

Statistical learning theory

(1)

Setup:

- Sample space X , (also called instance space).
Its points are called samples. (or instances)
Usually we have $X \subseteq \mathbb{R}^n$, where each component
represents a feature.
i.e. each sample is represented by a vector of
features.
- Label set y .
For multiclassification y is also finite.
For regression $y = \mathbb{R}$.
- Training data: ~~finite set~~
With $t = X \times Y$. Then training data will be
represented by ~~a some finite set~~ $S \subseteq \mathbb{Z}^m$.
- Thus the training data is a ph. $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$
consisting of labeled samples.

Prediction function

(2)

- A learner is requested to output
- the prediction function $h: X \rightarrow Y$, predicting the relationship between samples and labels -
For example, if there is a function relationship
 $f: X \rightarrow Y$, then the learner is asked to approximate this function.

Assumptions on data

- The training data is not completely random.
But rather, it is sampled from an underlying distribution

That is: we have a Random variable

$$(X, Y) : \Omega \rightarrow X \times Y \subset \mathbb{R}^n \times Y.$$

which induces a ~~not~~ ~~measure on~~ $X \times Y$ push-forward measure on $X \times Y$.

So, we have a joint distribution on samples and labels. called D

(3)

Intrinsically this joint distribution decomposes into

- - marginal distribution
- conditional distribution.

~~but~~

Assuming ~~disjoint~~ we have a joint density f_{xy}
we obtain the marginal density

$$f_x(x) = \int f_{x,y}(x,y) dy$$

and define $f_{y|x}(y) = \frac{f_{x,y}(x,y)}{f_x(x)}$

Note: In general we do not have access to
this distribution.

Interpretation/Example:

- The marginal distribution describes the distribution of the samples as they are not completely random
- The conditional distribution, describes the probability given a sample x , which values ~~which~~ it might have.

(4)

Say we want to classify Papayas into

- 1 = "tasty"
- 0 = "non-tasty"

We represent each papaya by the features

- color
- softness

Thus $X \subseteq \mathbb{R}^2$ is the space of all Papayas represented by color and softness.

The clarity: X does not contain arbitrary points as there are limits to color and softness.

The marginal distribution describes how likely it is to have Papayas with given color and softness.

The conditional distribution describes given a particular Papaya, how likely is it that it's tasty.

So the joint distribution describes the probability that a tasty papaya w/ given color and softness occurs.

(5)

In practice:

We expect outliers to occur less often!

• Error functions / loss functions.

loss functions:

The other's class:

Fix a set H of prediction functions.

This could be, all functions

• linear functions

Error functions

measure space

Given H and, or \mathcal{Z} , a loss function

• loss function is a function

$$\ell: H \times \mathcal{Z} \rightarrow \mathbb{R}_{\geq 0}$$

s.t. $\ell(\cdot, -): \mathcal{Z} \rightarrow \mathbb{R}_{\geq 0}$ is measurable

In our supervised learning setting we have

$f = h \times g$, with but this formulation includes
unsupervised learning.

The risk function / ^{the more} is the expected loss of
a prediction function with D over \mathcal{F} .

$$\text{i.e. } L_D(h) : \mathcal{H} \rightarrow \mathbb{R}^+$$

$$h \mapsto \mathbb{E} [\underset{\text{avg}}{\ell}(h, z)]$$

The empirical risk / ^{empirical} is the expected loss on
a given sample $S \in \mathcal{Z}^m$

$$L_S : \mathcal{H} \rightarrow \mathbb{R}^+$$

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

Example: 0-1 loss in classification:

$$\ell_{0-1}(h, (x, y)) = \begin{cases} 0 & h(x) = y \\ 1 & h(x) \neq y \end{cases} \Rightarrow \begin{aligned} L_0(h) &= P[h(x) \neq y] \\ &\stackrel{(x, y) \sim D}{=} D[h(x) \neq y] \end{aligned}$$

Square loss

$$\ell_2(h, (x, y)) = (h(x) - y)^2$$

$$L_2(h) = \int (h(x) - y)^2 dx dy$$

(2)

i.i.d.:

Remember that our training data is given by
 a pt $S \in \mathbb{Z}^m = (\mathbf{x}, y)^m$

We assume that the samples are all independent.

Notation $S \sim P^m$.

Empirical risk minimization

Remember that

Given a training set $S \in \mathbb{Z}^m$ sampled from an unknown \mathcal{D}

We want to output a predictor $h: \mathbf{x} \rightarrow y$

which depends on S .

Since \mathcal{D} is unknown we only have access
 to the empirical risk.

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

Note: That is simply the training error!

Clearly, minimizing wrt to the empirical risk when $H = \text{all}$
will lead to overfitting. ⑥

We want to search for conditions, which guarantee
that ERM does not overfit. That is also likely
to perform well on unseen data.

The obvious problem is that H ~~can~~ considers all
my functions at this point.

The simplest restriction is that H is finite.

We show that for H finite we can guarantee
not to overfit while retaining good generalizable
performance.

(9)

To this end we will make the following:

Definition (Realizability Assumption)

There exists $h^* \in \mathcal{H}$ s.t. $L_D(h^*) = 0$

Note that this implies w/ probability 1 that
over a random sample S we have

$$L_S(h^*) = 0$$

We moreover assume that $L_D(h_S)$ is a random
variable in S .

Note:

Analysis of link \mathcal{H}

Given $h_S \in \arg \min_{h \in \mathcal{H}} L_D(h)$
give a bound on $L_D(h_S)$

- It is not realistic to expect w/ full certainty
that S will suffice to direct the learner
to a good predictor.

(10)

It could always be that the sample is non-representative of \mathcal{D} .

Thus we can only speak of the probability that $L_D(h_S)$ is not too large.

Def: The probability that of getting a non-rep. sample by S .

Call $(1-\delta)$ the confidence parameter.

Moreover, we cannot guarantee perfect label prediction.

Def: so we have an accuracy parameter $\varepsilon > 0$.

Def: If $\varepsilon > 0$ we interpret $L_D(h_S) > \varepsilon$ as failure to learn.

$L_D(h_S) \leq \varepsilon$ is an approximately correct prediction.

(11)

We want to estimate the size

of the event $L_D(h_S) > \varepsilon$ in S .

More precisely:

~~What's next?~~

We want to estimate

$$\mathbb{P}^n \left(\{S \in \mathcal{Z}^n \mid L_D(h_S) > \varepsilon\} \right) \quad \begin{matrix} \text{remember} \\ L_D(h_S) \text{ random variable} \end{matrix}$$

Let $\mathcal{H}_B \subset \mathcal{H}$ be the set of "bad" hypotheses

$$\text{i.e. } \mathcal{H}_B = \{h \in \mathcal{H} \mid L_D(h) > \varepsilon\}$$

$$\text{Set } M = \{S \in \mathcal{Z}^n \mid \exists h \in \mathcal{H}_B, L_S(h) = 0\}$$

The set of "bad" samples.

i.e. for $S \in M$, the bad hypothesis works good.

(R)

The realizability assumption implies

$$L_S(h_S) = 0.$$

$\Rightarrow L_D(h_S) > \varepsilon$ can only happen, if

for some $h \in H_B$ we have $L_S(h) = 0$.

($\Rightarrow S \subset M$)

We have

$$\{ s \in \mathbb{F}^m \mid L_D(h_s) > \varepsilon \} \subset M.$$

Now write

$$M = \bigcup_{h \in H_B} \{ s \in \mathbb{F}^m \mid L_S(h) = 0 \}.$$

Hence:

$$\mathbb{D}^m(L_D(h_S) > \varepsilon) \leq \mathbb{D}^m(M) = \mathbb{D}^m\left(\bigcup_{h \in H_B} \{ s \in \mathbb{F}^m \mid L_S(h) = 0 \}\right)$$

$$\leq \sum_{h \in H_B} \mathbb{D}^m(\{ s \in \mathbb{F}^m \mid L_S(h) = 0 \})$$

Now we bound each summand.

(13)

Fix $h \in H_B$. Then $L_S(h) = 0$ if and only if $h(x_i) = y_i \quad \forall i$.

Hence

$$\begin{aligned} & D^m(\{s \in \mathbb{F}^m \mid L_S(h)=0\}) \\ &= \bigcup_{i=1}^n (\{s \in \mathbb{F}^m \mid h(x_i) = y_i \quad \forall i\}) \\ &= \overline{\bigcup_{i=1}^n \{((x_i, y_i)) \mid h(x_i) = y_i\}} \end{aligned}$$

$$L_D(h) = D(\text{irrig}(\{h(x_i) \neq y_i\}) \text{ w.r.t. } h).$$

We have

$$D(\{((x_i, y_i)) \mid h(x_i) = y_i\}) = 1 - L_D(h) \leq 1 - \varepsilon$$

since $h \in H_B$ and $L_D(h) \geq \varepsilon$

Now use $1 - \varepsilon \leq e^{-\varepsilon}$ we have,

$$D^m(\{s \in \mathbb{F}^m \mid L_S(h)=0\}) \leq (1 - \varepsilon)^m \leq e^{-\varepsilon m}$$

end

14

We conclude that

$$D^m \left(\{s \in \mathcal{Z}^m \mid L_D(h_S) > \varepsilon \} \right) \leq |H_B| c^{-\varepsilon m} \leq |H| c^{-\varepsilon m}$$

We have proven:

Prop: Let H be finite. $\delta \in (0, 1)$ and $\varepsilon > 0$

$$\text{Let } n \in \mathbb{N} \text{ s.t. } n \geq \frac{\log(|H|/\delta)}{\varepsilon}$$

Then for any D for which realizability holds

with probability $1-\delta$ on the choice of

Sample S of size n we have

$$L_D(h_S) \leq \varepsilon$$