

3.1 PAC Learning

- In the discussion in Statistical Learning Framework, we learned that for a finite hypothesis class, if the ERM rule with respect to that class is applied on a sufficiently large training sample, then the output hypothesis will be probably approximately correct.
- That is the Probably Approximately Correct (PAC) learning

Definition 3.1 (PAC Learnability)

- A hypothesis class H is PAC learnable if there exist:
 - function $m_H: (0,1)^2 \rightarrow \mathbb{N}$
 - learning algorithm with the following property:

For every $\epsilon, \delta \in (0,1)$
distribution D over X
labeling function $f: X \rightarrow \{0,1\}$,
- If the realizability assumption holds with respect to H, D, f , then when running the learning algorithm on $m \geq m_H(\epsilon, \delta)$ i.i.d. examples generated by D and labeled by f , the algorithm returns a hypothesis h such that, with probability of at least $1 - \delta$ (over the choice of the examples), $L_{D,f}(h) \leq \epsilon$

→ What do ϵ and δ stand for?

These are 2 approximation parameters:

- ϵ : accuracy parameter, determines how far the output classifier can be from the optimal one (approximately - correct)
- δ : confidence parameter, indicating how likely the classifier is to meet that accuracy requirement ("probably" in PAC)

→ continue next page

Why use these approximations (ϵ and δ)?

- Since the training set is randomly generated, there is always a chance that the set is noninformative (not giving information about distribution D). For example: it is possible to pick just 1 point in domain X , and sample it over and over again.
- Even if we get a training set S that represents distribution D , since S is finite, there is still some fine details of D that S fails to capture. That's where our accuracy parameter, ϵ , allows "forgiving" the learner's classifier for making minor mistakes.

Sample Complexity

How many examples in training set S are required to guarantee a probably approximately correct solution?

The function $m_H : (0, 1)^2 \rightarrow \mathbb{N}$ determines the sample complexity of learning H .

Corollary 3.2 (Sample complexity)

Every finite hypothesis class is PAC learnable with sample complexity

$$m_H(\epsilon, \delta) \leq \left\lceil \frac{\log(|H|/\delta)}{\epsilon} \right\rceil$$

Note: m (the size) is different from $m_H(\epsilon, \delta)$ (the threshold)

- Recall from corollary 2.3 of Statistical Learning Framework, we came to the conclusion that the size of training set should be:

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

- m_H is defined as the "minimal function" in corollary 3.2. That is, the value of m_H is the smallest possible value that satisfy all the PAC learnability with parameters ϵ and δ .

3.2 A More General Learning Model

A more generalized model can be achieved by considering 2 aspects:

- Removing the Realizability Assumption
- Learning Problems beyond Binary Classification

Releasing the Realizability Assumption - Agnostic PAC Learning

A more realistic model for the Data-Generating Distribution

- Recall that the realizability assumption said:

$$\exists h^* \in H \text{ s.t. } P_{x \sim D} [h^*(x) = f(x)] = 1$$

In practical, this is not always true. For our papayas example, we cannot guarantee that 2 papayas with the same color and softness are both tasty.

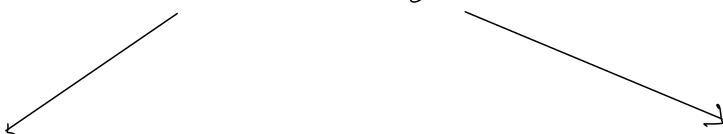
- Instead, we relax the assumption with a more flexible notion, a data-labels generating distribution, described as follows:

- D is a joint distribution over $\underbrace{X \times Y}$, denoted $D(x, y)$

domain set	label set
------------	-----------

- D can be considered to have 2 parts, by chain rule of joint prob:

$$D(x, y) = D(x) \cdot D(y | x)$$


marginal distribution D_x
over unlabeled domain points

conditional distribution $D(y | x)$
over labels for each domain point

The empirical and true error revised

Redefined true error

For distribution D over $X \times Y$, prediction rule h :

$$L_D(h) = \mathbb{P}_{(x,y) \sim D} [h(x) \neq y] = D(\{ (x,y) : h(x) \neq y \})$$

→ What this means in human language?

Probability that a data point-label pair (x, y) such that the prediction made by hypothesis h is incorrect.

→ What is the difference with realizability is assumed?

Things that contribute to true error

With Realizability Assumption

- Representativeness of S

Without Realizability Assumption

- Representativeness of S
- Chance of mislabeling S
- Choice of hypothesis class H

Training error is the same

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

→ Why is it the same regardless of realizability assumption.

- Recall that the realizability assumption assumes there is a "perfect" classifier that will make $L_{D,f}(h) = 0$
- Training error measure performance on training set S , so the existence of a "perfect" classifier should not matter.

$$L_S(h) = L_{D(\text{uniform over } S)}(h)$$

- Training loss equals the expected loss under distribution D where each points in S contribute equally to the error.

- Assume "uniform over S " essentially means consider each data points in S will contribute equally to the training error.
- $L_D(\text{Uniform over } S)(h)$ means ?
- True risk when estimated using sample S , where each data points are weighted equally
- This is common in ERM, where D is unknown, and the empirical distribution is used as an approximation

The goal

We wish to find some hypothesis, $h: X \rightarrow Y$, that minimizes the true risk, $L_D(h)$

The Bayes Optimal Predictor

Given any probability distribution D over $X \times \{0, 1\}$, the best predicting function $f: X \rightarrow \{0, 1\}$ is defined as:

$$f_D(x) = \begin{cases} 1 & \text{if } P(y=1|x) \geq \frac{1}{2} \\ 0 & \end{cases}$$

Proof (Exercise 7)

- Let $g: X \rightarrow \{0, 1\}$ be any other classifier.
We want to prove that with any g , $L_D(f_D) \leq L_D(g_D)$
- For any classifier g , we can prove that:

$$E_{X,Y}(f(x, y)) = E_x E_{Y|X=x}(f(x, y) | X=x)$$
 Q.

$$\quad \quad \quad \langle \text{Law of total expectation, Adam's law} \rangle$$

continue

$$\begin{aligned}
 \text{Therefore: } L_D(g) &= \mathbb{E}_{(x,y) \sim D} (1_{g(x) \neq y}) &< \text{definition} \\
 &= \mathbb{E}_{x \sim D_x} \left[\mathbb{E}_{y \sim D_{Y|x}} (1_{g(x) \neq y} \mid X=x) \right] &< \text{proven above} \\
 &= \mathbb{E}_{x \sim D_x} \left[P(\{g(x) \neq y \mid X=x\}) \right]
 \end{aligned}$$

We want to prove that, for all $x \in X$, we have:

$$P(\{g(x) \neq y \mid X=x\}) \geq P(\{f_D(x) \neq y \mid X=x\})$$

$$P(\{g(x) \neq y \mid X=x\})$$

$$= 1 - P(\{g(x) = y \mid X=x\}) \quad (1)$$

$$= 1 - P(\{g(x) = 1, Y=1 \mid X=x\}) - P(\{g(x) = 0, Y=0 \mid X=x\})$$

LOT P

$$P(\{g(x) = 1, Y=1 \mid X=x\})$$

$$= P(\{g(x) = 1 \mid X=x\}) \cdot P(\{Y=1 \mid X=x\})$$

equals $1_{g(x)=1}$, because

- When $g(x) = 1$: $P(\{g(x) = 1 \mid X=x\}) = 1$
- When $g(x) = 0$: $P(\{g(x) = 1 \mid X=x\}) = 0$

equals $1_{g(x)=0}$, because:

- When $g(x) = 1$: $P(\{g(x) = 0 \mid X=x\}) = 0$
- When $g(x) = 0$: $P(\{g(x) = 0 \mid X=x\}) = 1$

So (1) can be rewritten as:

$$\begin{aligned}
 P(\{g(x) = y \mid X=x\}) &= 1_{g(x)=1} P(\{Y=1 \mid X=x\}) + 1_{g(x)=0} P(\{Y=0 \mid X=x\}) \quad (2)
 \end{aligned}$$

Apply similar logic to $P(\{f_D(x) = y \mid X=x\})$, we have:

$$P(\{f_D(x) = y \mid X=x\}) \quad (3)$$

$$= 1_{f_D(x)=1} P(\{Y=1 \mid X=x\}) + 1_{f_D(x)=0} P(\{Y=0 \mid X=x\})$$

o Let (3) - (2), we have:

$$\begin{aligned}
 & P(\{ f_D(x) = y \mid X=x \}) - P(\{ g(x) = y \mid X=x \}) \\
 &= P(\{ Y=1 \mid X=x \}) (1_{f_D(x)=1} - 1_{g(x)=1}) + P(\{ Y=0 \mid X=x \}) (1_{f_D(x)=0} - 1_{g(x)=0}) \\
 &\quad \left. \begin{array}{l} \text{Since:} \\ \cdot P(\{ Y=0 \mid X=x \}) = 1 - P(\{ Y=1 \mid X=x \}) \\ \cdot 1_{f_D(x)=0} = 1 - 1_{f_D(x)=1} \\ \cdot 1_{g(x)=0} = 1 - 1_{g(x)=1} \end{array} \right\} \\
 &= [2P(\{ Y=1 \mid X=x \}) - 1] (1_{f_D(x)=1} - 1_{g(x)=1}) \geq 0
 \end{aligned}$$

Therefore:

$$\begin{aligned}
 & P(\{ f_D(x) = y \mid X=x \}) \geq P(\{ g(x) = y \mid X=x \}) \\
 \Leftrightarrow & P(\{ f_D(x) \neq y \mid X=x \}) \leq P(\{ g(x) \neq y \mid X=x \}) \\
 \Leftrightarrow & E_{x \sim D_X} [P(\{ f_D(x) \neq y \mid X=x \})] \leq E_{x \sim D_X} [P(\{ g(x) \neq y \mid X=x \})] \\
 \Leftrightarrow & E_{x \sim D_X} \left[E_{y \sim D_Y|x} (1_{f_D(x) \neq y} \mid X=x) \right] \leq E_{x \sim D_X} \left[E_{y \sim D_Y|x} (1_{g(x) \neq y} \mid X=x) \right] \\
 \Leftrightarrow & E_{(x,y) \sim D} \left[1_{f_D(x) \neq y} \right] \leq E_{(x,y) \sim D} \left[1_{g(x) \neq y} \right] \\
 \Leftrightarrow & L_D(f_D) \leq L_D(g)
 \end{aligned}$$

→ If Bayes predictor is so good, why not just use it?

Because we don't know distribution D , so we try to find a predictor whose is as close to Bayes predictor as possible

Definition 3.3 (Agnostic PAC Learnability)

A hypothesis class \mathcal{H} is agnostic PAC learnable if there exist:

- A function $m_{\mathcal{H}}: (0,1)^2 \rightarrow \mathbb{N}$

- A learning algorithm with the following property:

For every $\epsilon, \delta \in (0,1)$ and distribution D over $X \times Y$,

when running the learning algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ i.i.d examples generated by D , the algorithm returns a hypothesis h such that, with probability of at least $(1 - \delta)$, over the choice of the m training examples:

$$L_D(h) \leq \min_{h' \in \mathcal{H}} L_D(h') + \epsilon$$

How is this different from standard PAC Learning?

Standard PAC

. Best $h \in \mathcal{H}$ has true error = 0

$$L_D(h) = 0$$

. The learning goal is to find $h \in \mathcal{H}$, with high probability $1 - \delta$, such that true error is close to zero, by a margin of ϵ

$$L_D(h) \leq \epsilon$$

Agnostic PAC

. Best $h' \in \mathcal{H}$ has true error > 0

$$\min_{h' \in \mathcal{H}} L_D(h') > 0$$

. The learning goal is to find $h \in \mathcal{H}$, with high probability $1 - \delta$, such that true error not much worse than the best possible $h' \in \mathcal{H}$, so:

$$L_D(h) \leq \min_{h' \in \mathcal{H}} L_D(h') + \epsilon$$

3.2.2 The Scope of Learning Problems Modeled

Let's extend our model so that it can be applied to a wide variety of learning tasks, some of these tasks are:

Multiclass Classification:

Instead of having binary classification, the label set Y will be the set of categories that is large and finite.

Regression:

In this task, we try to find some simple pattern between domain set X and label set Y .

To accommodate wide range of learning tasks, we generalize our measure of success as follows:

Generalized Loss Functions

Given } any set H (that plays the role of our hypotheses)
some domain Z

Let l be any function from $H \times Z$ to the set of nonnegative real numbers,

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

Note:

- For prediction task, our $Z = X \times Y$
- But here we generalize beyond prediction task, and therefore allow Z to be any domain of examples.

Risk function

Is the expected loss of a classifier, $h \in H$, with respect to a probability distribution D over Z :

$$L_D(h) = \mathbb{E}_{z \sim D} [l(h, z)]$$

Empirical Risk

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z_i)$$

Some common loss functions, $l(h, z)$

- 0-1 loss: used in binary and multi class classification problems:

$$l_{0-1}(h, (x, y)) = \begin{cases} 0 & \text{if } h(x) = y \\ 1 & \text{if } h(x) \neq y \end{cases}$$

Consequently :

$$\begin{aligned} L_0(h) &= \underset{(x,y) \sim D}{E} [l(h, (x, y))] \\ &= \underset{(x,y) \sim D}{P} [h(x) \neq y] \\ &= D(\{ (x, y) : h(x) \neq y \}) \end{aligned}$$

fundamental bridge:

$$P(A) = E(I_A)$$

- Square loss: used in regression problems

$$l_{\text{sq}}(h, (x, y)) = (h(x) - y)^2$$

So, we have completed defining agnostic PAC Learning with general loss function, that is:

Definition 3.4 (Agnostic PAC Learnability for General Loss Functions)

A hypothesis class \mathcal{H} is agnostic PAC learnable with respect to a set Z and a loss function $l: \mathcal{H} \times Z \rightarrow \mathbb{R}_+$, if there exist a function $m_{\mathcal{H}}: (0,1)^2 \rightarrow \mathbb{N}$ and a learning algorithm with the following property:

For every $\epsilon, \delta \in (0,1)$
 distribution D over Z ,

when running the algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ i.i.d. examples generated by D , the algorithm returns $h \in \mathcal{H}$ such that, with probability of at least $1 - \delta$ (over the choice of the m training examples)

$$L_D(h) \leq \min_{h' \in \mathcal{H}} L_D(h') + \epsilon$$

$$\text{where } L_D(h) = \mathbb{E}_{z \sim D} [l(h, z)]$$