

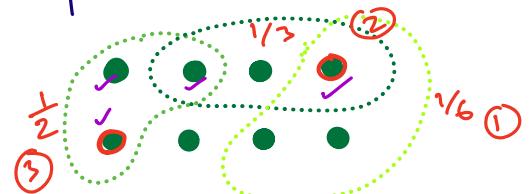
§ Sample Complexity and VC Dimension :Ch 20, Notes by
Har-Peled.

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.

- ① 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. $\frac{1}{3} : \frac{1}{2} : \frac{1}{6}$

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

Q. Assume we are given a range $[a, b]$ with an underlying unknown prob. distr. \mathcal{D} on $[0, 1]$ s.t.

$$\mathbb{P}_{x \in \mathcal{D}}(x \in [a, b]) \geq \epsilon.$$

Then how many samples do we need s.t. w.p. $\geq 1 - \delta$. we have at least one sample in $[a, b]$?

Let x_1, x_2, \dots, x_m be m indep. samples in \mathbb{R} from an unknown distr. \mathcal{D}



Given interval $[a, b]$, if $\mathbb{P}(x \in [a, b]) \geq \epsilon$, then prob that a sample of size $m = \frac{1}{\epsilon} \ln \frac{1}{\delta}$ intersects $[a, b]$ is $\geq 1 - (1 - \epsilon)^m \geq 1 - \delta$.

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

- 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 ten equidistant points would already work.
- Indeed, we'll see for any distribution, a sample size of $\leq 2(\frac{1}{\epsilon} \ln \frac{1}{\delta})$ intersects all intervals having prob $\geq \epsilon$, w.p. $\geq 1 - \delta$.

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

§ VC Dimension: (Vapnik - Chervonenkis dimension)

Definition 14.1: A range space is a pair (X, \mathcal{R}) where:

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

X is also called ground set

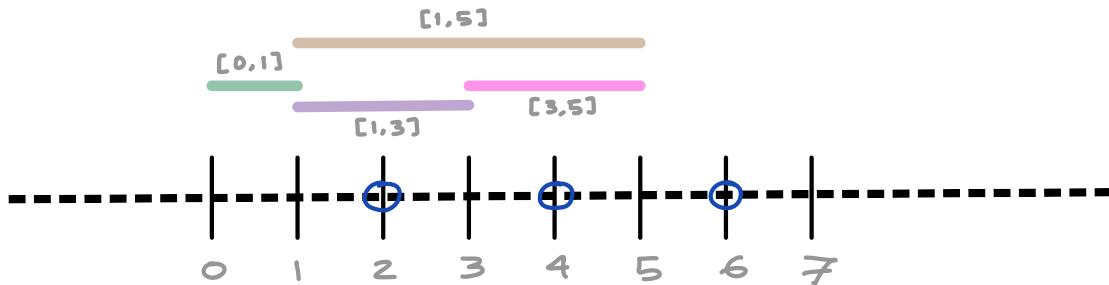
Example of range space:

$$X = \mathbb{R}, \quad \mathcal{R} = \{ [a, b] \mid [a, b] \subseteq \mathbb{R} \}.$$

set of all closed intervals

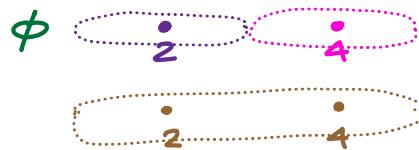
Definition 14.2: Let (X, \mathcal{R}) be a range space and let $S \subseteq X$. The projection of \mathcal{R} on S is

$$\mathcal{R}_S = \{R \cap S \mid R \in \mathcal{R}\}.$$



$$S = \{2, 4\}$$

Any two points can be shattered.



\mathcal{R}_S is the set of all possible subsets of S .

$$S = \{2, 4, 6\}$$

Any 3 points can't be shattered.

\mathcal{R}_S gives seven of the eight possible subsets of S , except $\{2, 6\}$.

- Any interval containing 2 & 6 must contain 4.

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.

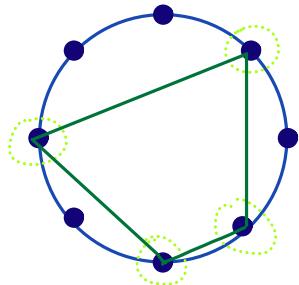
Note: $\text{VC dim}(\mathcal{R}) = d$ if there is some set of cardinality d that is shattered by \mathcal{R} . It does not say all sets of cardinality d are shattered by \mathcal{R} . To show $\text{VC-dim} \leq d$, we need to show all sets of cardinality $> d$ are not shattered by \mathcal{R} .

- More examples:

- **Convex sets:** $X = \mathbb{R}^2$, \mathcal{R} = the family of all closed convex sets on the plane.

Claim: This range space has infinite VC-dimension.

→ Need to show, for any $n \in \mathbb{N}$ there exists a set S with $|S| = n$, that can be shattered.



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

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

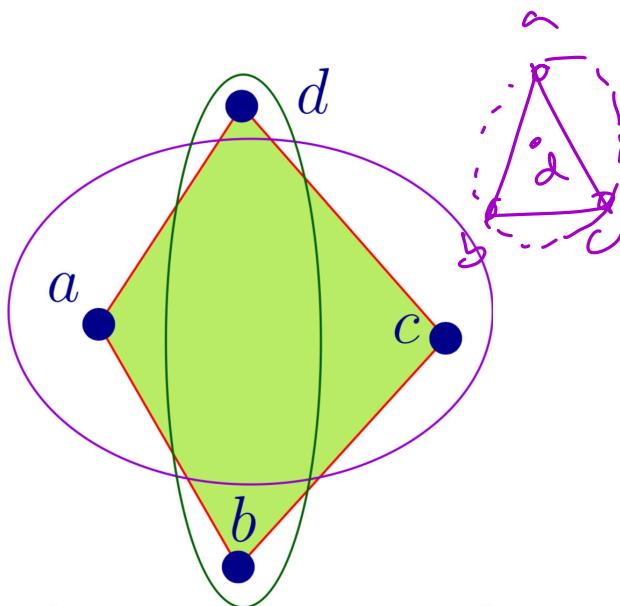
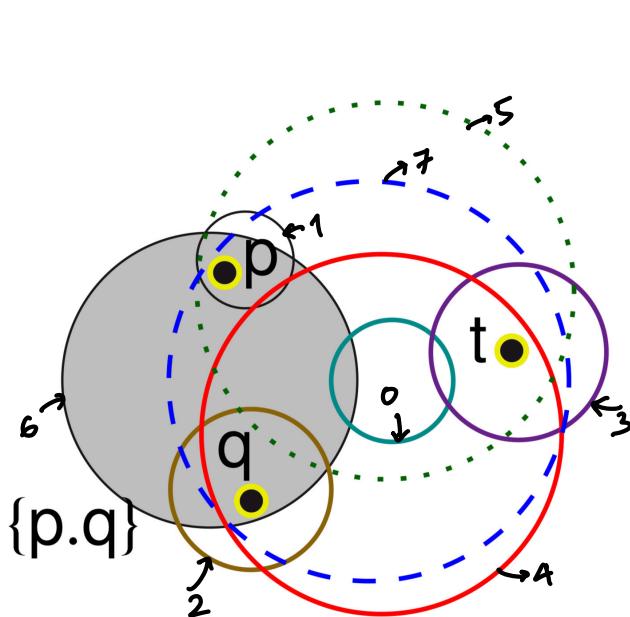
Hence, Y is included in the projection of \mathcal{R} on S_n .
Empty set is also a projection as well.

Hence, $\forall n \in \mathbb{N}$, S_n is shattered.

- **Disks:** $X = \mathbb{R}^2$, \mathcal{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: $\{a, b, c, d\}$?

→ Case 1: convex hull of P has only 3 points on its boundary, say a, b, c .

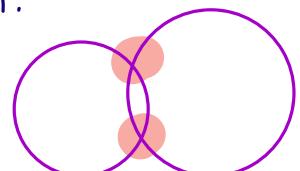
Then $X = \{a, b, c\}$ cannot be obtain as a projection.

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

→ Case 2: all 4 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 ≤ 2 times.



§ Growth Function:

different subsets of size $\leq d$ out of n elements : $g(d, n) = \sum_{i=0}^d \binom{n}{i}$.

For $n=d$, we have $g(d, n) = 2^d = 2^n$.

for $n > d \geq 2$, $g(d, n) \leq \sum_{i=0}^d \frac{n^i}{i!} \leq n^d$.

By convention $g(0, 0) = 1$.

- The number of ranges for a set of n elements grows polynomially in n (the power being dimension d), instead of exponentially.
- Observation: $g(d, n) = g(d, n-1) + g(d-1, n-1)$.

 \downarrow
not to include first element
 \downarrow
includes first elements
- Sauer-Shelah theorem:

Let (X, \mathcal{R}) be a range space with $|X|=n$, VC-dim d . Then, $|\mathcal{R}| \leq g(d, n) (\leq n^d)$.

Proof: 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 $g(d, n)=1$.

Take $x \in X$, Define

$$\mathcal{R}_x = \{r \setminus \{x\} \mid r \cup \{x\} \in \mathcal{R} \text{ and } r \setminus \{x\} \in \mathcal{R}\}.$$

$$\mathcal{R} \setminus x = \{r \setminus \{x\} \mid r \in \mathcal{R}\}.$$

claim: $|R| = |Rx| + |R \setminus x|$.

- We charge elements of R to their corrs. element in $R \setminus x$.

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

But such ranges contribute to exactly one element in Rx . ■

We'll show VC-dim of $(X \setminus \{x\}, Rx) \leq d-1$ & VC-dim of $(X \setminus \{x\}, R \setminus x)$ is $\leq d$.

Then we'll get:

$$\begin{aligned} \therefore |R| &= |Rx| + |R \setminus x| \stackrel{\text{claim}}{\leq} g(d-1, n-1) + g(d, n-1) \\ &\stackrel{\text{obs.}}{\rightarrow} = g(d, n). \end{aligned}$$

induction.

To conclude:

① VC-dim of $(X \setminus \{x\}, R \setminus \{x\}) \leq d$.

→ For contradiction, assume $R \setminus \{x\}$ shatters a set $S \subseteq X \setminus \{x\}$ of size $d+1$.

Now for every $R \in R \setminus \{x\}$, there is a corrs. $R' \in R$ s.t. either $R' = R$ or $R' = R \cup \{x\}$.

In either cases, $\mathcal{R}|_S$ contains $S \cap R' = S \cap R$.

Then \mathcal{R} shatters S as well, contradicting $\text{VC-Dim}(X, \mathcal{R})$ to be d .

ii) $\text{vc-dim of } (X \setminus \{x\}, \mathcal{R}_x) \leq d-1$

\rightarrow for contradiction, assume \mathcal{R}_x shatters a set $S \subseteq X \setminus \{x\}$ of size d .

Now for every $R \in \mathcal{R} \setminus \{x\}$, both R and $R \cup \{x\}$ are in \mathcal{R} .

In both cases, $\mathcal{R}|_S$ contains $S \cap (R \cup \{x\}) = S \cap R$.

Then \mathcal{R} also shatters S . \rightarrow A contradiction. \blacksquare

§ VC-dimension component bounds:

Bounds VC-dim of a complex range space as a fn of vc-dim of its simpler components.

Theorem 14.5: Let $(X, \mathcal{R}^1), \dots, (X, \mathcal{R}^k)$ be k range spaces each with VC dimensions at most d . Let $f : (\mathcal{R}^1, \dots, \mathcal{R}^k) \rightarrow 2^X$ be a mapping of k -tuples $(r_1, \dots, r_k) \in (\mathcal{R}^1, \dots, \mathcal{R}^k)$ 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 $O(kd \ln k)$.

circle
squares
etc.

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 $O(d)$.

VC dim of
△
↓
VC dim of
K gom.

§ Shattering dimension.

- Defn (Shatter function) : Given range space $S = (X, \mathcal{R})$, its shatter function $\pi_S(m)$ is the maximum number of sets that might be created by S when restricted to subsets of size m .

$$\pi_S(m) = \max_{\substack{B \subseteq X \\ |B|=m}} |\mathcal{R}_{|B|}|.$$

projection of
 \mathcal{R} on B

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

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

$$|\mathcal{R}_{|B|}| \leq \pi_S(|B|) \leq g(d, |B|).$$

Proof: $n := |B|$, so $|\mathcal{R}_{|B|}| \leq \pi_S(|B|)$ [By defn of π_S]

$$\begin{aligned} &\leq g(d, n) \quad (\text{By S.S. lemma}) \\ &\leq n^d. \end{aligned}$$

So, shattering dimension of a range space is bounded by its VC-dim.

shatt. dim & VC dim
are close.

Lemma A: If $S = (X, \mathcal{R})$ is a range space with shattering dim ρ , then $\text{VCdim}(S)$ is $O(\rho \log \rho)$.

Proof: Let $N \subseteq X$ be the largest set shattered by S and $\delta = |N|$.

We have, $2^\delta = |\mathcal{R}_{IN}| \stackrel{\text{shattering defn}}{\leq} \Pi_S(|N|) \stackrel{\text{corollary}}{\leq} c\delta^\rho$, where c is a constant.

$$\Rightarrow \delta \leq \lg c + \rho \lg \delta.$$

$$\Rightarrow \rho \geq (\delta - \lg c) / \lg \delta.$$

otherwise
 δ is already
 $O(1)$, i.e.
 $O(1/\log \delta)$.

Assuming, $\delta \geq \max(2, 2 \lg c)$, we have

$$\frac{\delta}{2 \lg \delta} \leq \rho \Rightarrow \frac{\delta}{\ln \delta} \leq \frac{2\rho}{\ln 2} \leq 6\rho$$

$$\Rightarrow \delta \leq 2(6\rho) \ln(6\rho). \quad \begin{array}{l} \text{Fact: for } u \geq \sqrt{e}, \text{ if } \frac{x}{\ln x} \leq u \\ \text{then } x \leq 2u \ln u. \end{array}$$

■

→ Advantage: shattering functions are sometimes easier to compute & gives good approximation of VC-dimension.

- Shattering dimension of disks.

Lemma: Consider range space $S = (X, \mathcal{R})$, where $X = \mathbb{R}^2$ and \mathcal{R} is the set of disks.

Then shattering dim of S is 3.

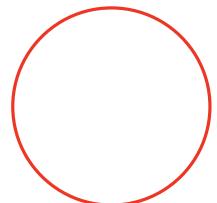
assume they
are in general
position.

Proof: Consider any set P of n points in the plane and set $\mathcal{F} = \mathcal{R}_{IP}$.

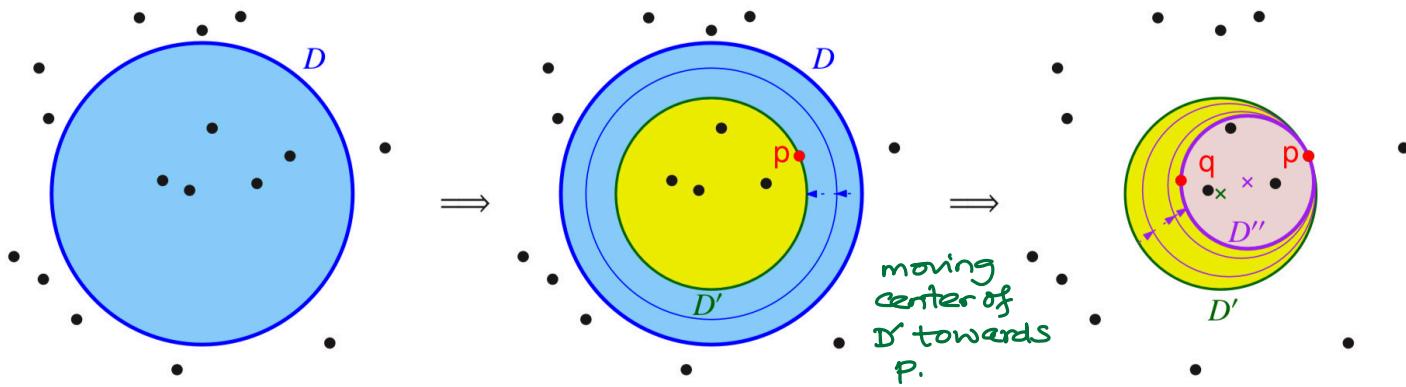
We claim, $|\mathcal{F}| \leq 4n^3$.

The set \mathcal{F} contains only n sets with a single point in them & $\leq \binom{n}{2}$ sets with 2 points in them.

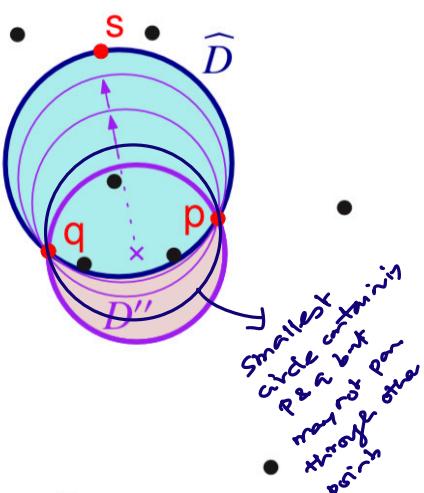
So, fix $Q \in \mathcal{F}$ s.t. $|Q| \geq 3$.



- ① Let disk D realizes Q , i.e. $P \cap D = Q$.
- ② Shrink D till its boundary passes through a point p .
- ③ Continue shrinking until the boundary passes through two points $p, q \in Q$.



④ We continuously deform D'' s.t. it has both p & q 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 $s \in P$.

\hat{D} is a unique circle passing through P, Q, S .

$$\text{Also, } D \cap (P \setminus \{s\}) = \hat{D} \cap (P \setminus \{s\}).$$

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 disks & 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.

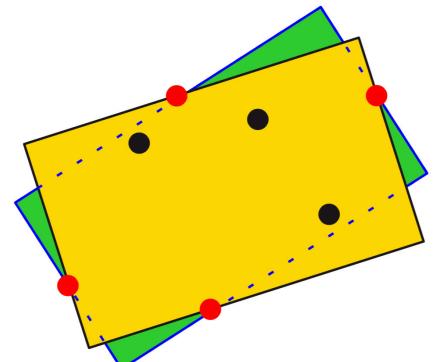
$$\therefore |\mathcal{R}| = 1 + n + 4 \binom{n}{2} + 8 \binom{n}{3} \leq 4n^3.$$

■

 Insight: Abore argumentation gives a powerful tool → shattering 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 shattering dim. instead of VC dimension.

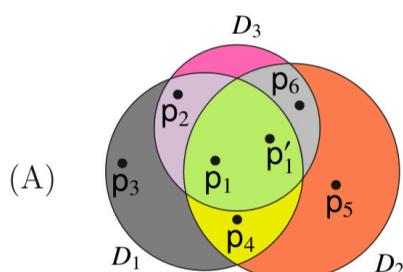
Example: shattering dimension of arbitrarily oriented rectangles is bounded by 5.



Dual shattering 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 $X^* = \{\mathcal{R}_p \mid p \in X\}$.

Ranges $\xleftrightarrow{} \text{Points}$

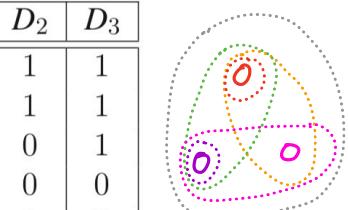


$$S^{**} = S$$

	0	0	0
D ₁	1	1	1
D ₂	1	1	1
D ₃	1	0	1
p ₁	1	1	1
p' ₁	1	1	1
p ₂	1	0	1
p ₃	1	0	0
p ₄	1	1	0
p ₅	0	1	0
p ₆	0	1	1

(B)

	p ₁	p' ₁	p ₂	p ₃	p ₄	p ₅	p ₆
D ₁	1	1	1	1	1	0	0
D ₂	1	1	0	0	1	1	1
D ₃	1	1	1	0	0	0	1



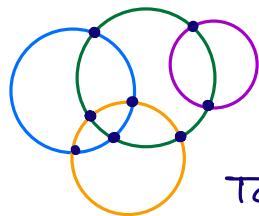
(C)

Let the **dual shatter function** of the range space S be $\pi_S^*(m) = \pi_{S^*}(m)$, where S^* is the dual range space to S .

Definition 20.2.9. The **dual shattering dimension** of S is the shattering dimension of the dual range space S^* . Alternately, dual shattering fn for S is the maximum number of points that are created when restricted to m sets \mathcal{F} . Dual shattering dim of S is the smallest ρ 's.t. $\pi_{\mathcal{F}}(m) = O(m^\rho) \nabla m$.

Claim : Dual shattering dim of disks is 2.

→ Disks intersect each other at most 2 times.



The complexity of arrangement of n disks is $O(2 \cdot \binom{n}{2})$, i.e., $O(n^2)$.

To maximize x^* , we need at least one point in every intersection combination of ranges in \mathcal{R} .

Hence, number of ranges in $x^* \leq$ the complexity of arrangement of ranges in $\mathcal{R} = O(n^2)$. ■

Lemma: Consider a range space $S = (X, \mathcal{R})$ 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 $S = (X, \mathcal{R})$ has dual shattering dimension δ , 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^* δ' is $O(\delta \log \delta)$.

Since the dual range space to S^* is S , we have from the previous lemma :

$$\text{VC-dim}(S) \leq 2^{\delta'+1} = \delta^{O(\delta)}.$$
 ■

→ This is very useful when shapes in \mathcal{R} are simple. If we show the dual shattering dim of S is $O(1)$, we also obtain VC-dim is $O(1)$.

§ ϵ -nets and ϵ -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 ϵ -net for A if N has a nonempty intersection with every set $R \in \mathcal{R}$ such that $|R \cap A| \geq \epsilon |A|$.

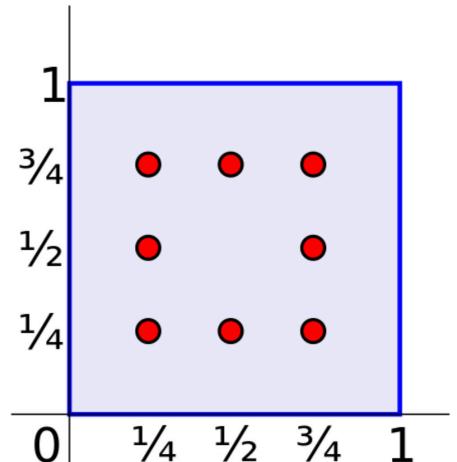
\uparrow $X \setminus A$ gets 0
p. m.e.s,
rest is uniform

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

$$\forall R \in \mathcal{R}, \quad \Pr_D(R) \geq \epsilon \Rightarrow R \cap N \neq \emptyset.$$

Here, $\Pr_D(R)$ 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 ϵ -net with $\epsilon = 1/4$ of the unit square in the range space where the ranges are closed filled rectangles.

- The minimum sample size that contains an ϵ -net (or ϵ -sample) can be bounded in terms of the VC-dimension of the range space.

§ The ε -net theorem.

Naive approach:

Let (X, \mathcal{R}) be a range space with $\text{VC-dim } d \geq 2$, let $A \subseteq X$ with $|A| = n$.
then there exists a combinatorial ε -net N for A of size at most $\lceil d \ln n / \varepsilon \rceil$.

Proof: Let \mathcal{R}' be the projection of \mathcal{R} on A .

By Sauer-Shelah theorem, $|\mathcal{R}'| \leq n^d$. \otimes

We take a sample of $k = \lceil d \ln n / \varepsilon \rceil$ points of A , independently & uniformly at random.

For each set $S \in \mathcal{R}$ with $|S \cap A| \geq \varepsilon |A|$, there is a corresponding set $S' \in \mathcal{R}'$.

So a point in our sample is in S' with prob $\geq \varepsilon$.

Then $\Pr[\text{our sample misses a given set } S'] \leq (1 - \varepsilon)^k$.

There are n^d such sets to consider. [From \otimes]

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 $S' \in \mathcal{R}'$.

Hence, N is an ε -net for A . ■

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

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

Proof:

Let M be a set of m independent samples from X acc. to \mathcal{D} .

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

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

For this, we go beyond the union bound approach.

Choose a second set T of m indep. samples from X acc. to \mathcal{D} .

Define, $E_2 := \{\exists S \in \mathcal{R} \mid \Pr_{\mathcal{D}}(S) \geq \varepsilon, |S \cap M| = 0, |S \cap T| \geq \frac{\varepsilon m}{2}\}$.

Following lemma shows E_1 & E_2 have similar probabilities.

Lemma 1: for $m \geq 8/\varepsilon$, $\Pr(E_2) \leq \Pr(E_1) \leq 2 \Pr(E_2)$.

Proof: If event E_1 holds, there is some S' s.t. $|S' \cap M| = 0$ and $\Pr_{\mathcal{D}}(S') \geq \varepsilon$.

↑ trivial as $E_2 \subseteq E_1$.

$$\text{Hence, } \frac{\Pr(E_2)}{\Pr(E_1)} = \frac{\Pr(E_1 \cap E_2)}{\Pr(E_1)} = \Pr(E_2 | E_1) \geq \Pr(|T \cap S'| \geq \varepsilon m/2). \quad \dots \textcircled{1}$$

Now for a fixed range S' and a random sample T , the random variable $|T \cap S'|$ has a binomial distr: $B(m, \Pr_{\mathcal{D}}(S'))$.

As $\Pr_{\mathcal{D}}(S') \geq \varepsilon$, using Chernoff bounds:

$$\Pr(|T \cap S'| < (1-\delta) \mathbb{E}[|T \cap S'|]) \leq e^{-\frac{\delta^2}{2} \cdot \mathbb{E}[|T \cap S'|]} \quad [\because m \geq 8/\varepsilon]$$

$$\Rightarrow \Pr(|T \cap S'| < (1-\frac{1}{2})m\varepsilon) \leq e^{-\frac{1}{4} \cdot \frac{1}{2} \cdot m\varepsilon} = e^{-m\varepsilon/8} \leq e^{-1} < \frac{1}{2} \quad \dots \textcircled{2}$$

Hence, from $\textcircled{1}$ & $\textcircled{2}$, $\Pr(E_1) \leq 2 \Pr(E_2)$. •

Now we bound $\Pr[E_2]$ by prob. of a larger event E'_2 :

$$E'_2 := \{ \exists S \in \mathcal{R} \mid |S \cap M| = 0 \text{ and } |S \cap T| \geq \varepsilon m/2 \}.$$

Lemma 2. $\Pr(E_1) \leq 2 \Pr(E_2) \leq 2 \Pr(E'_2) \leq 2 (2m)^d 2^{-\varepsilon m/2}$.

Proof: 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 $S \in \mathcal{R}$, $K = \varepsilon m/2$, let

$$E_S := \{ |S \cap M| = 0, |S \cap T| \geq K \}.$$

This event means $(M \cup T) \cap S \geq K$, but all these elements were placed in T and not in M .

So, out of $\binom{2m}{m}$ possible partitions of $M \cup T$, we choose one of $\binom{2m-K}{m}$ partitions where no element of S is in M .

$$\begin{aligned}
\text{Hence, } \mathbb{P}(E_S) &\leq \mathbb{P}(|M \cap S| = 0 \mid |S \cap (\text{MUT})| \geq k) \\
&= \binom{2m-k}{m} / \binom{2m}{m} \\
&= (2m-k)! m! / (2m)! (m-k)! \\
&= \frac{m(m-1)\dots(m-k+1)}{(2m)(2m-1)\dots(2m-k+1)} \leq 2^{-k} \leq 2^{-\varepsilon m/2}.
\end{aligned}$$

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

Hence, using union bounds:

$$\mathbb{P}(E_2') \leq (2m)^d 2^{-\varepsilon m/2}.$$

To complete the proof of ε -net theorem, we need to show $\mathbb{P}(E_1) \stackrel{\text{lem}}{\leq} 2(2m)^d 2^{-\varepsilon m/2} \leq \delta$, for $m \geq \frac{8d}{\varepsilon} \ln \frac{16d}{\varepsilon} + \frac{4}{\varepsilon} \ln \frac{2}{\delta}$.

equ., we need $\varepsilon m/2 \geq \ln(2/\delta) + d \ln(2m)$.

$$\text{As } m \geq \frac{4}{\varepsilon} \ln \frac{2}{\delta}, \quad \varepsilon m/4 \geq \ln(2/\delta).$$

To finish we show $\varepsilon m/4 \geq d \ln(2m)$.

[Fact: if $y \geq x \ln x \geq e$, then $2y/\ln y \geq x$]

use $y = 2m \geq \frac{16d}{\varepsilon} \ln \frac{16d}{\varepsilon}$, $x = \frac{16d}{\varepsilon}$, we have

$$\frac{4m}{\ln(2m)} \geq \frac{16d}{\varepsilon} \Rightarrow \frac{\varepsilon m}{4} \geq d \ln(2m).$$

§ Application :

Probably Approximately Correct (PAC) Learning.

We are given a set of items X and
a prob distr. \mathcal{D} is defined on X .

$$\begin{array}{c} +1 \quad -1 \\ \oplus \quad \ominus \\ -1 \quad +1 \\ \ominus \quad \ominus \end{array}$$

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

The concept class \mathcal{C} is the set of all possible
classifications defined by the problem.

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

Definition 14.9 [PAC Learning]: A concept class \mathcal{C} over input set X is PAC learnable¹
if there is an algorithm L , with access to a function $\text{ORACLE}(C, \mathcal{D})$, that satisfies the
following properties: for every correct concept $C \in \mathcal{C}$, 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
 $\text{ORACLE}(C, \mathcal{D})$ is polynomial in ϵ^{-1} and δ^{-1} , and with probability at least $1 - \delta$ the
algorithm L outputs a hypothesis h such that $\Pr_{\mathcal{D}}(h(x) \neq c(x)) \leq \epsilon$.

\downarrow
approximately probably

Theorem: Any finite concept class \mathcal{C} can be PAC-learned with $m = \frac{1}{\epsilon}(\ln |\mathcal{C}| + \ln \frac{1}{\delta})$ samples.

Proof: Let $c^* \in \mathcal{C}$ be the correct classification.

A hypothesis h is "bad" if $\Pr_x[h(x) \neq c^*(x)] \geq \epsilon$.

$$\begin{aligned} \Pr[\text{a bad } h \text{ is consistent with } m \text{ random samples}] \\ \leq (1 - \epsilon)^m. \end{aligned}$$

Using union bound,

$$\begin{aligned} \Pr[\exists \text{ bad } h \text{ is consistent with } m \text{ random samples}] \\ \leq |\mathcal{C}|(1 - \epsilon)^m \\ \leq \delta \quad [\text{as } m = \frac{1}{\epsilon}(\ln |\mathcal{C}| + \ln \frac{1}{\delta})]. \end{aligned}$$

So, we return whatever h is consistent with all m random sample. — all these are "non-bad".

By assumption, as correct classification $c^* \in \mathcal{C}$, at least one such h exists. ■

→ Note that X can be infinite.

So, it is interesting that we can PAC-learn \mathcal{C} , 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, $1/\epsilon$, $1/\delta$.

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 $|\mathcal{C}|$?
- 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} : $\mathcal{C} = \{[x, y] \mid x \leq y\} \cup \{\emptyset\}$.

Let $c^* \in \mathcal{C}$ be the concept to be learned.
 h be the hypothesis returned by our algo.

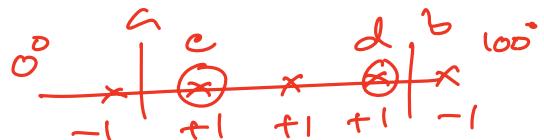
Training set T is collection of n points drawn from D .

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



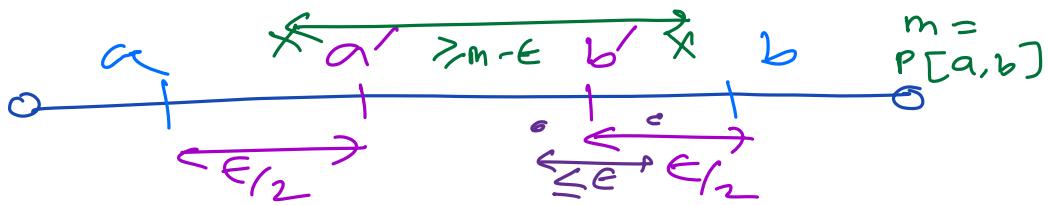
Q. What is the prob ALGO makes an error?

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

Case 1. $\Pr_D(x \in [a, b]) \leq \epsilon$.

From the above fact, prob of error $\leq \epsilon$

Case 2. $\Pr_{\mathcal{D}}(x \in [a, b]) > \epsilon$.



So, $a' \leq b'$.

For simplicity assume $a' < b'$.

If ALGO returns a bad hypothesis then error $\geq \epsilon$.
 Now if sample points fell in $[a, a']$ and $[b, b']$ then we would have returned $[c, d] \supseteq [a', b']$. \rightarrow a good hypothesis

So, $\Pr[\text{bad hypothesis}] \leq \Pr[\text{sample points didn't fall in } [a, a'] \text{ or } [b, b']]$.
 either

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

$$\leq 2 \left(1 - \frac{\epsilon}{2}\right)^n \leq 2 e^{-\epsilon n/2} \underset{\uparrow}{\cdot} \leq \delta.$$

by choosing $n \geq 2 \ln(2/\delta)/\epsilon$.

- ϵ -net & VC-dim generalizes this idea.

Theorem 14.13: Let \mathcal{C} be a concept class that defines a range space with VC dimension d . For any $0 < \delta, \epsilon \leq 1/2$, there is an

over a set of
points X

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

such that \mathcal{C} is PAC learnable with m samples.

Proof: Let $c \in \mathcal{C}$ be the correct classification.

For $c' \in \mathcal{C}$, $\Delta(c', c) = \{x \mid c(x) \neq c'(x)\}$.

$\Delta(c) := \{\Delta(c', c) \mid c' \in \mathcal{C}\}$ collection of all possible sets of points of disagreement with the correct classification.

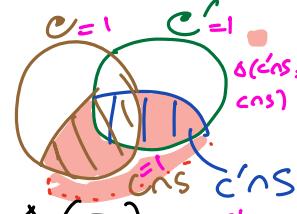
Surprising lemma!
motivate more!

Lemma: $\text{VC-Dim}(X, \Delta(c)) = \text{VC-Dim}(X, \mathcal{C})$.

for any $c \in \mathcal{C}$!

→ For any $S \subseteq X$, let $\mathcal{C}_S, \Delta(c)_S$ be projection of (X, \mathcal{C}) and $(X, \Delta(c))$ on S .

We define bijection $b: \mathcal{C}_S \rightarrow \Delta(c)_S$ by mapping $c' \cap S \in \mathcal{C}_S$ to $\Delta(c' \cap S, c \cap S) \in \Delta(c)_S$.



Then for $S \subseteq X$, $|\mathcal{C}_S| = |\Delta(c)_S|$, and S is shattered by \mathcal{C} iff it is shattered by $\Delta(c)$. Hence, the range spaces have same VC-dim.

HW: Show for $S \subseteq X$, S is shattered by \mathcal{C} iff S is shattered by $\Delta(c)$.

To complete the proof of lemma, we need to show b is 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.l.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 $c, c' \in \mathcal{C}$,

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 $\text{VC-Dim}[(X, \Delta(c))] = d$, there exists $m = O(d/\epsilon \ln d/\epsilon + 1/\epsilon \ln 1/\delta)$ s.t. sample of size $\geq m$ is ϵ -net for this range-space w.p. $\geq 1-\delta$.

Hence, w.p. $\geq 1-\delta$ 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 bad & thus can be excluded

As we only select hypothesis that are consistent w.r.t. in all the samples.

- ϵ -sample provides stronger guarantees.
- it maintains relative probability weight all sets $R \in \mathcal{R}$ within error of ϵ , and needs just additional $O(1/\epsilon)$ 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 ϵ -sample for X with respect to \mathcal{D} if for all sets $R \in \mathcal{R}$,

$$\left| \Pr_{\mathcal{D}}(R) - \frac{|S \cap R|}{|S|} \right| \leq \epsilon.$$

Relative freq.
of a range.

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 ϵ -sample for A if for all sets $R \in \mathcal{R}$,

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

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 ϵ -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 ϵ -sample for X with probability at least $1 - \delta$.

§ Application: Agnostic Learning.

In PAC learning, we assumed there is a $c^* \in \mathcal{C}$ 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 \mathcal{C} . In agnostic learning, the goal is find a nearly best classification c' s.t.

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

\downarrow
correct classif,
may not be in \mathcal{C}

If the training set define an $\epsilon/2$ -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 $\epsilon/2$.

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

Theorem 14.19: The following three conditions are equivalent:

1. A concept class \mathcal{C} over a domain X is agnostic PAC learnable.
2. The range space (X, \mathcal{C}) has the uniform convergence property.
3. The range space (X, \mathcal{C}) has a finite VC dimension.

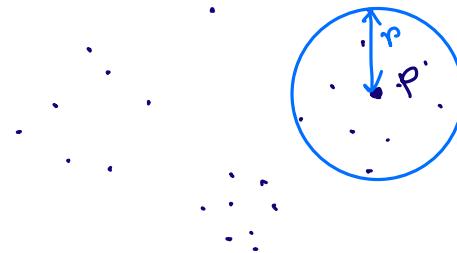
§ Applications: Data mining.

① Estimating dense neighborhoods.

Given: n points in \mathbb{R}^2 , $p = (x, y)$, $r \in \mathbb{R}$.
 \nearrow very large

Goal: What fraction of points are distance $\leq r$ from (x, y) .

[Application: opening new facility / business]



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

VC-Dim of the set of all disks = 3. \leftarrow constant

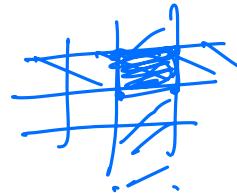
We can sample a random set of $O(\frac{1}{\epsilon^2} \ln \frac{1}{\delta})$ points and give fast approximate answers to all the queries by scanning only the sample.

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

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

- Similar example: Range searching.

(How many points are included in a query rectangle?)



② Mining frequent itemsets.

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

Goal: Find set of items that appear in $\geq \theta$ fraction of transactions.

Say we want to characterize $\text{freq} \geq \theta$ to be frequent items, $\text{freq} \leq \theta - \epsilon$ to be infrequent.

$\text{freq} : [\theta - \epsilon, \theta]$ can be ambiguous.

The number of possible transactions
= subsets of I is huge!

Even for transactions of size $\leq l$, there can be $O(|I|^l)$ of them which could be frequent.

(HW: Using Chernoff + union bound would give $\Omega\left(\frac{\theta}{\epsilon^2}(l \ln |I| + \ln \frac{1}{\delta})\right)$ samples are needed.)

ϵ -sample does better! → transactions that include s .

For each $s \subseteq I$, $T(s) := \{t \in T, s \subseteq t\}$.

let $R = \{T(s) \mid s \subseteq I\}$.

Claim: $\text{VC-Dim } [(T, R)] \leq l$.

- A transaction of size d has 2^d subsets, is therefore included in $\leq 2^d$ ranges. t is included in $T(s)$ if $s \subseteq t$

Now, consider $H \subseteq T$, $R_H = \{R \cap H \mid R \in R\}$.

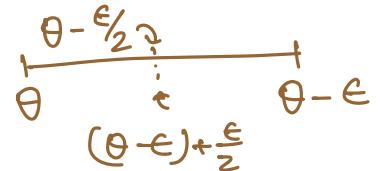
Then, $|R_H| \leq 2^l$, as H can belong to $\leq 2^l$ ranges.

Hence, no set of $\geq l$ transactions can be shattered.

$$\Rightarrow \text{VC-dim} = l.$$

Hence, by ϵ -sample theorem, w.p. $\geq 1 - \delta$, a sample of size $O\left(\frac{l}{\epsilon^2} \ln \frac{l}{\epsilon} + \frac{1}{\epsilon^2} \ln \frac{1}{\delta}\right)$ guarantee all itemsets are accurately determined to within $\epsilon/2$ 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.

- ϵ -net theorem: Let (X, \mathcal{R}) be a range space with VC-dim d and let \mathcal{D} be a prob. distribution on X . For any $0 < \delta, \epsilon \leq 1/2$, 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 ϵ -net for X with probability at least $1 - \delta$.

- $O(d \ln(d \text{OPT}))$ -approximation for hitting set with VC-dimension d .

Hitting set variant:

$$n = \# \text{ elements}, \quad m = \# \text{ sets}. \\ X := \{e_1, \dots, e_n\} \quad \mathcal{R} := \{S_1, \dots, S_m\}.$$

Algorithm:

→ Guess OPT (by binary search). $\epsilon = 1/2 \text{OPT}$.

Initialize:

→ Put $w(e_i) = 1 \quad \forall i \in [n]$. // start with uniform weights on each element

Loop:

→ Find ϵ -net N_ϵ of size $O\left(\frac{d}{\epsilon} \ln \frac{d}{\epsilon}\right)$.

→ If all sets are hit, return N_ϵ & stop.

→ Else $\exists S_j$ s.t. $S_j \cap N_\epsilon = \emptyset$ // if some set is not hit by ϵ -net

$w(e_i) = 2w(e_i) \quad \forall e_i \in S_j$ // Double weights of points in S_j

→ Goto Loop.

→ The algorithm is a variant of multiplicative weight update (MWU). Intuitively, total weight of points increases by a rate $(1+\epsilon)$, and OPT increases by a faster rate $(1 + \frac{1}{\text{OPT}}) = (1+2\epsilon)$. Thus the algorithm stops quickly w. good guarantee.

- Theorem : If \exists hitting set of size OPT , the doubling process can happen at most $O(\text{OPT} \cdot \log \frac{n}{\text{OPT}})$ times, and the total weight is at most n^4/OPT^3 .

→ Say, H be an optimal set.

For input X , say the set S_j is returned by an iteration.

$$\text{Then, } w(S_j) \leq \epsilon \cdot w(X)$$

Thus, in each iteration, $w(X)$ becomes at most $w(X) + w(S_j) \leq (1+\epsilon) w(X)$.

So, total weight of X after k iterations :

$$w(X) \leq n(1+\epsilon)^k \leq ne^{ek} \quad [\because 1+\epsilon \leq e^\epsilon] \quad \forall \epsilon > 0.$$

As H is a hitting set, $H \cap S_j \neq \emptyset$.

So, at least one element $x \in H$ is doubled in each iteration. Say, x is doubled totally z_x times.

$$w(H) = \sum_{h \in H} 2^{z_h}, \text{ where } \sum_{h \in H} z_h \geq k.$$

$$\geq (2^{2ek}) / 2^e.$$

[Here we have used convexity of exponential function.
from Jensen's inequality, $\sum_i p_i \phi(x_i) \geq \phi(\sum p_i x_i)$
where $p_i \geq 0$, $\sum p_i = 1$ & ϕ is convex.]

H is optimal i.e. $|H| = 1/2e$. Take $\phi(x) = 2^x$, $p_i = 2e^{-1}$ $\forall i \in [1, H]$.

$$\therefore \sum_{h \in H} 2^{z_h} = \frac{1}{2e} \left[\sum_{h \in H} 2e \cdot 2^{z_h} \right] \geq \frac{1}{2e} 2^{\sum 2e \cdot z_h} \geq \frac{1}{2e} 2^{2e \cdot k}.$$

As $w(H) \leq w(X)$,

$$(2^{2ek}) / 2^e \leq n e^{ek} \leq n 2^{\frac{3}{2} \cdot ek} \quad (\because e \leq 2^{\frac{3}{2}}) \approx 2.82$$

$$\Rightarrow 2^{2ek - (3ek/2)} \leq 2n e$$

$$\Rightarrow 2^{ek/2} \leq 2n e \Rightarrow ek/2 \leq \log(2n e)$$

$$\Rightarrow k \leq \frac{2}{e} \log(2n e) = O(\text{OPT} \cdot \log(n/\text{OPT})).$$

$$\therefore w(X) \leq n e^{ek} \leq n e^{2 \log(2n e)} \leq O(n^3/\text{OPT}^2). \blacksquare$$

\Rightarrow So we can stop the run if # iterations exceed k .

In special cases, if \exists small ϵ -nets of size $O(d/\epsilon)$, an $O(d)$ -approx is obtained.

- Matousek-Siedel-Welzl '90 : $\frac{1}{\epsilon} = \Theta(\text{OPT})$.
For disks in \mathbb{R}^2 , $\exists \epsilon$ -nets of size $O(1/\epsilon)$.

- LP-based approach (Even et al.) [hitting set]

Natural LP : (LP1)

$$\min \sum_{u \in X} x_u = J.$$

$$\text{s.t. } \sum_{u \in S} x_u \geq 1, \forall S \subseteq Q$$

$$x_u \geq 0, \forall u \in X$$

Equivalent LP : (LP2)

$$\max \varepsilon$$

$$\text{s.t. } \sum_{i \in S} \mu_i \geq \varepsilon, \forall S \subseteq R$$

$$\sum_{u \in X} \mu_u = 1$$

$$\varepsilon, \mu_u \geq 0, \forall u \in X$$

Equivalence proof:

$$\text{use substitution } \varepsilon = \sum_{u \in X} x_u \stackrel{\text{OPT}}{\approx}, \mu_u = \varepsilon \cdot x_u \quad \forall u \in X$$

$$\therefore J^* = 1/\varepsilon^*.$$

Algorithm:

1. Solve LP2 to obtain μ^*, ε^* .
2. Find ε^* -net H with $\text{weight}(u) = \mu_u, \forall u \in X$.

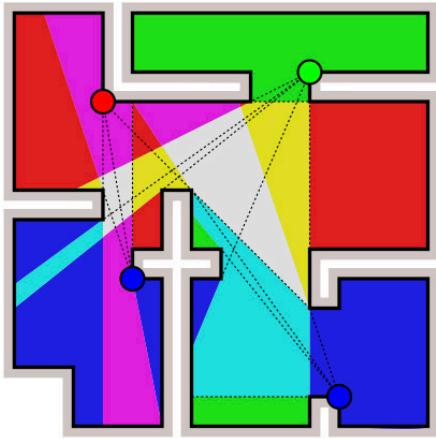
As $\sum_{i \in S} \mu_i \geq \varepsilon, \forall S \subseteq R$, H is a hitting set.

→ Can be extended to weighted setting as well.

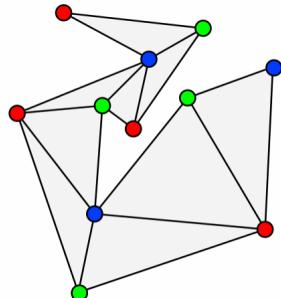
Choose $\varepsilon = 1/\text{OPT}$, then $\text{cost(Algo)} = d \text{OPT} \ln(d \text{OPT})$
 i.e. we get $d \ln(d \text{OPT})$ approximation.

• Application : Art Gallery Theorem.

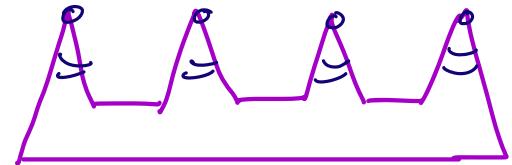
The **art gallery problem** or **museum problem** is a well-studied **visibility problem** in **computational geometry**. It originates from a real-world problem of guarding an **art gallery** with the minimum number of guards who together can observe the whole gallery. In the geometric version of the problem, the layout of the art gallery is represented by a **simple polygon** and each guard is represented by a **point** in the polygon. A set S of points is said to guard a polygon if, for every point p in the polygon, there is some $q \in S$ such that the **line segment** between p and q does not leave the polygon.



Four cameras cover this gallery.



A 3-coloring of the vertices of a triangulated polygon. The blue vertices form a set of three guards, as few as is guaranteed by the art gallery theorem. However, this set is not optimal: the same polygon can be guarded by only two guards.



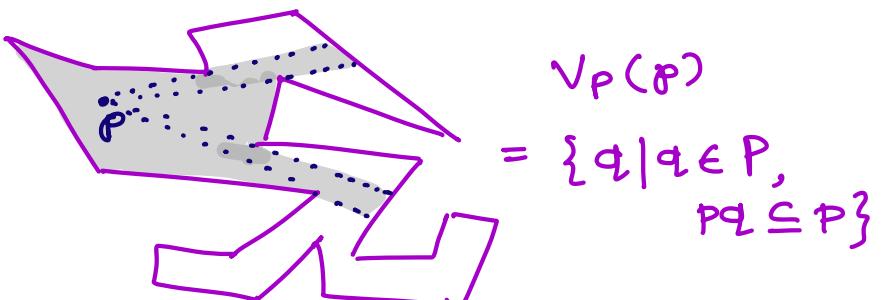
What about approximation?

- Consider the range space $S = (P, R)$ where R is the set of all possible visibility polygons inside P .

• **Theorem :**

$$\text{VC-dim}(S) = O(1).$$

[Ch 6.4, Har-Peled].



- We want to cover the entire polygon using min # of visibility polygons.

This is just geometric set cover.

Using prev algorithms we obtain $O(\log \text{OPT})$ -approximation.

