

... someone didn't show up again :(

## Learning Theory

We first go over some result in **sampling**. Suppose given an unknown distribution  $X$  we want to estimate  $\mu = \mathbb{E}X$ . Question: how many samples are sufficient for us to obtain a sufficiently good estimate of  $\mu$ , via the sample mean  $\hat{\mu} = \sum_{i=1}^m X_i/m$ ? A standard Chernoff bound argument gives

$$\mathbb{P}(|\hat{\mu} - \mu| \geq \epsilon\mu) \leq \delta \quad \text{for } m = \frac{3}{\mu\epsilon^2} \log \frac{2}{\delta}. \quad (*)$$

For example consider the estimation of  $\pi$ . Given a unit circle we tangentially inscribe it in a square with side length 2. We uniformly sample points from the square and count the fraction of points lying inside the circle, and by multiplying this fraction by 4 we obtain an estimate  $\hat{\pi}$ . By (\*) and some algebra, if  $m = 12\epsilon^{-2} \log(1/\delta)$ ,

$$\mathbb{P}(|\hat{\pi} - \pi| \geq \epsilon\pi) \leq \delta.$$

This approach is called **rejection sampling** that captures the fraction of desired samples.

**What is learning theory?** Given some data  $(x, y) \sim D$  for some distribution  $D$ , we want to learn a model to map  $x$  to  $y$ . More formally, given **concept class**  $\mathcal{C}$  of **hypotheses** (functions), in which each  $h_i$  maps  $x$  to a prediction  $h(x)$ , we want to find the most suitable one. For a hypothesis  $h$ , define

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

the probability of a wrong prediction. We will assume that there exists a perfect hypothesis  $h^*$  where  $\text{Error}(h^*) = 0$ .

A simple example of  $\mathcal{C}$  could be the set of **linear classifiers**  $h_{a,b}(x) = \text{sgn}(a^T x + b)$ , each of which partitions the space into two halfspaces. Suppose we sample  $m$  data points,  $(x_1, y_1), \dots, (x_m, y_m) \sim D$ , a distribution we do not have direct access to. Workaround and question: how large should  $m$  be so that the best hypothesis for the samples is close to  $h^*$ ?

Simple. Find large-error hypotheses and exclude them from consideration. By doing so, the final winner  $h$  will surely satisfy  $\text{Error}(h) \leq \epsilon$ . [?]

We now define VC-dimension and stuff. A **range space**  $(X, \mathcal{R})$  where  $X$  is a set of points and  $\mathcal{R} = 2^X$ , the collection of subsets of  $X$ . Given  $S \subset X$ , let  $R_S = \{R \cap S \mid R \in \mathcal{R}\}$ . A set  $S$  is said to be **shattered** by  $\mathcal{R}$  if  $|R_S| = 2^{|S|}$ , i.e., we recover all subsets of  $S$ . The **VC-dimension** is the cardinality of the largest set that can be shattered by  $\mathcal{R}$ .

The simplest example is the range space  $(\mathbb{R}, I)$  where the ranges are the set of intervals. A two-element set  $\{a, b\}$  can be shattered by for example intervals  $[a-2, a-1]$  (yielding  $\emptyset$ ),  $[a-\epsilon, a+\epsilon]$  (yielding  $\{a\}$ ),  $[b-\epsilon, b+\epsilon]$  (yielding  $\{b\}$ ), and  $[a-\epsilon, b+\epsilon]$  (yielding  $\{a, b\}$ ). However, a set of size 3  $\{a, b, c\}$  cannot be shattered: no element can cover  $\{a, c\}$  without covering  $\{b\}$ . Therefore,  $\text{VCdim}(\mathbb{R}) = \text{VCdim}(\mathbb{R}, I) = 2$ .

What about  $(\mathbb{R}^2, \mathcal{R})$  where  $\mathcal{R}$  is the set of all half spaces of  $\mathbb{R}^2$ ? It's easy to see that given three points, we can always find subspaces to split them arbitrarily. However, given 4 points, for example the four corners of a

square, no halfspace can include precisely the two diagonal elements while leaving the two other out. This implies  $\text{VCdim}(\mathbb{R}^2) = 3$ , and in general,  $\text{VCdim}(\mathbb{R}^d) = d + 1$ .

**Intuition: small VC dimension implies that the concept class has lower complexity**, so we need fewer samples to obtain the desired theoretical bounds. In terms of bias-variance tradeoff, halfspaces are simple (low variance) but may fail to fit everything (bias). Now we prove this claim.

### Theorem: The $\epsilon$ -Net Theorem

Given  $(X, \mathcal{R})$  and distribution  $D$ , a set  $N \subset X$  is an  $\epsilon$ -net if

$$|R \cap N| \neq \emptyset \quad \text{for all } R \in \mathcal{R} \text{ with } \mathbb{P}_D(R) \geq \epsilon.$$

In other words,  $N$  intersects with all elements of  $\mathcal{R}$  that is not too small. The  $\epsilon$ -net theorem states the following:

Given  $\epsilon, \delta$ , and  $\text{VCdim}(X, \mathcal{R}) = d$ , a random set of size  $\geq m = o(d\epsilon^{-1} \log(d/\epsilon) + \epsilon^{-1} \log(1/\delta))$  is an  $\epsilon$ -net with probability  $\geq 1 - \delta$ . Importantly, this bound is independent of the size of  $X$ .

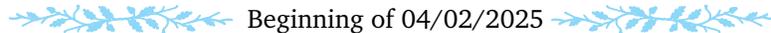
*Proof sketch.* STEP 1: Sauer's lemma. For finite  $X$ , if  $|X| = n$  and  $\mathcal{R}$  has VC-dimension  $d$ , then  $|\mathcal{R}| \leq n^d$ .

STEP 2. Suppose  $|X| = n$  and  $D$  is uniform. In this case, the ranges that we want to hit are those with  $|R| \geq \epsilon n$ . Choose  $N$ , a size  $m$  random sample of  $X$ . Then by definition and union bound,

$$\mathbb{P}(R \cap N = \emptyset) = (1 - \epsilon)^m \Rightarrow \mathbb{P}(\text{some } R \text{ satisfies } R \cap N = \emptyset) \leq n^d (1 - \epsilon)^m.$$

If we make the RHS  $< 1$  then the converse has nonzero probability (we want this to be  $1 - \delta$  by the theorem), i.e., some  $N$  will be an  $\epsilon$ -net. So we really just need  $n^d (1 - \epsilon)^m < \delta$ . Algebra shows  $m \leq d\epsilon^{-1} \log(n/\delta)$ .

STEP 3: generalizing to infinite  $X$ .



We use a technique called **double sampling** — we choose two sets  $M, T$ , both of size  $m$  from  $D$ . We want to bound the probability that  $M$  is not an  $\epsilon$ -net, i.e., the probability  $\mathbb{P}(E_1)$  where  $E_1 = \{\text{there exists } R \mid \mathbb{P}(R) \geq \epsilon, R \cap M = \emptyset\}$ .

In addition, we define  $E_2 = \{\text{there exists } R \mid \mathbb{P}(R) \geq \epsilon, R \cap M = \emptyset, \text{ and } |R \cap T| \geq \epsilon m/2\}$ . Clearly,  $E_2 \subset E_1$  so  $\mathbb{P}(E_2) \leq \mathbb{P}(E_1)$ . We now show that  $\mathbb{P}(E_1) \leq 2\mathbb{P}(E_2)$ , i.e.,  $E_2$  is not too unlikely to happen.

We use the definition of conditional probabilities:

$$\frac{\mathbb{P}(E_2)}{\mathbb{P}(E_1)} = \frac{\mathbb{P}(E_2 \cap E_1)}{\mathbb{P}(E_1)} = \mathbb{P}(E_2 \mid E_1).$$

If  $E_1$  has already happened, then there exists some heavy range  $R'$  with  $R' \cap M = \emptyset$ . Looking at this range  $R'$ ,

$$\mathbb{P}(E_2 \mid E_1) \geq \mathbb{P}(|R' \cap T| \geq \epsilon m/2).$$

Since  $\mathbb{P}(R') \geq \epsilon$ , the expected size of intersection is  $\mathbb{E}[|R' \cap T|] = \epsilon m$ . By Chernoff bound,

$$\mathbb{P}(|R' \cap T| < \epsilon m/2) \leq \exp(-\epsilon m/8).$$

By choosing a suitably large  $m$  this completes the claim that we can make  $\mathbb{P}(E_2) \leq \mathbb{P}(E_1) \leq 2\mathbb{P}(E_2)$ .

From now on we effectively only need to upper bound the probability of  $E_2$ . Clearly, an upper bound can be obtained by relaxing the first constraint (heavy set):

$$E_2 \subset E'_2 = \{\text{there exists } R \mid R \cap M = \emptyset, |R \cap T| \geq \epsilon m/2 = k\}.$$

**Claim.**  $\mathbb{P}(E'_2) \leq (2m)^d \exp(-\epsilon m/2)$ . Once this is proven, we immediately get  $\mathbb{P}(E_1) \leq 2(2m)^d \exp(-\epsilon m/2)$ , and the original claim on  $m$  follows by bounding this quantity by  $\delta$ .

*Proof.* Let  $S = M \cup T$ , and assume  $|S| = 2m$  for we are sampling from an infinite space  $X$ .

Fix one such  $S$ . The projection  $\{R \cap S, R \in \mathcal{R}\}$  of  $\mathcal{R}$  onto  $S$  has  $\leq (2m)^d$  ranges. Let  $E_R$  be the event that  $R \cap M = \emptyset$  and  $|R \cap S| \geq \epsilon m/2 = k$ . Clearly,

$$\mathbb{P}(E_R) \leq \frac{\mathbb{P}(R \cap M = \emptyset, |R \cap S| \geq k)}{\mathbb{P}(|R \cap S| \geq k)} = \mathbb{P}(R \cap M = \emptyset \mid |R \cap S| \geq k).$$

as the denominator  $\leq 1$ . Essentially, given  $S$  of size  $2m$ , we want to partition it into  $M, T$ , and ask what is the probability that all elements in  $|R \cap S|$  are not chosen by  $M$ . Algebra shows the above is

$$\leq \binom{2m-k}{m} \binom{2m}{m}^{-1} \leq 2^{-k} = 2^{-\epsilon m/2}$$

Finally, observe that (for a fixed  $S$ ),  $\mathbb{P}(E'_2) = \mathbb{P}(\bigcup_{R \in \mathcal{R}} E_R) \leq \sum_R \mathbb{P}(E_R) \leq (2m)^d \cdot 2^{-\epsilon m/2}$ .  $\square$

Going back to our original objective, where  $\mathcal{C}$  is the concept class,  $h \in \mathcal{C}$  a hypothesis, and  $\text{Error}(h) = \mathbb{P}_{(x,y) \sim D}(h(x) \neq y)$  the error rate. We assumed that there exists a perfect hypothesis  $h^*$  with  $\text{Error}(h^*) = 0$ . **How many samples  $m$  do we need to ensure that  $\text{Error}(\hat{h}) \leq \epsilon$ ?**

Given  $h$ , let  $\Delta(h) = \{x \in X : h(x) \neq h^*(x)\}$ , i.e., the collection of all points from the ambient space where  $h$  makes an error. Let  $\Delta = \{\Delta(h, h^*), h \in \mathcal{C}\}$ . Suppose  $\Delta$  has VC-dim  $d$ , applying the theorem to  $\Delta$  says by drawing sufficiently large  $m = \mathcal{O}(d\epsilon^{-1} \log(d/\epsilon) + \epsilon^{-1} \log(1/\delta))$  samples from  $D$ , then for all  $h$  with  $\mathbb{P}(\Delta(h) \geq \epsilon)$ , we can find a misclassified sample in  $\Delta(h)$ .

We were only interested in outputting a hypothesis  $\hat{h}$  that perfectly predicts all samples, i.e.,  $\hat{h}(x) = h^*(x)$  for all  $x \in \{x_1, \dots, x_m\}$ . Because our samples would hit anything heavy, i.e., exclude any hypotheses with  $\mathbb{P}(h(x) \neq h^*(x)) \geq \epsilon$ , our final winner  $\hat{h}$  must have error rate  $\epsilon$ .

Just one caveat: we need to show that  $\text{VCdim}(\Delta) = \text{VCdim}(\mathcal{C})$ . To see this, for a set  $S$ ,  $h_1 \cap S \neq h_2 \cap S$  if and only if  $\Delta(h_1) \cap S \neq \Delta(h_2) \cap S$ . Better seen by a picture. Thus, whatever shatters the concept class also shatters  $\Delta$ , and vice versa.