§ Sample Complexity and VC Dimension:

Ch 20, Notes by Har-Peled.

Sampling: using a small set of observations, estimate properties of an entire sample space.

sample complexity: minimum size sample to obtain the required result.

Let us consider two important problems. () Range detection, (2) Probability estimation,

- Range is just a subset of the underlying space. - Goal is to use one set of samples to detect a set of ranges or estimate the prob. of ranges [ the set of possible ranges can be really huge, even infinite].

For detection, we want the sample to intersect with each range in the set, while for prob. estimation, we want the fraction of points in the sample that intersect with each range in the set to approximate the assoc. prob. of the range,  $\frac{1}{3}:\frac{1}{2}:\frac{1}{2}$ 

Q. Can we obtain some sample whose size is independent of the number of ranges? and dependent on the structure of the space?

Ch 14 M-U Q. Assume we are given a range [a,b] with an underlying unknown prob. distr. D on [0,1] s.t.  $P_{x\in D}(x\in [a,b]) \ge E$ . Then how many samples do we need s.t. w.p.  $\ge 1-\delta$ . We have at least one sample in [a,b]?

Let  $x_1, x_2, \dots, x_m$  be mindep. samples in TR from an unknown distr. D

Given interval [a,b], if  $IP(x \in [a,b]) \ge \varepsilon$ . then prob that a sample of size  $m = \frac{1}{\varepsilon} e_n / \varepsilon$ intersects [a,b] is  $\ge 1 - (1-\varepsilon)^m \ge 1-\varepsilon$ .

Given k such intervals, union bound will show,  $IP[a \text{ sample of size } \frac{1}{\epsilon} \ln \frac{k}{5} \text{ intersects each of}$ the intervals]  $\geq 1 - k (1-\epsilon)^{\frac{1}{\epsilon} \ln \frac{5}{5}} \geq 1-\delta$ .

· Say, we want select few samples from [0,1] s.t. all intervals of length 1/10 contains at least one point in the sample.

- There are infinite such sets. So union bound won't help. But <u>ten</u> equidistant points would already work.

· Indeed, we'll see for any distribution, a sample size of  $\mathfrak{L}(\frac{1}{2} \ln \frac{1}{8})$  intersects all intervals having prob  $\mathfrak{Z}(\mathfrak{K},\mathfrak{W},\mathfrak{P},\mathfrak{Z}) = 1-\delta$ .

VC dimensions & Rademacher complexity helps in evaluation of sample complexity.

## § VC Dimension: (Vaprik - Chervonenkis dimension)

Definition 14.1: A range space is a pair  $(X, \mathcal{R})$  where: X is also called ground set

- 1. X is a (finite or infinite) set of points;
- 2.  $\mathcal{R}$  is a family of subsets of X, called ranges.

Example of range space: , intervals  

$$X = \mathbb{R}, \mathbb{R} = \{ [a, b] | [a, b] \subseteq \mathbb{R} \}.$$

**Definition 14.2:** Let  $(X, \mathcal{R})$  be a range space and let  $S \subseteq X$ . The projection of  $\mathcal{R}$  on *S* is



**Definition 14.3:** Let  $(X, \mathcal{R})$  be a range space. A set  $S \subseteq X$  is shattered by  $\mathcal{R}$  if  $|\mathcal{R}_S| = 2^{|S|}$ . The Vapnik–Chervonenkis (VC) dimension of a range space  $(X, \mathcal{R})$  is the maximum cardinality of a set  $S \subseteq X$  that is shattered by  $\mathcal{R}$ . If there are arbitrarily large finite sets that are shattered by  $\mathcal{R}$ , then the VC dimension is infinite.

So VC Dim of above range space (with infinite points & ranges) is only 2.

Note:  $VC \dim(R) = d$  if there is some set of cardinality d that is shattened by R. It does not say all sets of cardinality d are shattened by R. To show VC-dim  $\leq d$ , we need to show all sets of cardinality > dare not shattened by R.

- · More examples :
- Convex sets: X = IR<sup>2</sup>, R = the family of all closed convex sets on the plane.
  Claim: This range space has infinite vc-dimension.
  → Need to show, for any n ∈ IN there exists a set S with ISI=n, that can be shattered.



 $S_n = \{x_1, ..., x_n\}$  be n points on the boundary of a eincle.

Any subset  $Y \subseteq S_n$ ,  $Y \neq \phi$  defines a convex set that does not contain any points in  $S_n \setminus Y$ .

Hence, Y is included in the projection of R on Sn. Empty set is also a projection as well.

Hence, YnENS, Snis shattered.

• Disks: X = IR<sup>2</sup>, R = the family of all disks on the plane. Observation: For any 3 points on the plane (in general position) one can find eight disks so that the points are shattered. [Note: this is not true if the points are collinear. However, we just need to show one set of three points that can be shattered].



Can disks shatter a set P with four points : 1a,b,c,d? - Case 1: convex hull of P has only 3 points on its boundary, say a,b,c.

Then X = { a, b, c } can not be obtain as a projection.

Due to convexity, any disk containing a,b,c must contain d.

→ Case 2: all A points are on the convex hull. Then if we can realize {a,c} & {b,d} as projections, these two disks will intersect each other at 4 points — a contradiction.

> Note: pseudodisks are objects that can intersect at  $\in 2$  times.

§ Growth Function:  
# different subsets of  
size 
$$\leq d$$
 out of n elements:  $g(d,n) = \overset{d}{\underset{i=0}{\overset{d}{\underset{i=0}{\overset{d}{\atop}}}} \binom{n}{i}$ .  
For  $n=d$ , we have  $g(d,n) = 2^d$ ,  
for  $n > d > 2$ ,  $g(d,n) \leq \overset{d}{\underset{i=0}{\overset{d}{\underset{i=0}{\overset{d}{\atop}}}} \frac{n^i}{i!} \leq n^d$ .

- The number of ranges for a set of n elements grows polynomially in n (the power being dimension d), instead of exponentially.
- Observation: Y(d,n) = Y(d,n-1) + Y(d-1,n-1).
   not to include includes first first elements

• Saver-Shelah theorem: Let (X, R) be a range space with |X|=n, VC-dim d. Then,  $|R| \leq G(d, n) (\leq n^d)$ .

<u>Proof</u>: We prove the claim by induction on d, and for each d by induction on n.

Trivially holds for d=0 or n=0, as then  $e_j(d,n)=1$ .

Take 
$$x \in X$$
, Define  
 $R_x = \{r \mid \{x\} \mid r \cup \{x\} \in \mathbb{R} \text{ and } r \mid \{x\} \in \mathbb{R}\},$   
 $R \mid x = \{r \mid \{x\} \mid r \in \mathbb{R}\}.$ 

 $claim: | \mathcal{R}| = |\mathcal{R}_{\mathcal{X}}| + |\mathcal{R}_{\mathcal{X}}|.$ 

- We charge elements of R to their corres. element in R/x.

The only bad case is when  $\exists r s.t.$  both r U {x} and r {x} are in  $\mathbb{R}$ , as then these two distinct ranges get mapped to same range in  $\mathbb{R} \setminus \mathbb{X}$ .

But such ranges contribute to exactly one element in Rx.

We'll show vc dim of  $(X \setminus \{x\}, R_x) \leq d-1 \&$ vc-dem of  $(X \setminus \{x\}, R \setminus x)$  is  $\leq d$ . Then we'll get:  $\therefore |R| = |R_x| + |R| \approx |\leq G(d-1, n-1) + G(d, n-1)$ obs.  $\supseteq = G(d, n)$ .

In either cases, 
$$R|_{s}$$
 contains  $S \cap R' = S \cap R$ .  
Then  $R$  shatters  $S$  as well, contradicting  
 $VC$ -Dim $(X, R)$  to be  $d$ .

(i) vc-dim of 
$$(X \setminus \{x\}, \mathcal{R}_X) \leq d-1$$

### § VC-dimension component bounds:

Bounds VC-Din of a complex range space as a for of VC-din of its simpler components.

**Theorem 14.5:** Let  $(X, \mathcal{R}^1), \ldots, (X, \mathcal{R}^k)$  be k range spaces each with VC dimensions at most d. Let  $f : (\mathcal{R}^1, \ldots, \mathcal{R}^k) \to 2^X$  be a mapping of k-tuples  $(r_1, \ldots, r_k) \in (\mathcal{R}^1, \ldots, \mathcal{R}^k)$  to subsets of X, and let

$$\mathcal{R}^{f} = \{ f(r_{1}, \dots, r_{k}) \mid r_{1} \in \mathcal{R}^{1}, \dots, r_{k} \in \mathcal{R}^{k} \}.$$

The VC dimension of the range space  $(X, \mathcal{R}^f)$  is  $O(kd \ln k)$ .

This yields the following corollary.

**Corollary 14.6:** Let  $(X, \mathcal{R}^1)$  and  $(X, \mathcal{R}^2)$  be two range spaces, each with VC dimension at most d. Let

$$\mathcal{R}^{\cup} = \{ r_1 \cup r_2 \mid r_1 \in \mathcal{R}^1 \text{ and } r_2 \in \mathcal{R}^2 \},\$$

and

$$\mathcal{R}^{\cap} = \{r_1 \cap r_2 \mid r_1 \in \mathcal{R}^1 \text{ and } r_2 \in \mathcal{R}^2\}.$$

The VC dimensions of the range spaces  $(X, \mathcal{R}^{\cup})$  and  $(X, \mathcal{R}^{\cap})$  are O(d).

### § Shattering dimension.

• Defn (Shatter function): Given range space S=(X, R), its shatter function TS(m) is the maximum number of sets that might be created by S when restricted to subsets of size m.

- Shattering dimension of S is the smallest  $p \ s.t. \ \pi_s(m) = O(m^p)$ , for all m,
- · Note: In general  $TT_s(n) = 2^n$ .

Corollary: If S = (X, R) is a range space of VC-dim d. then V finite  $B \subseteq X$ , we have

 $|\mathcal{R}_{1B}| \leq \pi_s(1B1) \leq y(d, 1B1).$ 

$$\begin{array}{l} \begin{array}{l} Proof: n:= |B|, so |R_{1B}| \leq \pi_{s}(|B|) & By defn \\ of \pi_{s} \\ \leq G(d, n) & By s.s. \\ & \leq n^{d}. \end{array}$$

So, shattering dimension of a range space is bounded by its VC-din. Shatt. dime VC dim are close. Lemma (A): If S = (X, R) is a range space with shattening dim P, then VC dim (S) is  $O(P \log P)$ . Proof: Let  $N \subseteq X$  be the langest set shattened by S and S = 1 N I. We have.  $2^{S} = [R_{1N}] \leq \Pi S(1NI) \leq CS^{P}$ , where cis  $\Rightarrow S \leq lg c + P lg S$ .  $\Rightarrow P > (S - lg c) / lg S$ . Assuming,  $S \geq max(2, 2lg c)$ , we have  $\frac{\delta}{2lgS} \leq P \Rightarrow \frac{\delta}{lnS} \leq \frac{2P}{ln2} \leq 6P$   $\Rightarrow S \leq 2(6P) ln(6P)$ . [Fact: for  $u \geq JE$ , if  $\frac{\pi}{lnx} \leq u$ ]

- Advantage: shattering functions are sometimes easier to compute & gives good approximation of vc-dimension, · Shattering dimension of disks.

Lemma: Consider range space 
$$S = (X, R)$$
, where  $X = IR^2$  and  $R$  is the set of disks.  
Then shattening durn of  $S$  is  $3$ .  
Assume they are in general position.  
Proof: Consider any set  $P$  of  $n$  points in the plane and set  $\mathcal{P} = \mathcal{R}_{1P}$ .  
We claim,  $|\mathcal{P}| \leq 4n^3$ .

The set  $\mathcal{F}$  contains only n sets with a single point in them &  $\leq \binom{n}{2}$  sets with 2 points in them, So, fix  $Q \in \mathcal{F}$  s.t.  $|Q| \ge 3$ .

- (1) Let disk D realizes Q, i.e. POD = Q.
- (2) Shrink D till its boundary passes through a point p.
- 3 Continue shrinking until the boundary passes through two points p& dEQ.



(4) We continuously deform D' s.t. it has both p&d on its boundary.



This can be done by moving center of D" along the bisector line between p&q.

We continue till we hit a third point SEP.

D is a unique circle passing through p.q.s.

Also,  $D \cap (P \setminus \{s\}) = \widehat{D} \cap (P \setminus \{s\})$ ,

Thus we can specify the point set  $P \cap D$  by specifying  $(P,q,s,x_P,x_q,x_s)$ ; (P,q,s) are points defining D and  $x_* \in \{0,1\}$  states whether point \* is in Q or not. In the above case it is (1,1,0).

There are  $\leq 8\binom{n}{3}$  different subsets in  $\mathcal{F}$ containing more than 3 points, as each such subset maps to a 'canonical' disk, there are  $\binom{n}{3}$  such dists & each disk defines at most 8 different subsets.

Similar argumentation implies that there are at most  $4\binom{n}{2}$  subsets that are defined by a pair of points that realizes the diameter of the resulting disk.

 $\therefore |\mathcal{P}| = 1 + n + 4 \binom{n}{2} + 8 \binom{n}{3} \leq 4n^3.$ 

@ Insight: Above argumentation gives a powerful tool -> shattening dim of a range space defined by a family of shapes is always bounded by the number of points that determine a family.

This sometimes makes it more convenient to work with shattening dim, instead of vc dimension.

Example: shattening dimension of arbitrarily oriented rectangles is bounded by 5.



Dual shattening dimension,

Definition 20.2.8. The *dual range space* to a range space  $S = (X, \mathcal{R})$  is the space  $S^* = (\mathcal{R}, X^*)$ , where  $\mathsf{X}^{\star} = \{ \mathcal{R}_{\mathsf{p}} \mid \mathsf{p} \in \mathsf{X} \}.$ 



	0	0	0	
	$D_1$	$D_2$	$D_3$	
$p_1$	1	1	1	
$p_1'$	1	1	1	
$p_2$	1	0	1	0 0
<b>p</b> <sub>3</sub>	1	0	0	
<b>p</b> <sub>4</sub>	1	1	0	······································
<b>p</b> 5	0	1	0	
$p_6$	0	1	1	
(C)				

Let the *dual shatter function* of the range space S be  $\pi_{\mathsf{S}}^{\star}(m) = \pi_{\mathsf{S}}^{\star}(m)$ , where  $\mathsf{S}^{\star}$  is the dual range space to S. +  $H_{c}(m)$ Definition 20.2.9. The *dual shattering dimension* of S is the shattering dimension of the dual range space St. Alternately, dual shattering for S is the maximum number of points that are created when restricted to m sets I. Dual shakening dim of Sis the smallest p's.l. Trp(m) = O(mp') V m.

Clain : Dual shattening din of disks is 2. → Disks intersect each other at most 2 times.

The complexity of ampangement of n disks is  $O(2, \binom{n}{2})$ , i.e.,  $O(n^2)$ . To maximize Xt, we need at least one point in every intersection combination of ranges in R.

Hence, number of ranges in X\* < the complexity of arrangement of ranges in  $R = O(n^2)$ .

Lemma: Consider a range space S = (X, R) with vc-dim d. Then the dual range space  $S^* = (\mathcal{R}, X^*)$  has VC-dim  $\leq 2^{d+1}$ 

Lemma: If a range space S=(X, R) has dual shattening dimension 8, then its vc-dim is  $\leq \delta^{o(\delta)}$ .

Proof: The shattering dim of dual range space  $S^*$  is  $\leq \delta$ .

By Lemma A, VC-dim of S\* S' is O(Slog S). Since the dual range space to S\* is S, we have from the previous lemma: VC-dim  $(S) \leq 2^{\delta+1} = \delta^{O(\delta)}$ .

- This is very useful when shapes in R are simple. If we show the dual shattering dim of S is O(1), we also obtain Ve-dim is O(1),

- § E-nets and E-samples.
- ε-nets are combinatorial object that catches or intersects with every range of sufficient size.

**Definition 14.4** [combinatorial definition]: Let  $(X, \mathcal{R})$  be a range space, and let  $A \subseteq X$  be a finite subset of X. A set  $N \subseteq A$  is a combinatorial  $\epsilon$ -net for A if N has a nonempty intersection with every set  $R \in \mathcal{R}$  such that  $|R \cap A| \ge \epsilon |A|$ .

**Definition 14.5:** Let  $(X, \mathcal{R})$  be a range space, and let  $\mathcal{D}$  be a probability distribution on X. A set  $N \subseteq X$  is an  $\epsilon$ -net for X with respect to  $\mathcal{D}$  if for any set  $R \in \mathcal{R}$  such that  $\Pr_{\mathcal{D}}(R) \ge \epsilon$ , the set R contains at least one point from N, i.e.,

$$\forall R \in \mathcal{R}, \ \Pr_{\mathcal{D}}(R) \ge \epsilon \Rightarrow R \cap N \neq \emptyset.$$

Here, 
$$R_D(R)$$
 is the prob. that  
a point chosen according  
to D is in R.  
Note, combinatorial defn.  
corrs. to the setting  
when D is uniform over A.



An  $\varepsilon$ -net with  $\varepsilon$  = 1/4 of the unit square in the range space where the ranges are closed filled rectangles.

 The minimum sample size that contains an E-net (or E-sample) can be bounded in terms of the vC-dimension of the range space. § The E-net theorem.

Naive approach: Let (X, R) be a range space with vc-dim  $d \ge 2$ , let  $A \le X$  with |A| = n. then there exists a combinatorial  $\varepsilon$ -net N for A of size at most  $\lceil d \ln n / \varepsilon \rceil$ .

<u>Proof</u>: Let  $\mathcal{R}'$  be the projection of  $\mathcal{R}$  on  $\mathcal{A}$ . By Sauer-Shelah theorem,  $|\mathcal{R}'| \leq n^{d}$ . We take a sample of  $\mathcal{K} = \lceil d \ln n / \epsilon \rceil$  points of  $\mathcal{A}$ , independently & uniformly at random. For each set  $S \in \mathcal{R}$  with  $|S \cap \mathcal{A}| \geq \epsilon |\mathcal{A}|$ , there is a corresponding set  $S' \in \mathcal{R}'$ . So a point in our sample is in S' with prob  $\geq \epsilon$ . Then IP [our sample misses a given set  $S' \rceil$  $\leq (1-\epsilon)^{\mathcal{K}}$ .

There are nd such sets to consider. [From (B)] Applying union bound, the prob that the sample misses at least one such S' is

 $\leq n^{d}(1-\varepsilon)^{k} < n^{d} e^{-\varepsilon k} < n^{d} e^{-d \ln n} = 1.$ 

Thus by probabilistic method, there is a set N of size k that misses no set  $S' \in \mathbb{R}^{\prime}$ .

Hence, N is an E-net for A.

Now we prove the main theorem that shows existence of  $\varepsilon$ -net of size  $O(\frac{d}{\varepsilon} \ln \frac{d}{\varepsilon})$  independent of n.

• Theorem: Let (X, R) be a range space with VC-dim d and let  $\mathcal{D}$  be a prob. distribution on X. For any  $0 < \delta, \epsilon \leq \frac{1}{2}$ , there is an  $m = O\left(\frac{d}{\epsilon} \ln \frac{d}{\epsilon} + \frac{1}{\epsilon} \ln \frac{1}{\delta}\right)$  such that a random sample from  $\mathcal{D}$  of size m is an  $\epsilon$ -net for X with probability at least  $1-\delta$ .

#### Proof:

Let M be a set of m independent samples from X acc. to D.

Let  $E_1 := \{ \exists s \in \mathbb{R} \mid Pr_2(s) \geq \varepsilon \text{ and } |s \cap M| = 0 \}$ , i.e.  $E_1$  is the event that M is not an  $\varepsilon$ -net for X w.r.t. D.

We want to show,  $P(E_1) \leq \delta$ .

For this, we go beyond the union bound approach.

Choose a second set T of m indep. samples from X acc. 2.

Define,  $E_2 := \{ \exists S \in \mathcal{R} \mid \mathcal{R}_{\mathcal{D}}(S) \geq \varepsilon, |S \cap M| = 0, |S \cap T| \geq \varepsilon_{\mathcal{D}_2} \}$ 

Following lemma shows E1 & E2 have similar probabilities. strivial as E2 E1.

Lemma 1: for  $m \ge 8/\epsilon$ ,  $P(E_2) \le P(E_1) \le 2 P(E_2)$ . <u>Proof</u>: If event  $E_1$  holds, there is some S' s.t.  $IS' \cap M I = 0$  and  $Pr_D(S') \ge \epsilon$ .

Hence, 
$$\frac{\mathbb{P}(\mathbb{E}_2)}{\mathbb{P}(\mathbb{E}_1)} = \frac{\mathbb{P}(\mathbb{E}_1 \cap \mathbb{E}_2)}{\mathbb{P}(\mathbb{E}_1)} = \mathbb{P}(\mathbb{E}_2 | \mathbb{E}_1) \ge \mathbb{P}(|\mathcal{T} \cap S'| \ge \mathbb{E}_m/2).$$

Now for a fixed range S' and a random sample T, the random variable  $1 \pm nS'$  has a binomial distr.  $B(m, Fr_0(S'))$ .

As 
$$R_{\mathcal{D}}(S') \ge \varepsilon$$
, using Chernoff bounds:  
 $P(|T \cap S'| < (1 - \delta) \mathbb{E}[|T \cap S'|]) \le e^{-\frac{\delta^2}{2} \cdot \mathbb{E}[|T \cap S'|]}$   
 $\stackrel{f}{=} P(|T \cap S'| < (1 - \frac{1}{2}) \mathbb{m}\varepsilon) \le e^{-\frac{1}{4} \cdot \frac{1}{2} \cdot \mathbb{m}\varepsilon} = e^{\varepsilon m/8} \le e^{-1} < \frac{1}{2}$   
Hence, from (1) & (2),  $P(\varepsilon_1) \le 2 \cdot P(\varepsilon_2)$ .

Now we bound  $IP[E_2]$  by prob. of a larger event  $E_2'$ .  $E_2' := \{ \exists S \in \mathcal{R} \mid |S \cap M| = 0 \text{ and } |S \cap T| \ge Em/2 \}.$ 

• Lemma 2.  $P(E_1) \leq 2P(E_2) \leq 2P(E_2) \leq 2(2m)^d 2^{-\epsilon m/2}$ .

<u>Proof</u>: As M& T are random samples, we assume to choose 2m random samples and partition randomly into two equal sized sets M&T.

For a fixed 
$$S \in \mathbb{R}$$
,  $K = \frac{m}{2}$ , let  
 $E_s := \{1S \cap M = 0, 1S \cap T \mid \geqslant k\}$ 

This event means  $(M \cup T) \cap S \ge K$ , but all these elements were placed in T and not in M. So, out of  $\binom{2m}{m}$  possible partitions of MUT, we choose one of  $\binom{2m-k}{m}$  partitions where no element of S is in M.

Hence, 
$$IP(E_S) \leq IP(|MnS|=0||Sn(MUT)| \geq k)$$
  
=  $\binom{2m-k}{m} / \binom{2m}{m}$   
=  $(2m-k)!m!/(2n)!(m-k)!$   
=  $\frac{m(m-1)...(m-k+1)}{(2n)(2m-1)...(2m-k+1)} \leq 2^{-k} \leq 2^{-2m/2}$ .

By Saver-Shelah theorem, the projection of 
$$\mathcal{R}$$
  
on MUT has  $\leq (2m)^d$  ranges.

Hence, using union bounds:  $\mathbb{P}(\mathbb{E}_2') \leq (2m)^d 2^{-\frac{2m}{2}}$ 

To complete the proof of  $\varepsilon$ -net theorem, we need to show  $P(\varepsilon_1) \leq 2(2m)^d 2^{-sm/2} \leq \delta$ ; for  $m \geq \frac{8d}{\varepsilon} \ln \frac{16d}{\varepsilon} + \frac{4}{\varepsilon} \ln \frac{2}{\delta}$ . equ., we need  $\varepsilon m/2 \geq \ln(2/\delta) + d\ln(2m)$ . As  $m \geq \frac{4}{\varepsilon} \ln \frac{2}{\delta}$ ,  $\varepsilon m/4 \geq \ln(2/\delta)$ . To finish we show  $\varepsilon m/4 \geq d\ln(2m)$ . [Fact: if  $y \geq x \ln x \geq \varepsilon$ . then  $\frac{2y}{\ln y} \geq x$ ] Use  $y = 2m \geq \frac{16d}{\varepsilon} \ln \frac{16d}{\varepsilon}$ ,  $x = \frac{16d}{\varepsilon}$ , we have  $\frac{4m}{\ln(2n)} \geq \frac{16d}{\varepsilon} \Rightarrow \frac{sm}{4} \geq d\ln(2m)$ . § Application:

Probably Approximately Correct (PAC) Learning. We are given a set of items X and a prob distr. D is defined on X.

A binary classification is a subset C = X s.t. all items in C are labeled 1 and in X\C are labeled -1.

The concept class & is the set of all possible classifications defined by the problem.

- · Learning algo has access to ORACLE (C, D) that produces a pair (x, c(x)), where  $x \sim D$ and c(x) = 1 if  $x \in C$  and -1 otherwise.
- · We also assume the classification problem is realizable, i.e.  $\exists h \in \mathcal{C}, P_{\mathcal{P}_{\mathcal{D}}}(h(x) \neq c(x)) = 0$ .

**Definition 14.9** [PAC Learning]: A concept class C over input set X is PAC learnable<sup>1</sup> if there is an algorithm L, with access to a function ORACLE(C, D), that satisfies the following properties: for every correct concept  $C \in C$ , every distribution  $\mathcal{D}$  on X, and every  $0 < \epsilon, \delta \leq 1/2$ , the number of calls that the algorithm L makes to the function ORACLE(C, D) is polynomial in  $\epsilon^{-1}$  and  $\delta^{-1}$ , and with probability at least  $1 - \delta$  the algorithm L outputs a hypothesis h such that  $\Pr_{\mathcal{D}}(h(x) \neq c(x)) \leq \epsilon$ . approximately probably

- Theorem: Any finite concept class & can be PAC-learned with  $m = \frac{1}{e} (lnlel + ln \frac{1}{5})$  samples. Proof: let  $c^* \in \mathcal{C}$  be the correct classification. A hypothesis h is "bad" if  $\operatorname{Pr}_{\mathcal{D}}[h(z) \neq c^*(z)] \ge \varepsilon$ .
- P[a bad h is consistent with m random samples] $<math>\leq (1-\epsilon)^{m}$

Using union bound,

- $P[\exists bad h is consistent with m random samples]$  $<math>\leq |\mathcal{L}| (1-\varepsilon)^{m}$ 
  - $\leq \delta \left[ as m = \frac{1}{\epsilon} \left( \ln \left| \epsilon \right| + \ln \frac{1}{\delta} \right) \right]$

so, we return whatever to is consistent with an m random sample. - all these are "non-bad".

By assumption. as connect classification et e. e., at least one such h exists.

- Note that X can be infinite. So, it is interesting that we can PAC-learn R, with sample complexity independent of n.

Note: A concept class is efficiently PAC learnable if the algorithm runs in time polynomial in the size of problem, 1/E, 1/S.

Here, we are only interested in sample complexity, computational complexity may not necessarily be polynomial in sample size.

- · can we make sample complexity indep of 1,e1?
- · Can we extend this to infinite concept classes?

Say, we are learning interval  $[a, b] \in \mathbb{R}$ . Concept class is collection of all closed intervals in  $\mathbb{R}$ :  $\mathcal{L} = \{[z, y] \mid z \leq y\} \cup \phi$ .

Let  $c^* \in C$  be the concept to be learned. In be the hypothesis returned by our algo.

Training set T is collection of n points drawn from D.

Let  $x \in T$ , if  $x \in [a, b]$  it is a tre example, else a -re example.

Algo: If no sample is positive return trivial hypothesis.

else return [c,d] where c,d are smallest and largest +ve examples among the samples.

Q. What is the prob ALGO makes an error?

ALGO can only make an error on an input xif  $x \in [a, b]$ . For  $x \notin [a, b]$ , it always returns - 1.

Case 1.  $\operatorname{Pr}_{\mathcal{D}}(x \in [a, b]) \leq \varepsilon$ . From the above fact, prob of error  $\leq \varepsilon$  Case 2.  $Pr_{\mathcal{D}}(x \in [a, b]) > C$ .



Say,  $a' \ge a$  be smallest ral s.t.  $\operatorname{Re}_{a}([a,a']) \ge \frac{e}{2}$ , and  $b' \le b$  be largest val s.t.  $\operatorname{Re}_{a}([b',b]) \ge \frac{e}{2}$ .

so, a' < b. For simplicity assume a' < b.

9f ALGO returns a bad hypothesis then error ≥ €, Now if sample points fell in [a, a'] and [b, b'] then we would have returned [c, d] ⊇ [a', b']. → a good hypothesis So, IP[bad hypothesis] ≤ IP( sample points didn't fall, in [a, a'] on [b, b']]. either

The prob. that a training set of a points does not contain any examples from either [a, a'] or [b, b'] is

$$\leq 2\left(1-\frac{\varepsilon}{2}\right)^n \leq 2e^{-\varepsilon n/2} \leq \delta.$$
  

$$\int f$$
  
by choosing  $n \geq 2\ln(2/\delta)/\varepsilon$ 

· E-net & VC-din generalizes this idea.

**Theorem 14.13:** Let C be a concept class that defines a range space with VC dimension d. For any  $0 < \delta$ ,  $\epsilon \le 1/2$ , there is an over a set of

$$m = O\left(\frac{d}{\epsilon}\ln\frac{d}{\epsilon} + \frac{1}{\epsilon}\ln\frac{1}{\delta}\right)$$

points X

such that C is PAC learnable with m samples.

To complete the proof of lemma, we need to show bis bijection. Take  $c', c'' \in \mathcal{C}$  with  $c' \cap S \neq c'' \cap S$ . Then  $\exists y \in S \ s.t. \ c'(y) \neq c''(y)$ . w.1. o.g. assume  $c'(y) \neq c(y)$  but c''(y) = c(y). Then,  $y \notin \Delta(c' \cap S, c) \cap S$ but  $y \in \Delta(c'' \cap S, c \cap S)$ . Hence,  $\Delta(c' \cap S, c \cap S) \neq \Delta(c'' \cap S, c \cap S)$ .

For the other direction, if for  $c, c' \in \mathcal{R}$ , s.t.  $\Delta(c' \cap s, c \cap s) \neq \Delta(c'' \cap s, c \cap s)$ , then  $\exists y \in S$  s.t.  $c'(y) \neq c''(y)$ , so  $c' \cap S \neq c'' \cap S$ .

Thus as  $VC-Din[(X, \Delta(c))] = d$ , there exists  $m = O(d/\epsilon \ln d/\epsilon + /\epsilon \ln /\delta) \text{ s.t. sample of size }m$ is  $\epsilon$ -net for this range-space  $w.p. \ge 1-\delta$ . Hence,  $w.p. \ge 1-\delta$  it has nonempty intersection with every set  $\Delta(c', c)$  that has prob >  $\epsilon$ ., i.e. ALGO can exclude any hyp. w. error/prob >  $\epsilon$ .

> So any c' far from correct c is anyway bit. & thus can be excluded

As we only select hypothesis that are consistent with an the samples.

· E-sample provides stronger guarantees.

- it maintains relative probability weight all sets  $R \in R$  within error of e, and needs just additional  $O(\gamma_{e})$  factor in sample size.. **Definition 14.6:** Let  $(X, \mathcal{R})$  be a range space, and let  $\mathcal{D}$  be a probability distribution on *X*. A set  $S \subseteq X$  is an  $\epsilon$ -sample for *X* with respect to  $\mathcal{D}$  if for all sets  $R \in \mathcal{R}$ ,

$$\left| \Pr_{\mathcal{D}}(R) - \frac{|S \cap \overline{R}|}{|S|} \right| \le \epsilon.$$
 Relative freq.

Again, by fixing the distribution  $\mathcal{D}$  to be uniform over a finite set  $A \subseteq X$ , we obtain the combinatorial version of this concept.

**Definition 14.7** [combinatorial definition]: Let  $(X, \mathcal{R})$  be a range space, and let  $A \subseteq X$  be a finite subset of X. A set  $N \subseteq A$  is a combinatorial  $\epsilon$ -sample for A if for all sets  $R \in \mathcal{R}$ ,

$$\left|\frac{|A \cap R|}{|A|} - \frac{|N \cap R|}{|N|}\right| \le \epsilon.$$

**Definition 14.8:** A range space  $(X, \mathcal{R})$  has the uniform convergence property if for every  $\epsilon, \delta > 0$  there is a sample size  $m = m(\epsilon, \delta)$  such that for every distribution  $\mathcal{D}$  over X, if S is a random sample from  $\mathcal{D}$  of size m then, with probability at least  $1 - \delta$ , S is an  $\epsilon$ -sample for X with respect to  $\mathcal{D}$ .

E-sample theorem :

**Theorem 14.15:** Let  $(X, \mathcal{R})$  be a range space with VC dimension d and let  $\mathcal{D}$  be a probability distribution on X. For any  $0 < \epsilon, \delta < 1/2$ , there is an

$$m = O\left(\frac{d}{\epsilon^2}\ln\frac{d}{\epsilon} + \frac{1}{\epsilon^2}\ln\frac{1}{\delta}\right)$$

such that a random sample from  $\mathcal{D}$  of size greater than or equal to m is an  $\epsilon$ -sample for X with probability at least  $1 - \delta$ .

# § Application : Agnostic Learning.

In PAC learning, we assumed there is a  $c^* \in C$ that is correct on all items in X and so conforms with all examples in training set.

→ But training set can have error & there may not be any correct classification in C. In agnostic learning, the goal is find a nearly best classification c' s.t.

$$R_{\mathcal{D}}(c'(x) \neq c(x)) \leq \inf_{h \in \mathcal{R}} (h(x) \neq c(x)) + \epsilon$$
  

$$\int_{h \in \mathcal{R}} (h(x) \neq c(x)) + \epsilon$$
  
connect classif,  
may not be in  $\mathcal{R}$ 

If the training set define an  $\frac{\epsilon}{2}$ -sample for (X,  $\Delta(c)$ ) then algo has sufficiently many examples to estimate the error prob of each c'  $\in \mathcal{C}$  to within an additive error  $\frac{\epsilon}{2}$ .

Using E-sample theorem, agnostic learning of a concept class with VC-dim d requires  $O\left(\frac{d}{\varepsilon^2}\ln\frac{d}{\varepsilon^2} + \frac{1}{\varepsilon^2}\ln\frac{1}{\varepsilon}\right)$  samples.

### **Theorem 14.19:** *The following three conditions are equivalent:*

1. A concept class C over a domain X is agnostic PAC learnable.

- **2.** The range space (X, C) has the uniform convergence property.
- **3.** The range space (X, C) has a finite VC dimension.

§ Applications: Data mining. ① Estimating dense neighborhoods. Given: n points in  $\mathbb{R}^2$ , p = (x, y),  $r \in \mathbb{R}$ . Goal: What fraction of points are distance  $\leq r$  from (x, y).

[Application: opening new facility/business]



We define range space  $(\mathbb{R}^2, \mathbb{R})$  where  $\mathbb{R}$ includes  $\forall (x, y) \in \mathbb{R}^2$  and  $r \in \mathbb{R}^+$ , set of all points inside the disk of radius rcentered at (x, y).

vc-Din of the set of all disks = 3. constant we can sample a random set of  $O\left(\frac{1}{\varepsilon^2}\ln\frac{1}{\varepsilon\delta}\right)$ points and give fast approximate answers to all the queries by scanning only the sample.

E-sample theorem guarantees w.p.  $\geq 1-\delta$ , we can answer all queries within  $\in$  of the correct value.

we can also it for other purposes, such as (approx) identifying k densest disks.

- similar example: Range searching. (How many points are included in a quepy rectangle?) 2 Mining frequent itemsets.

→ Given: A set of items I, a collection of transactions J, where each transaction  $t \in J$ is  $\subseteq I$ . - Both |I|, |J| are large.

Coal: Find set of items that appear in > 0 fraction of transactions.

Say we want to characterize freq  $\geq 0$  to be frequent items, freq  $\leq 0 - \epsilon$  to be infrequent. Freq.  $[0 - \epsilon, 0]$  can be ambriguous.

The number of possible of transactions = subsets of I is huge! Even for transactions of size  $\leq l$ , there can be  $O(1 \times l^k)$  of them which could be frequent:

(HW: Using Chernoff + union bound would give  $\mathcal{Q}\left(\frac{Q}{E^2}\left(1 \ln |\mathcal{I}| + \ln \frac{1}{5}\right)\right)$  samples are needed,)

E-sample does better! transactions that includes. For each  $s \subseteq I$ ,  $T(s) := \{I \in J, S \subseteq I\}$ Let  $R = \{T(s) | S \subseteq I\}$ , the max size of any transaction claim: VC-Dim  $[(T, R)] \leq l$ . in the data set - A transaction of size l has  $2^{d}$  subsets. is there fore included in  $\leq 2^{d}$  ranges. Now, consider  $S \subseteq J$ ,  $R_{S} = \{R \cap S | R \in R\}$ . Then,  $|R_{S}| \leq 2^{d}$ , as can belong to at nost  $2^{d}$ ranges. Thus no more than I transactions can be shattered. > VC-dein=L.

Hence, by  $\varepsilon$ -sample theorem,  $v.p. > 1-\delta$ , a sample of size  $O(\frac{l}{\varepsilon^2} \ln \frac{l}{\varepsilon} + \frac{l}{\varepsilon^2} \ln \frac{1}{\delta})$ guarantee all itemsets are accurately determined to within  $\frac{\sigma_2}{2}$  of their true proportion

- This is enough to identify frequent itemsets.

§ Rademacher complexity.

- Bounds can depend on the training set distribution
- Generalizes to nonbinary functions.