The Probably Approximately Correct (PAC) Learning

Prof. Dan A. Simovici

UMB

1/22

Outline

1 The Formal Model

2 The Agnostic PAC Learning

3 The Scope of Learning Problems

Problem framework: Suppose that you wish to determine if a fruit is good for eating and you base this decision of its color and softness.

Also, suppose that you tasted a sample of fruits and you got the following results:

F#	color	softness	good to eat
1	green	hard	0
2	dark green	mushy	1
3	green	soft	1
4	yellow	hard	0
5	orange	soft	1

A learning algorithm \mathcal{A} starts with a hypothesis class \mathcal{H} and a sample S, under certain conditions, it returns a hypothesis h that has a small true error. A specific definition is given next. We assume that the data set \mathcal{X} is equipped with a probability distribution \mathcal{D} .

What is the PAC Model?

Definition

A hypothesis class \mathcal{H} is PAC learnable if there exists a function $m_{\mathcal{H}}: (0,1)^2 \longrightarrow \mathbb{N}$ and a learning algorithm \mathcal{A} such that for every $\epsilon, \delta \in (0,1)$, every distribution \mathcal{D} over \mathcal{X} , and for every labeling function $f: \mathcal{X} \longrightarrow \{0,1\}$, if realizability assumption holds with respect to $\mathcal{H}, \mathcal{D}, f$, then when running the algorithm \mathcal{A} on a sample S that consists of $m \ge m_{\mathcal{H}}(\epsilon, \delta)$ generated by \mathcal{D} and labeled by f, \mathcal{A} returns a hypothesis h such that, with probability at least $1 - \delta$ (over the choice of examples), we have for the true error $L_{(\mathcal{D},f)}(h)$:

$$P(L_{(\mathcal{D},f)}(h) \leq \epsilon) \geq 1 - \delta.$$

 ϵ is the accuracy parameter and δ is the confidence parameter

Approximation Parameters

- the accuracy parameter ϵ determines how far the output classifier can be from the optimal one, and
- the confidence parameter δ indicates how likely is the classifier is to meet that accuracy requirement.

What is Agnostic PAC Learning?

- The realizability assumption, the existence of a hypothesis h^{*} ∈ H such that P_{x∼D}(h^{*}(x) = f(x)) = 1 is not practical in many cases.
- Agnostic learning replaces the realizability assumption and the targeted labeling function *f*, with a distribution *D* defined on pairs (data, labels), that is, with a distribution *D* on *X* × *Y*.

When the probability distribution D was defined on X, the generalization error of a hypothesis was defined as:

$$L_{\mathcal{D},f}(h) = \mathcal{D}(\{x \mid h(x) \neq f(x)\}).$$

■ For agnostic learning *D* is defined over *X* × *Y*, so we redefine the generalization error as:

$$L_{\mathcal{D}}(h) = \mathcal{D}(\{(x, y) \mid h(x) \neq y\}).$$

We seek a predictor for which $L_D(h)$ is minimal.

• The definition of the empirical risk remains the same:

$$L_{\mathcal{S}}(h) = \frac{|\{i \mid h(x_i) \neq y_i \text{ for } 1 \leq i \leq m\}|}{m}$$

8 / 22

The Bayes Classifier and Its Optimality

Let \mathcal{D} be any probability distribution over $\mathcal{X} \times \mathcal{Y}$, where $\mathcal{Y} = \{0, 1\}$. Let X be a random variable ranging over \mathcal{X} and Y be a random variable ranging over $\mathcal{Y} = \{0, 1\}$. The Bayes predictor is the function $f_{\mathcal{D}}$ defined as

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \text{if } P(Y=1|X=x) \geqslant rac{1}{2} \\ 0 & \text{otherwise.} \end{cases}$$

Theorem

Given any probability distribution \mathcal{D} over $\mathcal{X} \times \{0,1\}$ the best label predicting function $f : \mathcal{X} \longrightarrow \{0,1\}$ is the Bayes predictor $f_{\mathcal{D}}$.

In other words, we need to prove that for hypothesis $g : \mathcal{X} \longrightarrow \{0, 1\}$ we have $L_{\mathcal{D}}(f_{\mathcal{D}}) \leq L_{\mathcal{D}}(g)$.

Proof

Let X be a random variable ranging over \mathcal{X} , Y be a random variable ranging over $\mathcal{Y} = \{0, 1\}$, and let α_x be the probability of a having a label 1 given x, that is, $\alpha_x = P(Y = 1 | X = x)$. With this notation the Bayes predictor is

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \text{if } \alpha_x \ge \frac{1}{2} \\ 0 & \text{if } \alpha_x < \frac{1}{2}. \end{cases}$$

11/22

Proof (cont'd)

We have:

$$L_{\mathcal{D}}(f_{\mathcal{D}}) = P(f_{\mathcal{D}}(X) \neq y | X = x)$$

= $P(f_{\mathcal{D}}(x) = 1 | X = x) P(Y = 0 | X = x)$
 $+ P(f_{\mathcal{D}}(x) = 0 | X = x) P(Y = 1 | X = x)$
= $P\left(\alpha_x \ge \frac{1}{2}\right) P(Y = 0 | X = x)$
 $+ P\left(\alpha_x < \frac{1}{2}\right) P(Y = 1 | X = x)$

Note: when we write $P(f_{\mathcal{D}}(X) \neq y | X = x)$ we mean the probability that a pair (x, y) is such that $f_{\mathcal{D}}(X) \neq y$ assuming that X = x. Similar conventions apply to all probabilities listed above.

イロト (四) (日) (日) (日) (日) (日)

Proof (cont'd)

If $\alpha_x \ge \frac{1}{2}$, then

$$\min\{\alpha_x, 1-\alpha_x\} = 1-\alpha_x, P\left(\alpha_x \ge \frac{1}{2}\right) = 1, P\left(\alpha_x < \frac{1}{2}\right) = 0,$$

and

$$P\left(\alpha_{x} \geq \frac{1}{2}\right) (1 - \alpha_{x}) + P\left(\alpha_{x} < \frac{1}{2}\right) \alpha_{x}$$
$$= 1 - \alpha_{x} = \min\{1 - \alpha_{x}, \alpha_{x}\}.$$

If $\alpha_x < \frac{1}{2}$, then $\min\{\alpha_x, 1 - \alpha_x\} = \alpha_x$, $P\left(\alpha_x \ge \frac{1}{2}\right) = 0$, $P\left(\alpha_x < \frac{1}{2}\right) = 1$ and

$$P\left(\alpha_{x} \ge \frac{1}{2}\right) (1 - \alpha_{x}) + P\left(\alpha_{x} < \frac{1}{2}\right) \alpha_{x}$$
$$= \alpha_{x} = \min\{1 - \alpha_{x}, \alpha_{x}\}.$$

Proof (cont'd)

Let g be any other classifier. We have:

$$P(g(X) \neq Y | X = x) = P(g(X) = 0 | X = x) P(Y = 1 | X = x) + P(g(X) = 1 | X = x) P(Y = 0 | X = x) = P(g(X) = 0 | X = x) \alpha_x + P(g(X) = 1 | X = x)(1 - \alpha_x) \ge P(g(X) = 0 | X = x) \min\{\alpha_x, 1 - \alpha_x\} + P(g(X) = 1 | X = x) \min\{\alpha_x, 1 - \alpha_x\} \ge (P(g(X) = 0 | X = x) + P(g(X) = 1 | X = x)) \cdot \min\{\alpha_x, 1 - \alpha_x\} = \min\{\alpha_x, 1 - \alpha_x\} = P(f_{\mathcal{D}}(X) \neq y | X = x).$$

Agnostic PAC-Learnability

Definition

A hypothesis class \mathcal{H} is agnostic PAC learnable if there exists a function $m_{\mathcal{H}} : (0,1)^2 \longrightarrow \mathbb{N}$ and a learning algorithm \mathcal{A} with the following property:

For every

- $\epsilon, \delta \in (0, 1)$ and
- for every distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$,

when running \mathcal{A} on $m \ge m_{\mathcal{H}}(\epsilon, \delta)$ iid examples generated by \mathcal{D} , \mathcal{A} returns a hypothesis h such that with probability at least $1 - \delta$ (over the choice of the m training examples) we have

$$L_{\mathcal{D}}(h) \leqslant \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon.$$

- If the realizability assumption holds, agnostic PAC learning provides the same guarantees as PAC learning.
- When the realizability assumption does not hold, no learner can guarantee an arbitrary small error.
- A learner A can declare success if the error is not much larger than the smallest error achievable by a hypothesis from H.

The Probably Approximately Correct (PAC) Learning

└─ The Scope of Learning Problems

Multiclass Classification

Example

Let \mathcal{X} be a set of document features, and \mathcal{Y} a set of topics (sports, politics, health, etc.)

(sports, politics, health, etc.).

By document features we mean counts of certain key words, size, or origin of the document.

The loss function will be the probability of the event that occurs when the predictor suggest a wrong label.

The Probably Approximately Correct (PAC) Learning

└─ The Scope of Learning Problems

Regression

Example

In *regression* we seek to find a functional relationship h between the \mathcal{X} and \mathcal{Y} components of the data.

For example, to predict the weight of a baby at birth ${\mathcal X}$ can be a set of triplets in ${\mathbb R}^3$

(head circumference, abdominal circumference, femur length) and \mathcal{Y} is is the weight at birth. We seek *h* that will minimize the loss $L_{\mathcal{D}}(h) = E_{(x,y)\sim\mathcal{D}}(h(x) - y)^2$. └─ The Scope of Learning Problems

Generalized Loss Functions

Definition

Given a set of hypotheses \mathcal{H} , a domain Z, a loss function is a function $\ell : \mathcal{H} \times Z \longrightarrow \mathbb{R}_+$.

For prediction problems we have $Z = \mathcal{X} \times \mathcal{Y}$.

Definition

The risk function is the expected loss of the classifier $h \in \mathcal{H}$ with respect to a probability distribution \mathcal{D} over Z, namely

$$L_{\mathcal{D}}(h) = E_{z \sim \mathcal{D}}(\ell(h, z)).$$

The empirical risk is the expected loss over the sample $S = (z_1, \ldots, s_m) \in Z^m$ as

$$L_{S}(h) = \frac{1}{2} \sum_{i=1}^{m} \ell(h, z_{i}).$$

└─ The Scope of Learning Problems

0-1 Loss

The random variable z ranges over $\mathcal{X} \times \mathcal{Y}$ and the loss function is

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

This is used in binary or multiclass classification problems. For the 0/1 loss the definition of $L_D(h) = E_{z \sim D}(\ell(h, z))$ coincides with the previous definition in the agnostic PAC, $L_D(h) = D(\{(x, y) \mid h(x) \neq y\}).$ The Probably Approximately Correct (PAC) Learning

└─ The Scope of Learning Problems



The random variable z ranges over $\mathcal{X} \times \mathcal{Y}$ and the loss function is

$$\ell_{sq}(h,(x,y)) = (h(x) - y)^2.$$

└─ The Scope of Learning Problems

Agnostic PAC Learnability for General Loss Functions

Definition

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

For every $\epsilon, \delta \in (0, 1)$ and for every distribution \mathcal{D} over Z, when running \mathcal{A} on $m \ge m_{\mathcal{H}}(\epsilon, \delta)$ iid examples generated by \mathcal{D} , \mathcal{A} returns a hypothesis h such that with probability at least $1 - \delta$ (over the choice of the m training examples) we have

$$L_{\mathcal{D}}(h) \leqslant \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon,$$

where $L_{\mathcal{D}}(h) = E_{z \sim \mathcal{D}}(\ell(h, z))$.

イロト イボト イヨト イヨト