

Lecture 23      11/23/2025

## Sampling in Geometric Range Spaces

Set systems arise in many applications.

A set system consists of a pair

$(P, R)$  where  $P$  is a set

and  $R$  is a collection of subsets of

$P$ . When  $P$  is finite we

have a finite set system. Sometimes

finite set systems are thought of

as hypergraphs with  $P$  as vertices

and each  $e \in R$  as a hyper edge.

Here we will be concerned with

set systems that arise in geometric

settings where  $P$  is typically  
all of  $\mathbb{R}^d$  for some dimension  $d$ ,  
or a finite subset of  $\mathbb{R}^d$ . We also  
consider  $R$  to be sets that  
are induced by "structured" shapes  
that intersect with  $P$ .

Examples of shapes include convex  
sets such as intervals, disks,  
half spaces, convex polygons etc.

In the geometric setting  $(P, R)$  are  
often called range spaces and each  
 $r \in R$  is called a range. Typically  
we associate  $r$  with a shape such as  
an interval  $I$  and ~~the~~ then



$$\mathcal{R} = \mathcal{I} \cap \mathcal{P}.$$

Geometric large spaces have additional properties that lead to a number of applications and the notion of VC-dimension,  $\epsilon$ -sample theorem,  $\epsilon$ -net theorem and others have had striking influence on many areas, in particular machine learning where the notion of VC-dimension arose.

## VC-Dimension

VC-dimension of a set system is one important measure of the complexity of a set system.

Defn: Let  $(P, R)$  be a range space.

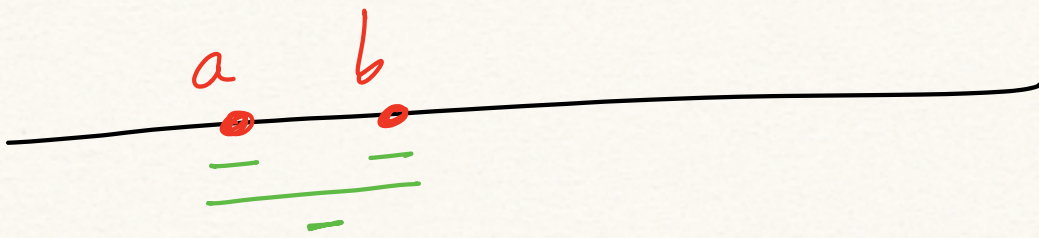
A finite subset  $Q \subseteq P$  is said to be shattered by  $R$  if  $\forall Q' \subseteq Q$

$\exists r \in R$  such that  $Q' = Q \cap r$ .

In other words  $\{Q \cap r \mid r \in R\}$   
 $= 2^Q$  the powerset of  $Q$ .

Ex: Suppose  $P$  is the real line and  $R$  is the collection of all closed intervals.





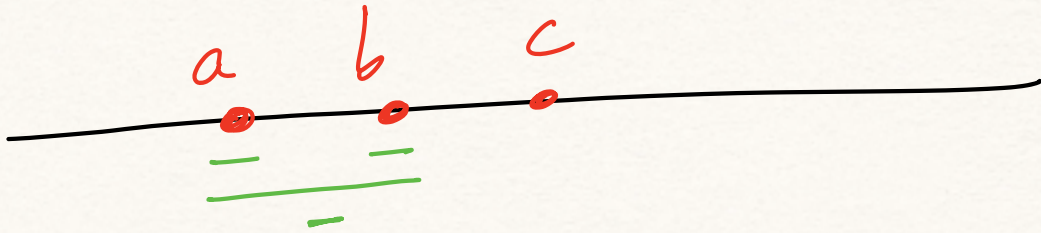
It can be seen from the figure that  $Q = \{a, b\}$  can be shattered by collection of intervals.

Defn: The VC-dimension of a set system  $(P, R)$  is the maximum cardinality of a finite set  $Q \subseteq P$  such that  $Q$  is shattered by  $R$ .

Ex: Let  $P = \mathbb{R}$  and  $R$  be the collection of intervals. Then VC-dim  $= 2$ .

Why. We saw that it is at least 2. Can it be  $\geq 3$ .

Suppose  $\mathcal{Q} = \{a, b, c\}$  where  $a < b < c$



Can we get the set  $\{a, c\}$  as an intersection of  $\{a, b, c\}$  and an interval? No.

Ex:  $P = \mathbb{R}^2$  the 2-d plane and

$\mathcal{R} = \{D \mid D \text{ is a closed disk in the plane}\}$

VC-dimension is 3.

3 points can be shattered but not 4.





Ex:  $P = \mathbb{R}^d$  and  $R =$  set of half-spaces

Recall a half space is defined by  
an inequality  $\sum_{i=1}^d a_i x_i \leq b$  for some

$$a_1, a_2, \dots, a_d, b \in \mathbb{R}$$

Claim  $VC\text{-dim} = d+1$

It is easy to see that  $VC\text{-dim} \geq d+1$

Take the  $d+1$  points  $(0, 0, \dots, 0)$

and  $(1, 0, \dots, 0), (0, 1, 0, \dots, 0) \dots (0, 0, \dots, 0, 1)$ .

This set can be shattered. Why?

However  $d+2$  points cannot be shattered  
and this follows from Radon's theorem.

## Theorem [Radon's theorem]

Let  $Q$  be a set of  $d+2$  points in  $\mathbb{R}^d$ . Then one can partition  $Q$  into  $S_1$  and  $S_2$  such that  $\text{convexhull}(S_1) \cap \text{convexhull}(S_2) \neq \emptyset$ .

The preceding theorem  $\Rightarrow Q$  cannot be shattered by half spaces.



Now that we have seen the definition of VC-dimension we state and prove a key technical lemma about set systems with bounded VC-dimension.

### Sauer's Lemma

Theorem [Sauer's Lemma]

Suppose a set system  $(P, R)$  has VC-dimension at most  $d$ . Let

$Q \subseteq P$  be a finite set of cardinality  $n$ .

Then  $\left| \left\{ Q \cap r \mid r \in R \right\} \right|$

$$\leq \sum_{i=0}^d \binom{n}{i} \leq n^d.$$

Proof: by induction on  $n$ .

If  $n=0$  it is trivial.

Let  $Q$  be a set of  $n$  points  $n > 0$ .

We can restrict attention to  $x \cap Q$

$\forall x \in R$ . Hence we can work with a finite range space. Now all we need to do is count  $|R|$ .

Fix some  $p \in Q$ .

Let  $R_1 = \{x - \{p\} \mid x \in R\}$

be the set of all ranges obtained by

removing  $p$  from the original ranges.



Suppose  $\exists r$  such that  $p \in r$  and  
 also  $r - \{p\} \in R$ . Then both  $r \cup \{p\}$   
 and  $r - \{p\}$  project to same range  
 in  $R_1$ . So to count  $|R|$  we create  
 a separate range space.

Let  $R_2 = \{r - \{p\} \mid r \cup \{p\} \in R \text{ and } r - \{p\} \in R\}$ .

From this explanation we have

Claim:  $|R| = |R_1| + |R_2|$ .

Now we consider the two range  
spaces  $(Q - \{p\}, R_1)$  and  
 $(Q - \{p\}, R_2)$ .

Claim:  $\text{VC-dim}(Q - \{p\}, R_1) \leq d$

Proof: Removing a point does not increase  
VC-dim

Claim:  $\text{VC-dim}$  of  $(Q - \{p\}, R_2) \leq d-1$ .

Proof: If  $Q' \subseteq Q - \{p\}$  is shattered  
by  $R_2$  then since every range



$\mathcal{R} \in \mathcal{R}_2$  satisfies the property that

$$\mathcal{R} \cup \{p\} \text{ and } \mathcal{R} - \{p\} \in \mathcal{R}$$

we would have  $\mathcal{Q}' \cup \{p\}$  is shattered  
by  $\mathcal{R}$ . Thus  $|\mathcal{Q}'| \leq d-1$ .  $\square$

Now by induction

$$|\mathcal{R}_1| \leq \sum_{i=0}^d \binom{n-1}{i}$$

$$\text{and } |\mathcal{R}_2| \leq \sum_{i=0}^{d-1} \binom{n-1}{i}$$

Thus

$$\begin{aligned} |\mathcal{R}| &\leq \sum_{i=0}^d \binom{n-1}{i} + \sum_{i=0}^{d-1} \binom{n-1}{i} \\ &\leq \sum_{i=0}^d \binom{n-1}{i} + \sum_{i=0}^d \binom{n-1}{i-1} \end{aligned}$$

$$\leq \sum_{i=0}^d \binom{n}{i}.$$

□.

In many settings the only way VC-dim is used is via the bound given by Sauer's lemma. So it makes sense to define the following.

Defn: The shattering dimension of a range space  $(P, R)$  is  $d$  if  $\forall Q \subseteq P$  with  $|Q| \leq n$  the size of  $R|Q \leq n^d$ .

$\text{VC-dim}(P, R) \leq d \Rightarrow \text{shattering dim}(P, R) \leq d$

Converse is also true with weaker parameter



$$\text{shattering-dim}(P, R) \leq d \Rightarrow \text{VC-dim}(P, R) \leq O(d \log d).$$

One important aspect of VC-dim is kind of closure when combining.

Theorem: Suppose  $(P, R_1)$  and  $(P, R_2)$  are range spaces with VC-dim  $d_1$  and  $d_2$  respectively. Then VC-dim of  $(P, R)$  where  $R = \{ (x_1 \cup x_2) \mid x_1 \in R_1, x_2 \in R_2 \}$  is  $O(d_1 + d_2)$ . Similarly for  $(P, R)$  where  $R = \{ x_1 \cap x_2 \mid x_1 \in R_1, x_2 \in R_2 \}$ .

## $\epsilon$ -Sampling and $\epsilon$ -net Theorems

We now discuss two theorems about how a random sample of a set from a set system  $(P, R)$  can approximate it.

For the following discussion it is useful to think of  $P$  as a finite set. Some of the concepts can be lifted to infinite sets with appropriate generalizations.



For a given system  $(P, R)$

$$\text{let } \mu(x) = \frac{|x \cap P|}{|P|} \text{ denote}$$

the measure of  $x$ .

Suppose we take a "small" random sample  $Q$  from  $P$ . Does  $Q$  preserve the measure of all  $x \in R$ ?

For this define

$$\mu_Q(x) = \frac{|x \cap Q|}{|Q|}.$$

Clearly a small sample can't touch all ranges so we need to allow some "additive" error  $\epsilon$ .

Defn: A subset  $Q \subseteq P$  is an  $\varepsilon$ -sample for  $(P, R)$  if

$$|\mu(x) - \mu_Q(x)| \leq \varepsilon \quad \forall x \in R.$$

A related notion is the following.

Defn:  $Q$  is an  $\varepsilon$ -net for  $(P, R)$  if

$$|Q \cap x| \neq \emptyset \quad \forall x \in R \text{ where } \mu(x) > \varepsilon.$$

Note that an  $\varepsilon$ -sample is automatically an  $\varepsilon$ -net.



Theorem Let  $(P, R)$  be a range space with VC-dimension  $\leq d$ .

Let  $n \geq \frac{C}{\epsilon^2} \left( d \log \frac{d}{\epsilon} + \log \frac{1}{\delta} \right)$ .

Then a random sample of  $n$  points with repetition from  $P$  is an  $\epsilon$ -sample with probability  $\geq 1 - \delta$ .

Note: The sample size does not depend on  $|P|$ . Could be infinite!

Note: The theorem relies only on the growth rate of the number of distinct ranges of a given size that follows from Sauer's lemma.

Hence the proof is not so tied to VC-dimension itself.

A stronger bound is known for  $\varepsilon$ -nets.

Theorem: Let  $(P, R)$  be a range space with  $\text{VC-dim} \leq d$ .  
Let  $l \geq \frac{C}{\varepsilon} (d \log \frac{d}{\varepsilon} + \log \frac{1}{\delta})$

Then a random sample of  $l$  points with repetition from  $P$  is an  $\varepsilon$ -net with probability  $\geq (1-\delta)$ .



## Proof of $\epsilon$ -Sample Theorem

It is a clever argument using a "double sample" argument.

First we recall the additive Chernoff bound.

Theorem: Let  $X_1, X_2, \dots, X_n \in [0, 1]$  and independent. Let  $Y = \frac{1}{n} \sum_{i=1}^n X_i$  and  $\mu = \mathbb{E}[Y]$ . Then

$$(i) \Pr[Y > \mu + \epsilon] \leq e^{-2\epsilon^2 n}$$

$$(ii) \Pr[Y < \mu - \epsilon] \leq e^{-2\epsilon^2 n}$$

To see how to use the above theorems in our setting. Fix a range  $\mathcal{R}$ .  $\mu(\mathcal{R}) = \frac{|\mathcal{R} \cap P|}{|P|}$ . Suppose we take an  $l$ -sample  $Q$  with repetition.

What is  $E[\mu_Q(\mathcal{R})]$ ?

Let  $X_i$  be indicator random variable for sample  $i$  being in  $\mathcal{R}$ .

Let  $Y = \frac{1}{l} \sum_{i=1}^l X_i$ .

$$E[X_i] = \frac{|\mathcal{R} \cap P|}{|P|} = \mu(\mathcal{R}).$$

Hence  $E[Y] = \mu(\mathcal{R})$ .

$$\text{Therefore } |E[Y] - \mu(\mathcal{R})| > \varepsilon \leq 2e^{-2\varepsilon^2 l}.$$



## Proof of $\epsilon$ -Sample Theorem

It is a clever argument using a "double sample" argument.

Let  $Q_1$  be a sample of size  $l$ .

Let  $Q_2$  be an independent sample, also of size  $l$ .

Let  $A$  be the event that

$\exists$  some range  $x$  s.t.

$$|\mu_1(x) - \mu(x)| > \epsilon \quad \text{Here } \mu_1(x)$$

$$= \mu_{Q_1}(x) \text{ for short.}$$

Let  $B$  be the event that  $\exists x$  s.t.

$$|\mu_2(x) - \mu_1(x)| > \epsilon/2$$

Claim:  $P_2[A] \leq 2 P_2[B]$ .

Let  $D$  be the event that

$\exists x$  s.t.  $| \mu_1(x) - \mu(x) | > \varepsilon$  and

$$| \mu_2(x) - \mu_1(x) | > \frac{\varepsilon}{2}$$

We have

$$\begin{aligned} P_2[B] &\geq P_2[D] = P_2[D \text{ and } A] \\ &= P_2[D|A] P_2[A]. \end{aligned}$$

We claim that  $P_2[D|A] \geq \frac{1}{2}$

which would imply that  $P_2[A] \leq 2 P_2[B]$ .

To see this, suppose event  $A$  happens.

$\Rightarrow \exists x$  such that  $| \mu_1(x) - \mu(x) | > \varepsilon$ .

For this  $x$ , via the additive Chernoff bound



$$P_2 [ |\mu_2(x) - \mu(x)| ] < \frac{\varepsilon}{2} \cdot . \text{ This is}$$

because it is a fixed  $x$  and sample is big enough. and  $\mathcal{Q}_2$  is indep of  $\mathcal{Q}_1$ .

$$\text{If } |\mu_1(x) - \mu(x)| > \varepsilon \text{ and } |\mu_2(x) - \mu(x)| < \frac{\varepsilon}{2}$$

then by triangle inequality

$$\begin{aligned} |\mu_2(x) - \mu_1(x)| &\geq |\mu(x) - \mu_1(x)| - |\mu(x) - \mu_2(x)| \\ &\geq \varepsilon - \frac{\varepsilon}{2} > \frac{\varepsilon}{2}. \end{aligned}$$

$$\Rightarrow |\mu_2(x) - \mu_1(x)| > \frac{\varepsilon}{2} \Rightarrow D \text{ happens.}$$

$$\text{Thus } P_2 [ D | A ] > \frac{1}{2}.$$

□.

Thus we can focus on event B

which is  $P_2 [ |\mu_2(x) - \mu_1(x)| ] > \frac{\epsilon}{\sqrt{2}}$

for some range  $x$ . With a factor of 2 we  
will get  $P_2 [A]$ .



To analyze  $B$  we think of the process differently. Instead of picking  $Q_1$  and  $Q_2$  independently as two separate steps we think of picking 2 elements  $Q_0$  and then splitting  $Q_0$  into two halves  $Q_1$  and  $Q_2$ .

$$P_x[B] = \sum_{Q_0} P_x[Q_0] P_x[B | Q_0] \\ \leq \max_{Q_0} P_x[B | Q_0].$$

What is  $\max_{Q_0} P_x[B | Q_0]$ ?

We think of  $Q_0$  as an arbitrary

multiset of  $2l$  points and  $Q_1$  and  $Q_2$  are obtained by even splitting of  $Q_0$  into  $l$  points each.

What is the advantage of this?

If we fix  $Q_0$  then  $R|Q_0$  has  $\leq (2l)^d$  ranges.

Thus we need to only worry about a "small" number of ranges.

For any fixed range  $r$  in  $R|Q_0$ , via the additive Chernoff bound,

$$(i) \quad \Pr \left[ |\mu_1(r) - \mu_0(r)| \geq \frac{\epsilon}{4} \right] \leq \frac{\delta}{4 (2l)^d}.$$

$$(ii) \quad \Pr \left[ |\mu_r(r) - \mu_0(r)| \geq \frac{\epsilon}{4} \right] \leq \frac{\delta}{4 (2l)^d}$$



This is because l.s.,  $\frac{C}{\varepsilon^2} (d \log \frac{d}{\varepsilon} + \log \frac{1}{\delta})$ .

Since  $|\mu_1(x) - \mu_2(x)| \leq |\mu_1(x) - \mu_0(x)| + |\mu_2(x) - \mu_0(x)|$

$$P_x \left[ |\mu_2(x) - \mu_1(x)| > \frac{\varepsilon}{2} \right] \leq \frac{\delta}{2(2d)^d}.$$

By the union bound we

$$\Rightarrow P_x [B] \leq \frac{\delta}{2}.$$

$$\Rightarrow P_x [A] \leq \delta.$$

$$\Rightarrow P_x [\bar{A}] \geq 1 - \delta.$$

□.

