

LP - techniques

Set-Cover: NP-hard to approximate beyond $c \log n$.

- $O(\log n)$ -approximation algorithms.

Independent Set: NP-hard to approximate beyond $n^{1-\epsilon}$.

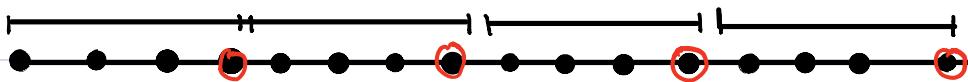
- We obtain better algorithms for geometric settings.
- Our algorithms will be via LP-rounding.
- Before we do this, we introduce a seemingly unrelated combinatorial problem, called ϵ -nets.
- These have applications in data structures, etc.

ε -nets

- Let (X, \mathcal{R}) be a set system. $|X|=n$.
- For $\varepsilon > 0$, an ε -net is a set $Y \subseteq X$, such that
 $R \cap Y \neq \emptyset \quad \forall R \in \mathcal{R} : |R| > \varepsilon n$.

ie: Y is a subset of points that "hits" all large subsets.

- Example: Let X be a set of n points on \mathbb{R} .
Let \mathcal{R} be a disjoint, contiguous sequence of εn points:



- Here $n = 16$.
- Let $\varepsilon = \frac{1}{4}$.
- Each interval drawn above is therefore "large".
- Choosing every 4th point yields an ε -net for this instance.
- $|Y| = 4$.
- The example also shows that we require at least $1/\varepsilon$ points.

More generally, let X be a set of n points on the real line, and let R be a collection of intervals.

For $\varepsilon > 0$, how can we construct an ε -net?

Construct Y by picking every $\lfloor \frac{1}{\varepsilon n} \rfloor^{\text{th}}$ point from X .

Every interval of length $\geq \varepsilon n$ will be hit by Y .

We obtain an ε -net of size $\left\lceil \frac{1}{\varepsilon} \right\rceil$

Note that the size of the ε -net is independent of n .

Let us try to construct an ε -net for other geometric set systems:

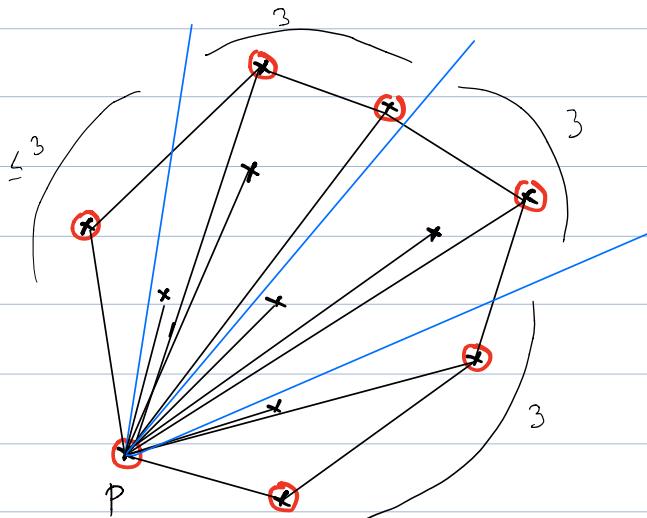
Let $X = n$ points in \mathbb{R}^2

$R = \text{half-spaces in } \mathbb{R}^2$.

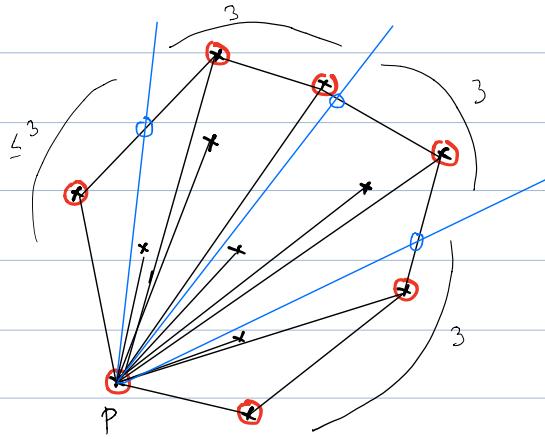


- The circled points are on the convex hull of X .
- Every halfspace containing $\geq \varepsilon n$ points must be hit.

- Start with an arbitrary point, P , on the convex hull, and sort the points of X radially with respect to P .



- Partition into groups containing $\lceil \varepsilon n \rceil - 1$ points in this order.
- # Groups $\leq n / (\lceil \varepsilon n \rceil - 1) \approx 1/\varepsilon$



Pick P into the ϵ -net; and restrict attention to those half-spaces not containing P . Then,

Obs 1:

- Any half-space containing $\geq \epsilon n$ points must contain points from at least 2 groups.

Obs 2:

- Any half-space containing $\geq \epsilon n$ points, & thus points from 2 groups must contain the point intersecting the ray defining some group boundary, & the convex hull (depicted as blue circles in the figure).

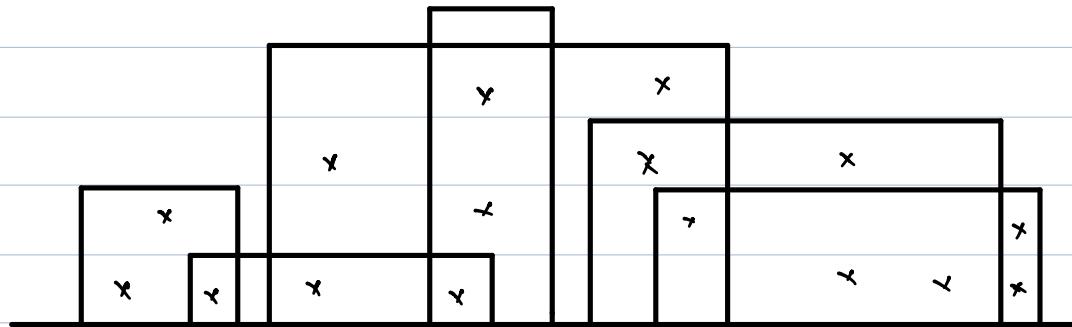
Obs 3: For each blue point, add the two vertices of the edge of the convex hull on which it lies.

This forms an ϵ -net of size at most : $\lceil 2/\epsilon \rceil + 1$.

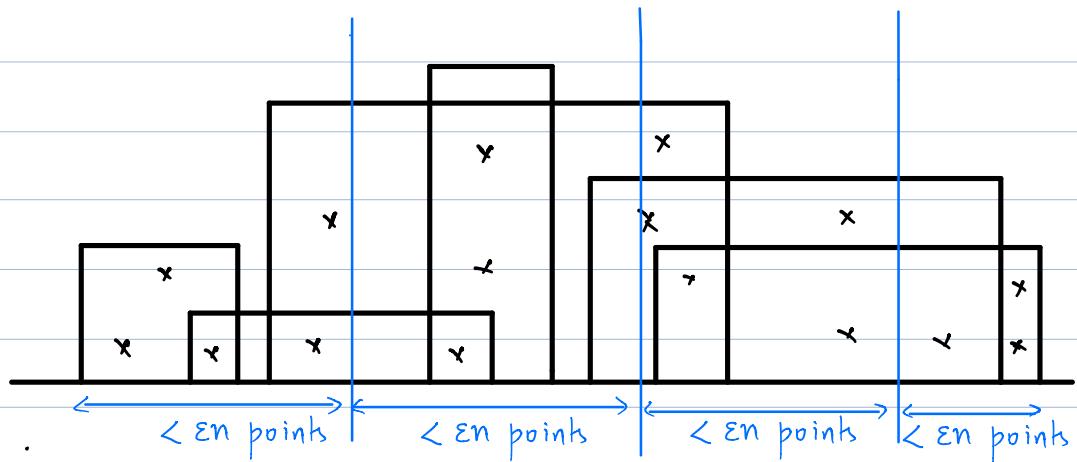
Example:

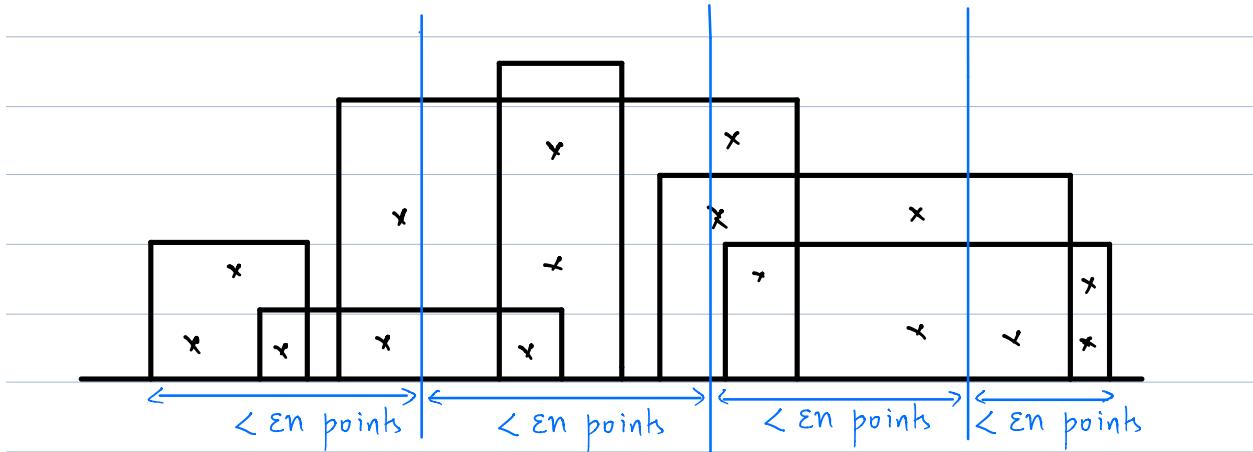
Let $X =$ points in \mathbb{R}^2 .

$R =$ Sky-lines / rectangles on the x -axis.

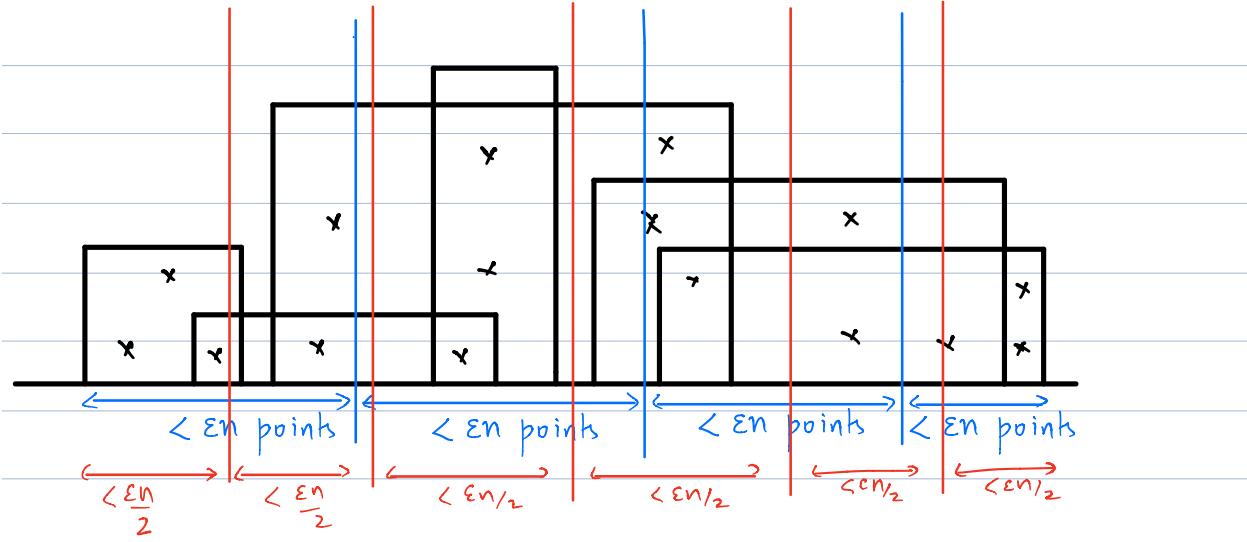


- Let us use the earlier idea of splitting into slabs, each containing $< \epsilon n$ points.

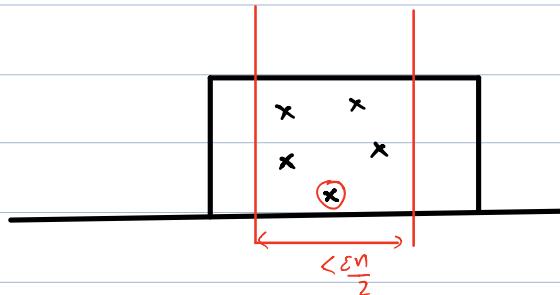




- Any slab containing $\geq \epsilon n$ points cannot lie entirely in a single slab & must cross a slab boundary.
 - Now it is not clear how to proceed.
 - Suppose we make a finer partition: each slab containing $< \frac{\epsilon n}{2}$ points:



Obs 1: Every rectangle containing $>\epsilon n$ points must contain a point from slab it completely crosses.



Claim: Choosing the lowest point in each slab is an ϵ -net of size $\left\lceil \frac{2}{\epsilon} \right\rceil$.

Random Sampling

• The techniques we developed to compute ϵ -nets thus far seem very specialized to the geometric objects.

• They do not seem to generalize for objects such as disks, or rectangles, or other more complex shapes.

• We now revisit the ϵ -net for intervals & give a different construction of an ϵ -net that generalizes.

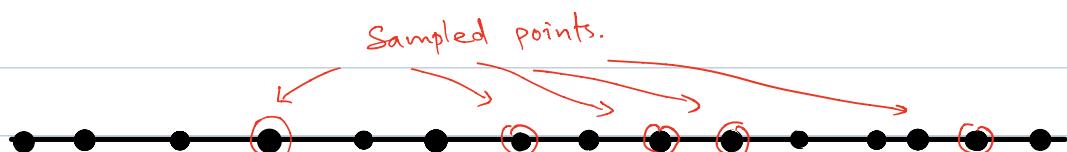
• Given: $X = \text{Points on } \mathbb{R}$.

$R = \text{Intervals of } \mathbb{R}$.

• We give a randomized construction.

• Pick each point in X uniformly & independently with probability p (to be determined later).

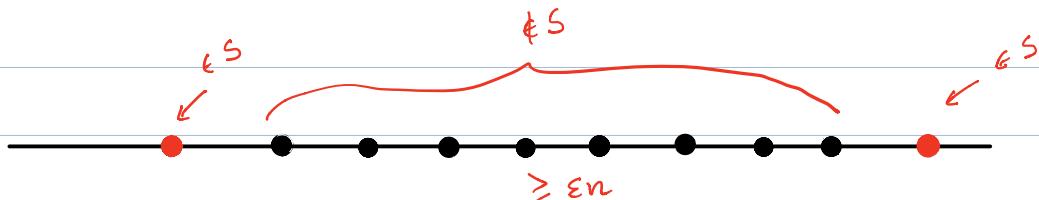
• Let S be the set of sampled points.



The expected number of points sampled is:

$$\mathbb{E}[|S|] = np \quad \text{--- (1)}$$

- Note that if we hit every interval of length ϵn , we obtain an ϵ -net for (X, \mathcal{S}) .
- The sample S fails to be an ϵ -net, if there is a contiguous sequence of at least ϵn points between two consecutive sampled points of S :



- Call such an empty interval a canonical interval.
- The probability of obtaining such a canonical interval is that
 - (i) the two extreme points of the defined interval are chosen in S .
 - (ii) No point between them is chosen in S .
- The probability that we get an empty interval of length t is: $p^2(1-p)^t \rightarrow$ none of the points in-between are chosen in S .
↓ the two end-points
are chosen in S

- For $t \geq \varepsilon n$, this probability is:

$$p^2(1-p)^t \leq p^2 e^{-pt} \leq p^2 e^{-p\varepsilon n} \quad (1)$$

- The expected number of empty intervals of length at least εn is therefore:

$$n^2 p^2 e^{-p\varepsilon n} \quad (2)$$

- Choose $p = \frac{c}{\varepsilon n}$.

Then, the expected number of empty intervals is at most:
(plugging in the value of p into (2))

$$\cancel{n^2} \left(\frac{c}{\varepsilon^2 n^2} \right) \exp \left(-\frac{c \varepsilon n}{\varepsilon n} \right) = \frac{c}{\varepsilon^2} e^{-c}.$$

- What we would like is the expectation to be < 1 .

- So, let us try $p = \frac{c}{\varepsilon n} \log \frac{1}{\varepsilon}$. Then, the expected number becomes:

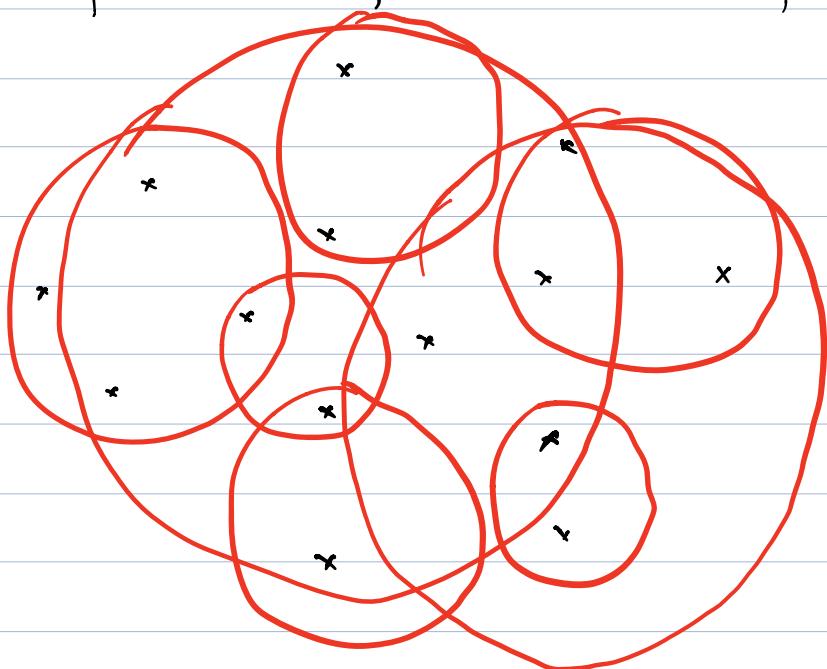
$$\begin{aligned} & \cancel{n^2} \left(\frac{c^2}{\varepsilon^2 n^2} \right) \log^2 \frac{1}{\varepsilon} \exp \left(-\frac{c \varepsilon n}{\varepsilon n} \cdot \log \left(\frac{1}{\varepsilon} \right) \right) \\ &= \frac{c^2}{\varepsilon^2} \log^2 \frac{1}{\varepsilon} \exp \left(-c \log \frac{1}{\varepsilon} \right) \end{aligned}$$

$$= \frac{c^2}{\varepsilon^2} \log^2 \frac{1}{\varepsilon} \cdot \varepsilon^c < 1 \text{ for } c \text{ large enough.}$$

- The expected size of the sample we picked is obtained by plugging the value of p into (1).
- This is: $np = O\left(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon}\right)$.

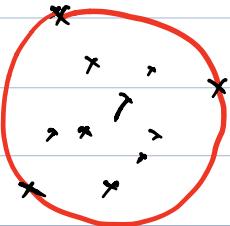
B

Let us try to use this technique to obtain ε -nets for the set system $X \subseteq \mathbb{R}^2$, $R = \text{disks in the plane}$.



- We proceed as before: Let $S \subseteq X$ be a set of points chosen uniformly & independently with probability p (to be determined later).

- If our set S hits every disk containing $\geq \varepsilon n$ points, then we obtain an ε -net.
- Just as we defined a canonical interval, as one starting at a sampled point, we can define a canonical disk.
- A canonical disk can be defined by 3 input points.



Claim: Any input disk is contained in the union of 2 canonical disks.

Obs:

- If we construct an $\varepsilon/2$ -net hitting the canonical disks, we obtain an ε -net for the original set of disks.
- Let us use ε as before.

- The probability of obtaining an empty canonical disk with t points inside is:

$p^3(1-p)^t$ \rightarrow No point in the interior is chosen.
 The three points defining the canonical disk are chosen.

- Hence, the expected number of empty canonical disks is at most:

$$n^3 p^3 e^{-p\varepsilon n}$$

- Put $p = \frac{c}{\varepsilon n} \log\left(\frac{1}{\varepsilon}\right)$, & we obtain that the expected number of empty disks is:

$$n^3 \left(\frac{c^3}{\varepsilon^3 n^3} \right) \log^3\left(\frac{1}{\varepsilon}\right) \exp\left(-c \frac{\varepsilon n}{\varepsilon n} \log \frac{1}{\varepsilon}\right)$$

$$\leq \frac{c^3}{\varepsilon^3} \log^3 \frac{1}{\varepsilon} \varepsilon^c < 1 \text{ for } c \text{ large enough.}$$

- Therefore, there is a sample S of size: $np = O\left(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon}\right)$ that is an ε -net.

- The key reason why the above argument worked for disks is that we could define n^3 canonical disks of large size.

- Suppose we have a set system, where we could likewise define n^d canonical sets, for some constant d .

Taking $p = \frac{cd}{\varepsilon n} \log\left(\frac{d}{\varepsilon}\right)$, we get that the expected number

$$\begin{aligned} \text{of empty sets of size } \geq \varepsilon n \text{ is: } & p^{\varepsilon n} n^d \exp(-\varepsilon n p) \\ &= O\left(\frac{d}{\varepsilon} \log \frac{d}{\varepsilon}\right) \end{aligned}$$

- Thus, what we require for the argument to work is that a canonical object can be determined by a constant # points.
- In fact, we will only require that for any large subset S of points, only an $O(|S|^d)$ subsets can be obtained by intersecting with R .

Shatter dimension: Let (X, R) be a set system, $Y \subseteq X$.

Let $R|_Y = \{Y \cap R : R \in R\}$. ie: the induced sets in Y .

- Suppose there exists $d \in \mathbb{N}$ such that $|R|_Y = O(|Y|^d)$, $\forall Y \subseteq X$
then we say that the shatter dimension of $(X, R) = d$.

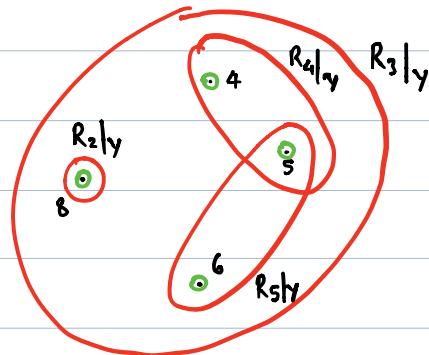
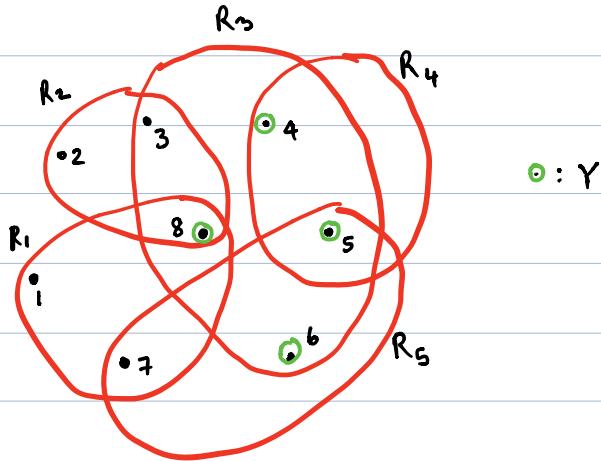
- We only require that large sets are not shattered.

We next introduce a notion that will bound the shatter dimension:

VC-dimension

- VC-dimension, or Vapnik-Chervonenkis dimension is a measure of the complexity of a set system
- The measure was introduced by Vapnik & Chervonenkis in the context of learning theory in '71.
- This notion has since found wide applicability in several areas, most notably geometry.
- Let (X, \mathcal{R}) be a set-system.
- Here we allow both X & \mathcal{R} to be infinite.
- Set-systems are also called range-spaces, & \mathcal{R} , the set of ranges.
- Let $Y \subseteq X$ be a subset of size n .
- $\mathcal{R}|_Y = \{ R \cap Y : R \in \mathcal{R} \}$. i.e. $\mathcal{R}|_Y$ is the induced set-system on Y .

- For example:



Observe that R_1 & R_2 map to the same set in $R|_Y$.

- Given a set $Y \subseteq X$, $|Y|=n$, the total number of potential subsets is $2^{|Y|}$.

- A set $Y \subseteq X$ is said to be shattered if $R_{1Y} = 2^Y$.

i.e.: the sets R induce all possible subsets of Y .

- In the example above, if we had set $Y = \{5, 8\}$, then

$$R_{1Y} = \{\emptyset, \{8\}, \{5\}, \{8, 5\}\}$$

$\nwarrow R_4$ $\swarrow R_2, R_1$ $\nearrow R_5$ $\searrow R_3$

- For a set system (X, \mathcal{R}) , the largest set $Y \subseteq X$ that can be shattered is the VC-dimension of the set system.

Example: $X = \mathbb{R}$, $\mathcal{R} = \{[a, \infty), a \in \mathbb{R}\}$.

i.e: The ground set is the real line, & the sets are all semi-infinite lines to the right.

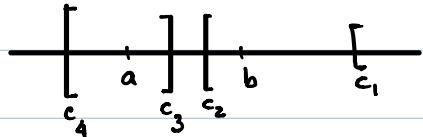
Let $Y = \{a\}$: This set can easily be shattered.

Claim: No set of size 2 can be shattered.

Thus VC-dimension of this set system = 1.

Example: $X = \mathbb{R}$, $\mathcal{R} = \{[a, \infty), (-\infty, b]\}$.

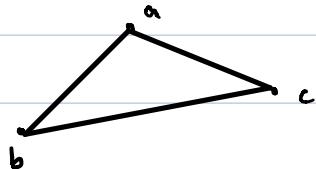
- A set of size 2 can be shattered:



Claim: No set of size 3 can be shattered.

Example: $X = \mathbb{R}^2$, \mathcal{R} = all halfspaces.

VC-dim ≥ 3 : Given 3 non-collinear points a, b, c .



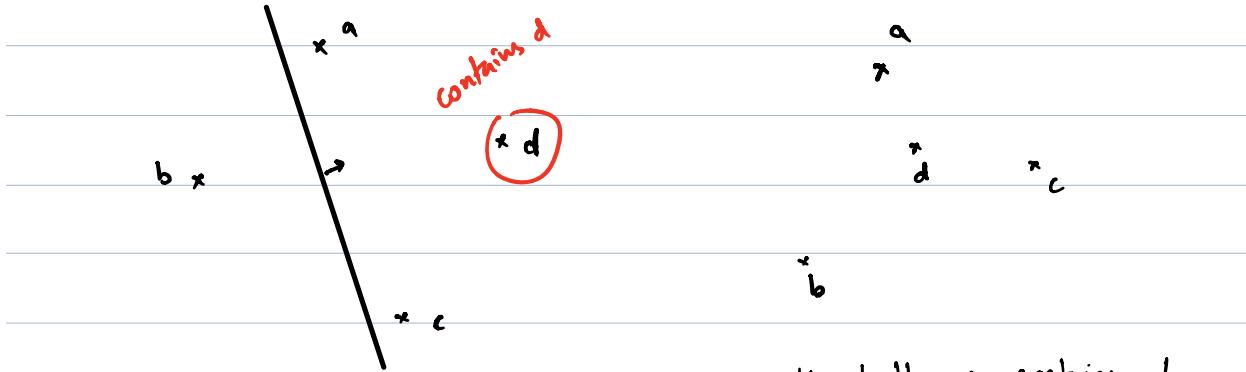
it is easy to find half-spaces forming all possible subsets.

Claim: $\text{VC-dim} \leq 3$.

Let $\{a, b, c, d\}$ be 4 points.

Either the points are in convex position, or one point, say d is contained in the convex hull of $\{a, b, c\}$.

If the four points are in convex position:



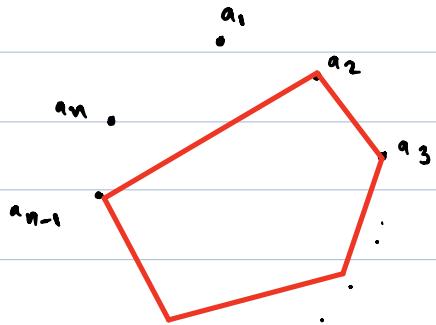
No halfspace contains only $\{d\}$.

Any halfspace containing the diagonally opposite points contains one of the other points.

Example: $X = \mathbb{R}^2$, \mathcal{R} = all convex polygons.

Claim: VC-dim. of this set system is ∞ .

Consider n points in cvx. position.



To form a set

$$\{a_{i_1}, a_{i_2}, \dots, a_{i_k}\}$$

construct a convex polygon

whose vertices are

$$\{a_{i_1}, a_{i_2}, \dots, a_{i_k}\}.$$

• Sets of any size n can be shattered.

• Hence VC-dimension = ∞ .

• In many geometric situations we will see, VC-dimension will be bounded.

Shatter function & Shatter dimension.

Consider the following recurrence:

$$\Phi_d(n) = \begin{cases} 1, & n > 0, d = 0 \\ 1, & n = 0, d \geq 0 \\ \Phi_{d-1}(n-1) + \Phi_d(n-1), & n > 0, d > 0 \end{cases}$$

Lemma: $\Phi_d(n) = \sum_{i=0}^d \binom{n}{i} = O(n^d)$, $d < n$ & $\Phi_d(n) = 2^n$ otherwise.

Lemma (Sauer-Shelah) Let $\mathcal{H} = (X, \mathcal{Q})$ be a set-system, $\dim(\mathcal{H}) = d$, $|X| = n$. Then, $|\mathcal{Q}| = O(n^d)$.

Proof: By induction on d, n . (We only show the inductive step)

Suppose the statement is true for all n & $\dim < d$, and for all $n' < n$ & all dim. d .

- Let $a \in X$ be an arbitrary element.

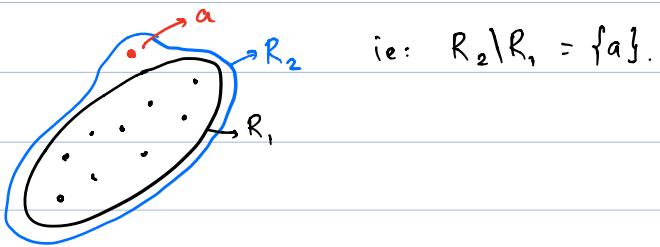
- We can define two set systems:

$\mathcal{H}_1 = (X - a, \mathcal{Q}_{-a})$, where $X - a = X \setminus \{a\}$.

$\mathcal{Q}_{-a} = \{R \setminus \{a\} : R \in \mathcal{Q}\}$.

- Note that two sets in \mathcal{R} can be mapped to the same set in \mathcal{H} .

Eg:



Let $\mathcal{R}' = \{R \in \mathcal{R} : a \in R, \text{ and } R \setminus \{a\} \in \mathcal{R}\}$, ie: all sets of type R_2 in the example above.

$$\text{Thus: } |\mathcal{R}| = |\mathcal{R}_{-a}| + |\mathcal{R}'|$$

Define $\mathcal{H}_2 = (X-a, \mathcal{R}')$.

- $\dim(\mathcal{H}_1) \leq d$.
- We show that $\dim(\mathcal{H}_2) \leq d-1$.

Suppose not: There is a set $A \subseteq X-a$, $|A|=d$, that is shattered by \mathcal{R}' .

- This implies: $A \cup \{a\}$ is shattered by \mathcal{R} .
- But this implies there is a set $A \cup \{a\}$, $|A \cup \{a\}|=d+1$ that is shattered by \mathcal{R} , contradicting the assumption that $\dim(\mathcal{H})=d$.
- Thus, the number of sets: $\Phi_d(n) = \Phi_d(n-1) + \Phi_{d-1}(n-1)$
 $= O(n^d)$

Set Cover for Set Systems with bounded VC-dimension.

- Let (X, \mathcal{R}) be a set system with VC-dimension d .
- Recall that for the Set Cover problem, we could obtain an $O(\log n)$ -approximation.
- Here, we will obtain a better approximation factor for sets of bounded VC-dimension:
 - For the sake of uniformity with the earlier presentation, we will consider instead the hitting set problem.
 - For a set system (X, \mathcal{R}) , the hitting set problem is the problem of picking $Y \subseteq X$ of smallest cardinality such that $\bigcap_{R \in \mathcal{R}} R \cap Y \neq \emptyset$.
 - We define a dual set system as $(\mathcal{R}, \mathcal{R}_X)$, where \mathcal{R} is the ground-set, and $\mathcal{R}_x = \{R \in \mathcal{R} : x \in R\}$, $x \in X$ & $\mathcal{R}_X = \bigcup_{x \in X} \mathcal{R}_x$.

- The hitting set problem is the set cover problem for the dual set system.

- Further, if the VC-dimension of (X, \mathcal{A}) is d , then the VC-dimension of $(\mathcal{R}, \mathcal{R}_X)$ is $O(2^d)$.

- Consider the LP-relaxation for the hitting set problem:

$$\min_{v \in X} \sum x_v$$

subject to,

$$\sum_{v \in S} x_v \geq 1 \quad \forall S \in \mathcal{A}$$

$$x_v \geq 0. \quad \forall v \in X.$$

- Let x be an optimal solution to the LP, and
- Let $W = OPT_{LP}$ be the optimal solution of this linear program.

- Assign a weight of $\frac{x_v}{W}$ to each $v \in X$

- Since x is a feasible solution to the LP, $\sum x_v \geq 1$.

- Thus, $\sum \frac{x_v}{w} \geq \frac{1}{w} \forall R \in \mathcal{R}$.

ie: Each set has weight $\geq \frac{1}{w}$.

- Since the VC-dimension is d , \exists an ϵ -net of size

$$O\left(\frac{d}{\epsilon} \log \frac{d}{\epsilon}\right).$$

- Put $\epsilon = \frac{1}{w}$, & we obtain a hitting set of size

$$O(dw \log dw) = O(d \text{OPT} \log \text{OPT})$$

ie: an $O(\log \text{OPT})$ -approximation.

Remark: Note that the algorithm also works in the weighted setting -

- Here, each set $R \in \mathcal{R}$ has a weight $w_R > 0$, and we want to minimize the total weight of the sets chosen that is still a cover of the points.

- This is because each set in the sample is chosen with uniform probability. [Show this]

- The constructions of ϵ -nets we saw earlier did not pick each set with uniform probability.
- Therefore, the rounding algorithm presented does not yield an $O(1/\epsilon)$ -approximation even when we have an ϵ -net of size $O(1/\epsilon)$.
- Quasi-Uniform Sampling yields an $O(1/\epsilon)$ -approx. in some settings.

Packing problems:

- Given a set of disks \mathcal{D} in the plane, disk $D \in \mathcal{D}$ has weight $w_D > 0$.

- We want to compute a maximum weight independent set.

i.e.: A set $S \subseteq \mathcal{D}$ of disks that are pairwise disjoint, and: $\sum_{D \in S} w(D)$ is as large as possible.

- Let us write an LP-relaxation for the problem:

$$\max \sum_{D \in \mathcal{D}} w_D x_D$$

subject to:

$$\sum_{\substack{D: p \in D}} w_D \leq 1 \quad \forall p \in P$$

$$w_D \in [0, 1].$$

- Let $x = (x_1, \dots, x_n)$ be an optimal LP solution.

$$\text{OPT}_{LP} = \sum_{D \in \mathcal{D}} w_D x_D.$$