

PAC learning framework: basic definitions.

①

We introduce first the PAC learning framework. PAC means "Probably Approximately correct".

Let X = domain set of "feature vectors".

(also called instances/patterns) $\underline{x} \in X$

Y = "label set" $y \in Y$

$\left\{ \begin{array}{l} \text{for classification probl } Y = \{0, 1\} \\ \text{for regression } Y = \mathbb{R} \\ \text{etc ...} \end{array} \right.$

$Z = X \times Y$ $z = (x, y)$ a "sample"

$\left\{ \begin{array}{l} \underline{x} \text{ is often seen as an "input" vector/pattern} \\ \underline{y} \text{ is often seen as an "output" label/value.} \end{array} \right.$

$S = \{ (\underline{x}_1, y_1), \dots, (\underline{x}_m, y_m) \}$
 $= \{ z_1, \dots, z_m \}$
 $=$ "Training set" (might appear with repetitions)

The training set will be the input to learning algorithms. We must model where do the samples come from?

$$\mathcal{D}(x, y) = \mathcal{D}(y | x) \mathcal{D}(x) .$$

conditional marginal .

We assume there exists some underlying distribution \mathcal{D} such that $z_i \sim \mathcal{D}$ (say iid for simplicity) $i = 1, \dots, m$. This distribution is assumed to exist in "Nature" or in the "world" producing the samples but is generally unknown to the learner/us.

[Note: we could try to determine some estimate of \mathcal{D} but this is highly sub-optimal a priori. Too many samples might be needed; one might be interested in some simple function of the samples. For example one might be interested in a classification rule; given a new x^{new} determine the label $y^{\text{new}} \in \{0, 1\}$.]

Learning Task :

Given a training sequence S we want to have an algorithm A that outputs an "hypothesis" or a "rule" $h = A(S)$, here an "hypothesis" or "rule" is a function :

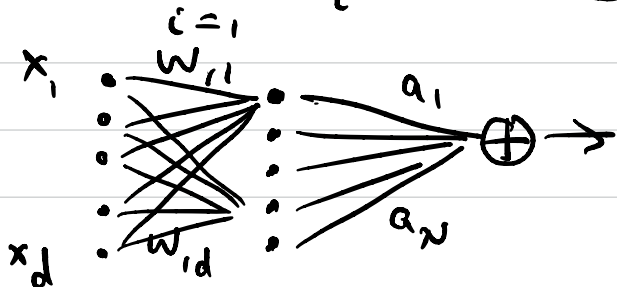
$$h : X \rightarrow Y$$

from feature vectors to label set.

- In classification $h(x) \in \{0, 1\}$ say.
- In regression $h(x) \in \mathbb{R}$ say.

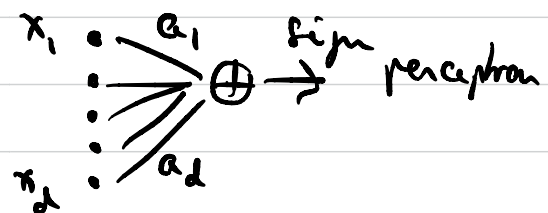
$$h(x) = \underline{a}^T \cdot \underline{x} \quad , \quad h(x) = \text{sign}(\underline{a}^T \cdot \underline{x})$$

$$h(x) = \sum_{i=1}^N a_i \sigma((Wx)_i)$$



two layer network.

ect ...



(4)

One often distinguishes two versions of the learning task:

Proper learning: \mathcal{H} is a class of hypothesis given a priori. You ask which $h \in \mathcal{H}$ is the "best" in a sense to be specified (see later on).

Note: \mathcal{H} plays the role of some prior knowledge / we might have specific constraints for possible hypothesis.

Improper learning: We still have some class of hypothesis \mathcal{H} but the "best" h (picked by the algorithm say) might be outside $h \notin \mathcal{H}$.

Now we discuss what we mean by the "best" h .

We usually have a "loss function":

$$l : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$$

$$(h, z) \mapsto l(h, z)$$

Popular loss fcts are:

$$l(h, z) = l(h, \underline{x}, y) = (h(\underline{x}) - y)^2$$

"square loss"

$$l(h, z) = l(h, \underline{x}, y) = \mathbb{1}(h(\underline{x}) \neq y)$$

$$= \begin{cases} 1 & \text{if } h(\underline{x}) \neq y \text{ unit cost} \\ 0 & \text{if } h(\underline{x}) = y \text{ no cost.} \end{cases}$$

"error fct or indicator loss fct".

Definition: True loss, true risk, true error, generalization error, ...

$$L_{\mathcal{D}}(h) \equiv \mathbb{E}_{z \sim \mathcal{D}} [l(h, z)]$$

6

This measures how well on average does an hypothesis do on samples $z \sim \mathcal{D}$.

We may define the "best" or "optimal" hypothesis as

$$h^* \in \operatorname{argmin} L_{\mathcal{D}}(h)$$

(or say $h^* = \operatorname{argmin} L_{\mathcal{D}}(h)$ for definiteness)

However there is a conceptual problem here because the real of the game is that \mathcal{D} is unknown. One possible remedy as we will be using is to minimize

a proxy:

Empirical Risk: $L_S(h) \equiv \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$

associated to a training set S .

ERM rule or algorithm: The empirical risk minimization rule consists to take

$$A(S) = h_S \in \operatorname{argmin} L_S(h).$$

Now we are ready to address the question of learnability of an hypothesis class. The philosophy is to take a "competitive point of view" where one tries to achieve something nearly as good as h^* . Note that the definition here is independent of a specific rule such as ERM or other rule.

Definition : PAC learning.

An hyp class \mathcal{H} is "agnostic PAC learnable" with respect to a set \mathcal{Z} and loss $l : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$ if \exists a function $m_{\mathcal{H}, l} : [0, 1]^2 \rightarrow \mathbb{N}$; $(\epsilon, \delta) \mapsto m_{\mathcal{H}, l}(\epsilon, \delta)$ and a learning rule or algo A such that :

• $\forall (\epsilon, \delta) \in [0, 1]^2$ and $\forall \mathcal{D}$ on \mathcal{Z} , when running

A on $m > m_{\mathcal{H}, l}(\epsilon, \delta)$ iid samples $S \sim \mathcal{D}^m$

$$\mathbb{P}_S \left\{ L_{\mathcal{D}}(A(S)) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon \right\} \geq 1 - \delta$$

[i.e we succeed with accuracy ϵ with prob $1 - \delta$]
 [i.e we succeed probably approximately correctly]

Remarks:

1) if $\mathcal{H} = \{h_0\}$ $\min_{h \in \mathcal{H}} L_{\infty}(h) = L_{\infty}(h_0)$ and the algorithm $A(S) = h_0$ that trivially outputs h_0 succeeds.

So the singleton class is trivially learnable.

if \mathcal{H} grows $\min_{h \in \mathcal{H}} L_{\infty}(h)$ might diminish and we might need more and more samples.

Determining if a class is learnable is non-trivial and depends on many things:

- size and quality of S
- size of \mathcal{H}
- loss ℓ
- types of A

all these objects interact and it is not easy to determine their interactions.

2) Perhaps the most popular learning algo is ERM.

$$A(S) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} L_S(h)$$

which minimizes empirical risk. How is this minimized or even better we really want to minimize or somehow approximately minimize is a modern question of ML.

3) Above definition of learnability does not address the question of complexity and computational resources.

This is an important issue even with ERM if \mathcal{H} is large and/or S is large.

4) Vocabulary: we say that we are in the Realizable case if $\exists h_* \in \mathcal{H}$ such that $L_{\mathcal{D}}(h_*) = 0$.

Then \mathcal{H} is learnable if $\Pr_S \{ L_{\mathcal{D}}(A(S)) \leq \epsilon \} \geq 1 - \delta$

Agnostic learnable means that we do not know if we are in the realizable case.

Where do we go from here?

→ we will first analyse \mathcal{H} finite. We will show that finite classes are learnable by ERM. We will use a notion of uniform convergence of $L_S(h)$ to $L_D(h)$ somehow and will obtain an estimate of $M_{\mathcal{H}}(\epsilon, D)$

→ in the exercises you will look at natural \mathcal{H} given by separating planes, circles, threshold sets and see that infinite \mathcal{H} 's can also be learnable. This will be understood in terms of the VC dimension which is an effective measure of the "size" of \mathcal{H} . We will see that classes with finite VC dim are learnable.

→ But it is not true that all infinite classes are learnable. We will prove a "No free lunch theorem" which roughly states

that there is no universal A : i.e. given A one can find D s.t. A will fail. (if $|X|$ is big enough).

