rightarrow Y$. Both D, f are unknown to the agent. The training data is obtained by obtaining a sample $x_i \in X$ and using f to get its correct label.

* Measure of Success !-

The error of a classifier is the probability that the label given by the agent doesn't align with the "true label". It is represented by L .

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

Loss of estimator 'h' Prob. that By LOTUS, write
 D-distr. of x $h \neq f$ in terms of D.
 f - "true function"

L is also known as risk, generalization error, true error, etc.

- However, the agent has no idea what the values of D , f are; so the true error cannot be calculated by the agent. We therefore define **Empirical error** using the training data as follows:-

$$\text{empirical risk of } h = \frac{\sum_{i=1}^m \mathbb{I}_{\{h(x_i) \neq y_i\}}}{m} \rightarrow \begin{array}{l} \text{no. of } x_i \text{ where } h \\ \text{is wrong} \end{array}$$

We thus try to minimize the empirical risk. The process of coming up with an estimator to minimize empirical risk is Empirical Risk Minimization (ERM). However, this does have some problems.

Problem

Overfitting - "If it's too good to be true, it probably is."

- Consider the estimator $\hat{h}(x) = \begin{cases} y_i & \text{if } \exists x_i \in S \text{ and } x = x_i \\ 0 & \text{otherwise} \end{cases}$

The estimator can be seen to have empirical error as zero, meanwhile its true error is 0.5, which is horrible. Such a phenomenon where the estimator performs amazingly well with train data, but poorly outside, is called as **Overshooting**.

* ERM using Inductive bias :-

- Suppose that you define a set of functions, H , using prior knowledge of the problem. Also, suppose that $\forall h \in H$, overfitting does not occur. From this bias, we can choose an estimator $h \in H$ such that $\forall h' \in H; L_H(h') \geq L_H(h)$, i.e., the ERM for h is the least.
- Therefore, the problem now becomes choosing a proper set H , and proving that overfitting cannot occur for any element.
- Also understand that reducing the size of H protects us against overfitting at the cost of learning becoming stiff. (Trade-off!)
- The simplest restriction on the set H would be to make it finite. We will now show that ERM_H will not overfit if H is finite, provided that the training sample is sufficiently large. We make the following reasonable assumption to aid us.

* Realizability Assumption :-

$$\forall H, \exists h^* \in H \text{ such that } L_{D_f}(h^*) = 0.$$

In words, there exists an estimator for every set whose true risk is zero. Note that this in turn implies that the empirical risk is zero as well, for every ERM.

* The i.i.d assumption :-

All points in S are sampled independently from the true distribution, D .

- Because the sample S is a multivariate R.V, by extension, our ERM's estimator would be a random variable as well, which means that the true risk is R.V!

- As the true risk is a RV, we can only talk in terms of probabilities now. We define a few terms-

1) Confidence Parameter - $(1-\delta)$

If δ is the probability of obtaining a bad sample, $(1-\delta)$ is the conf. parameter

2) Accuracy parameter - ϵ

h_s with $L_{D,f}(h_s) > \epsilon$ is said to be a failure, whereas $\leq \epsilon$ is said to be approximately correct.

- We would like the upper bound of $P\left[\{S_x : L_{D,f}(h_s) > \epsilon\}\right]$ to be low.
prob. of getting bad S .

Let the set of bad hypotheses be given by $H_B = \{h \in H : L_{D,f}(h) > \epsilon\}$ empirical = 0
overfitting!
additionally define misleading samples, $M = \{S : \exists h \in H_B, L_S(h) = 0\}$
 $= \bigcup_{h \in H_B} \{S : L_S(h) = 0\}$

We wanted to find $P\left[\{S : L_{D,f}(h_s) > \epsilon\}\right]$, from realization assumption

we have that $L_s(h_s) = 0$, $L_{D,f}(h_s) > \epsilon$

$$\Rightarrow \{S_x : L_{D,f}(h_s) > \epsilon\} \subseteq M$$

$$P\left[\{S_x : L_{D,f}(h_s) > \epsilon\}\right] \leq P\left[\bigcup_{h \in H_B} \{S : L_S(h) = 0\}\right]$$

$$P(A \cup B) \leq P(A) + P(B)$$

$$\leq \sum_{h \in H_B} P\left[\{S : L_S(h) = 0\}\right]$$

Simplifying RHS :-

$$\begin{aligned} \sum_{h \in H_B} P[\{S : L_{S,h} = 0\}] &= \sum_{h \in H_B} P[\{S : \forall x \in S, h(x) = f(x)\}] \\ &= \sum_{h \in H_B} \prod_{x \in S} P[h(x) = f(x)] = \sum_{h \in H_B} \prod_{x \in S} 1 - L_{D,f}(h) \leq \sum_{h \in H_B} (1-\epsilon) \\ &\leq (1-\epsilon)^m \end{aligned}$$

Now, use that $(1-\epsilon) \leq e^{-\epsilon}$

$$\Rightarrow \sum_{h \in H_B} (1-\epsilon)^m \leq \sum_{h \in H_B} e^{-m\epsilon} = |H_B| e^{-m\epsilon} \leq |H| e^{-m\epsilon}$$

Result The probability of getting a bad sample is capped by $|H| e^{-m\epsilon}$ ($m = |S|$)

$$\Rightarrow P[\{S : L_{D,f}(h_S) > \epsilon\}] \leq |H| (1-\epsilon)^m \leq |H| e^{-m\epsilon}$$

Corollary Let $\delta \in (0, 1)$ and $\epsilon > 0$. Let 'm' be an integer which satisfies :-

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

Then, for any labelling function 'f' and any distribution 'D', where the realizability assumption holds, for a sample S of size m, we have that

$$L_{D,f}(h_S) \leq \epsilon \quad \text{with probability atleast } (1-\delta)$$

That is, the minimum size of the sample that is to be taken to have a confidence of $(1-\delta)$ with error ϵ is now known!