



More Learning Theory

Rishabh Iyer

University of Texas at Dallas

Based on the slides of Vibhav Gogate, David Sontag, and Nick Rouzzi

Last Time

- Probably approximately correct (PAC)
 - The only reasonable expectation of a learner is that with high probability it learns a close approximation to the target concept
 - Specify two small parameters, $0 < \epsilon, 0 < \delta < 1$
 - ϵ is the error of the approximation
 - $(1 - \delta)$ is the probability that, given M i.i.d. samples, our learning algorithm produces a classifier with error at most ϵ

$\text{test-error} - \text{train error} \leq \epsilon$ $M = \text{Min number of samples req}$
 to ensure test-error $\leq \epsilon$ with
 prob $1 - \delta$

Learning Theory

ken - Bo



- We use the observed data in order to learn a classifier
- Want to know how far the learned classifier deviates from the (unknown) underlying distribution
 - With too few samples, we will with high probability learn a classifier whose true error is quite high even though it may be a perfect classifier for the observed data
 - As we see more samples, we pick a classifier from the hypothesis space with low training error & hope that it also has low true error
 - Want this to be true with high probability – can we bound how many samples that we need?

Generalization Bound

$$\epsilon^{\text{test}} \leq \epsilon^{\text{train}} + \sqrt{\frac{1}{2M} \left[(\log |H| + \log \frac{1}{\delta}) \right]}$$

Bias low if $M \gg \log |H|$
Variance

w. p $1-\delta$

- * test & train data points belong to same underlying distribution
- * samples are iid

- What we proved last time:

$$\text{train-error} = \emptyset$$



Theorem: For a finite hypothesis space, H , with M i.i.d. samples, and $0 < \epsilon < 1$, the probability that any consistent classifier has true error larger than ϵ is at most $|H|e^{-\epsilon M}$



- We can turn this into a sample complexity bound

Sample Complexity

- Let δ be an upper bound on the desired probability of not ϵ -exhausting the sample space
 - The probability that the version space is not ϵ -exhausted is at most $|H|e^{-\epsilon M} \leq \delta$
 - Solving for M yields

$$M \geq -\frac{1}{\epsilon} \ln \frac{\delta}{|H|}$$

$$= \left(\ln |H| + \ln \frac{1}{\delta} \right) / \epsilon$$

$a_{\text{test}} \leq \frac{1}{M} \left[\log |H| + \log \frac{1}{\delta} \right]$

PAC Bounds

Not nec. zero-train error



Theorem: For a finite hypothesis space H , M i.i.d. samples, and $0 < \epsilon < 1$, the probability that true error of any of the best classifiers (i.e., lowest training error) is larger than its training error plus ϵ is at most $|H|e^{-2M\epsilon^2}$

- Sample complexity (for desired $\delta \geq |H|e^{-2M\epsilon^2}$)

$$M \geq \left(\ln|H| + \ln \frac{1}{\delta} \right) / 2\epsilon^2$$

Sample Complexity
Bound

\mathcal{E} = gap between test & train error

PAC Bounds [Gen-Bound]



- If we require that the previous error is bounded above by δ , then with probability $(1 - \delta)$, for all $h \in H$

$$\epsilon_h \leq \epsilon_h^{train} + \sqrt{\frac{1}{2M} \left(\ln |H| + \ln \frac{1}{\delta} \right)}$$

- For small $|H|$
 - High bias (may not be enough hypotheses to choose from) ↪ Train error is large (underfitting)
 - Low variance

PAC Bounds



- If we require that the previous error is bounded above by δ , then with probability $(1 - \delta)$, for all $h \in H$

$$\epsilon_h \leq \epsilon_h^{train} + \sqrt{\frac{1}{2M} \left(\ln |H| + \ln \frac{1}{\delta} \right)}$$

The diagram illustrates the decomposition of the PAC bound. It shows the total error ϵ_h as the sum of two terms: ϵ_h^{train} and a square root term. The square root term is further decomposed into two components: "bias" (the first part under the square root) and "variance" (the second part under the square root). Brackets below the terms are labeled "bias" and "variance".

- For large $|H|$
 - Low bias (lots of good hypotheses) [train - error \downarrow]
 - High variance \leftarrow Test error could be high

VC Dimension

[1972]

Vapnic - C' 1972



- Our analysis for the finite case was based on $|H|$ finite
- If H isn't finite, this translates into infinite sample complexity
- We can derive a different notion of complexity for infinite hypothesis spaces by considering only the number of points that can be correctly classified by some member of H
- We will only consider the binary classification case for now

d
OCEM

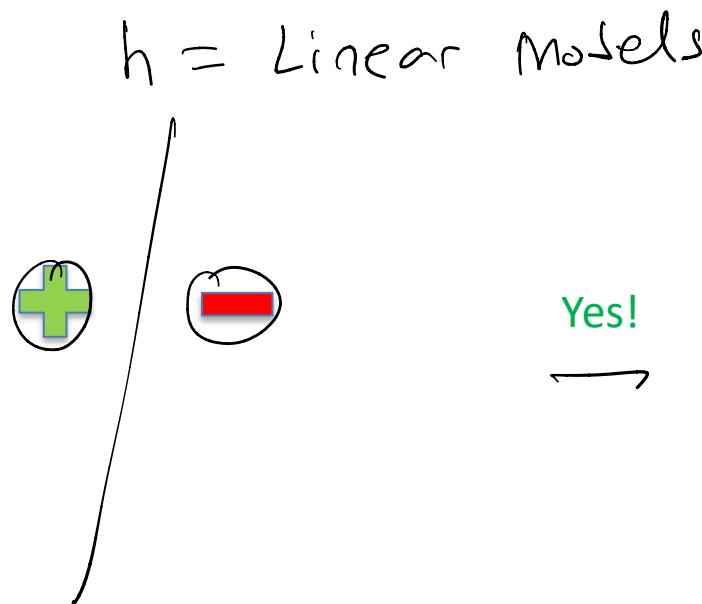
- Any Lin model
- DT / NW
- NN

$H \rightarrow$ Linear model
 $|H| =$ All possible values of θ

VC Dimension

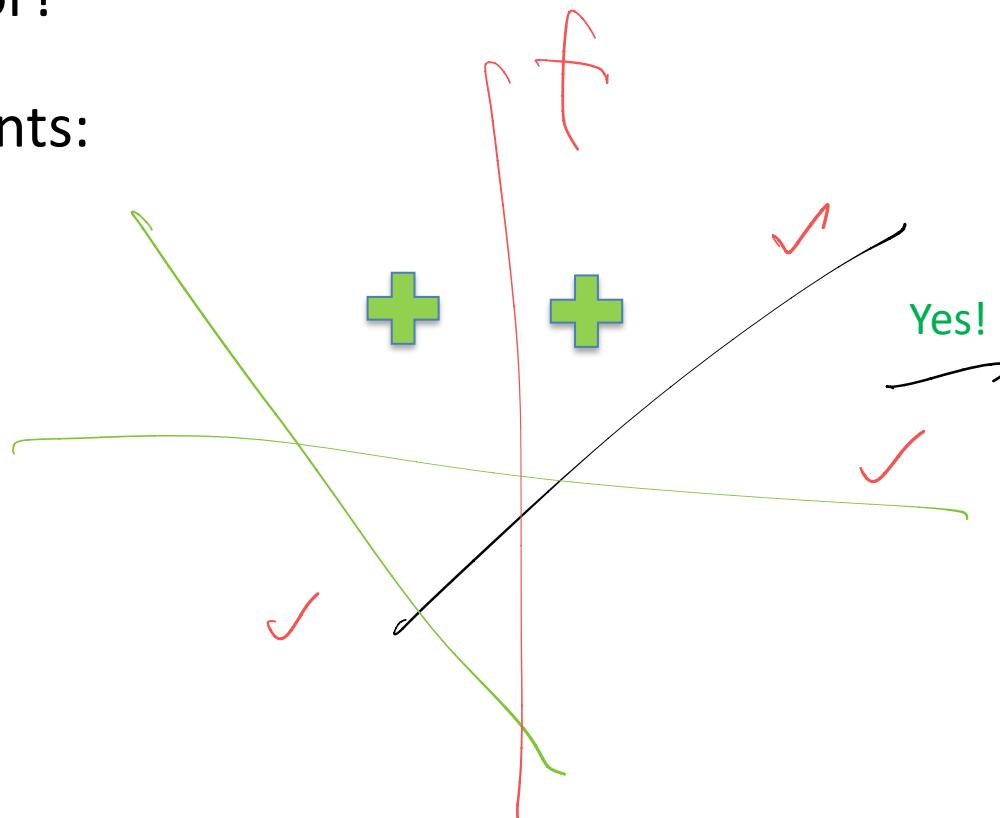


- How many points in 1-D can be correctly classified by a linear separator?
- 2 points:



VC Dimension

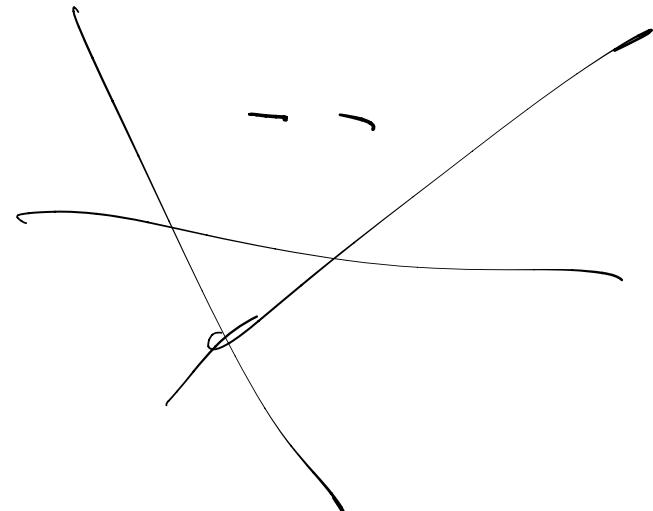
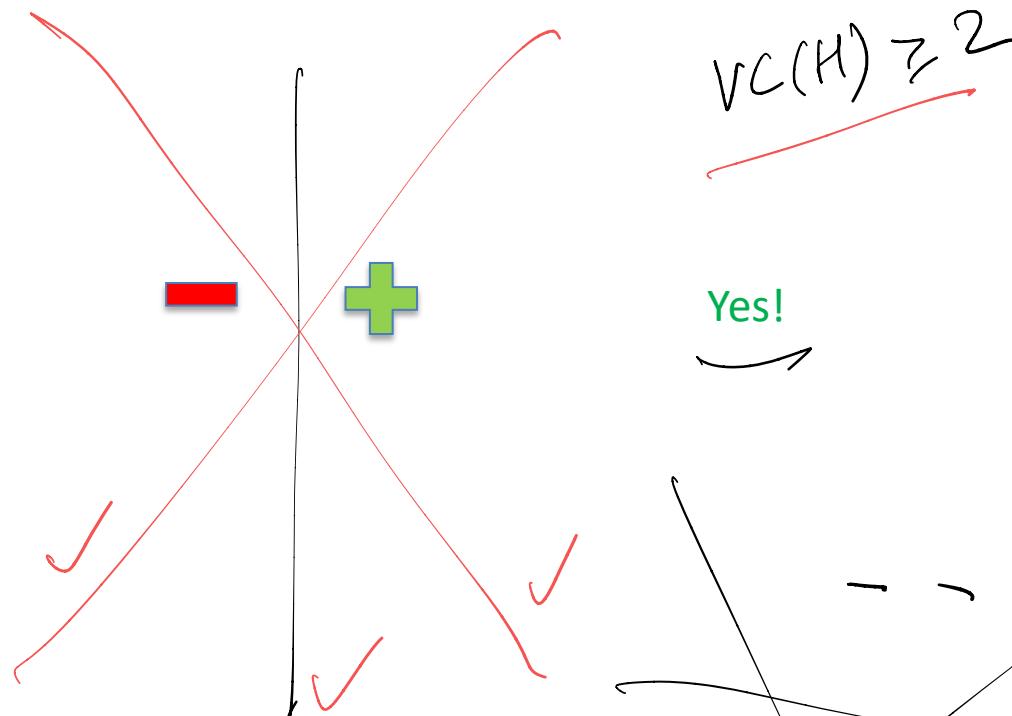
- How many points in 1-D can be correctly classified by a linear separator?
 - 2 points:



VC Dimension

- How many points in 1-D can be correctly classified by a linear separator?

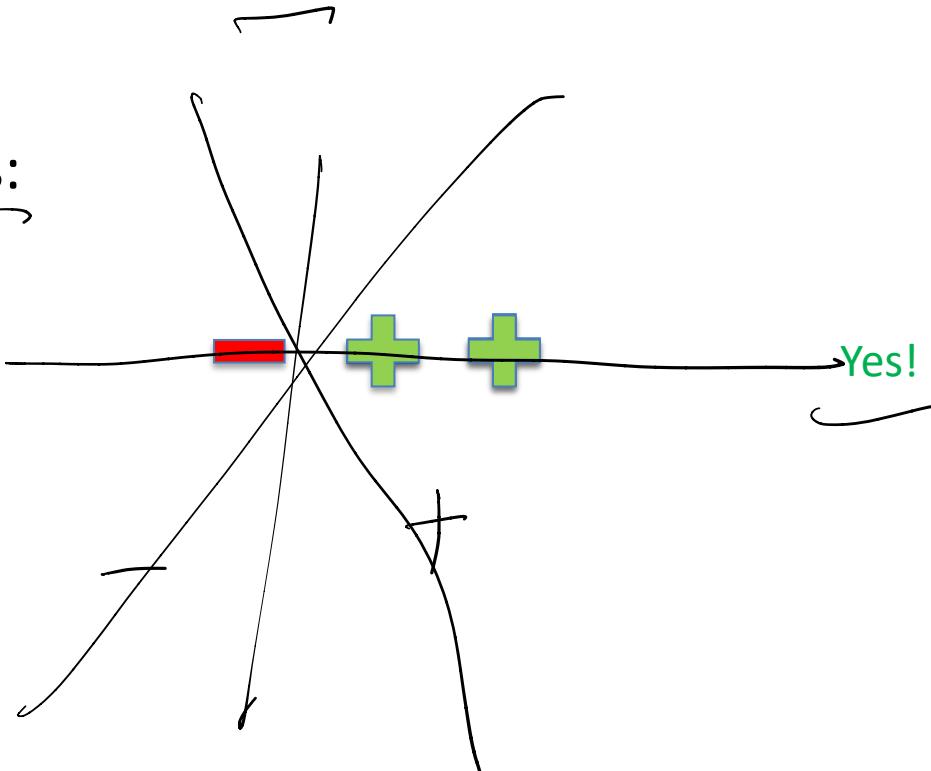
- 2 points:



VC Dimension

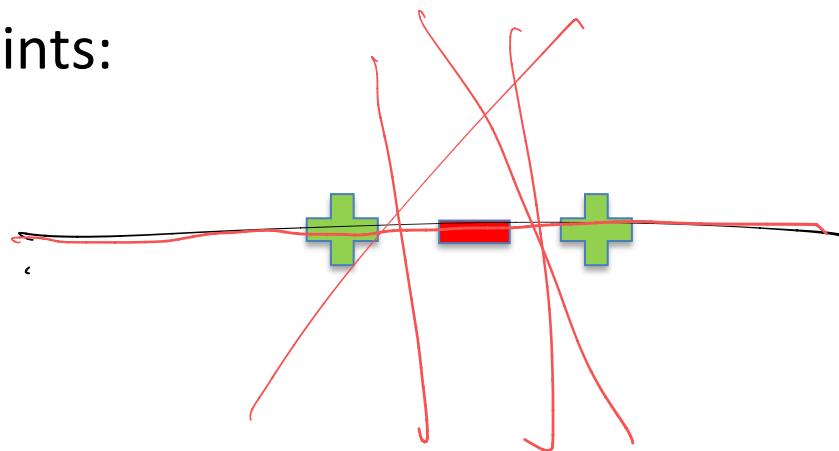
- How many points in 1-D can be correctly classified by a linear separator?

- 3 points:



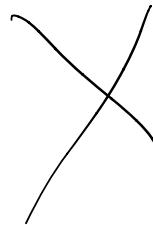
VC Dimension

- How many points in 1-D can be correctly classified by a linear separator?
 - 3 points:



$$VC(H) \leq 3$$

NO!

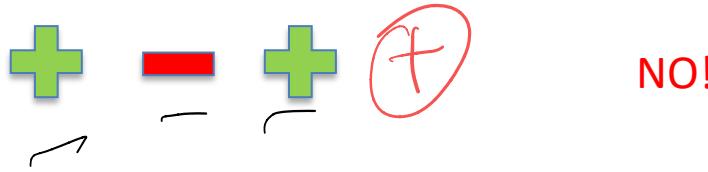


"Largest # points in d-dim space such that Hypothesis h can Shatter / Separate points is VC Dm."

VC Dimension

- How many points in 1-D can be correctly classified by a linear separator?
 - 3 points:

$$1D, VC(h) = 2.$$



- 3 points and up: for any collection of three or more there is always some choice of pluses and minuses such that the points cannot be classified with a linear separator (in one dimension)

Given dimensionality " m "

$$VC_n(h) = d \quad \text{if}$$

① \exists set C of points $(C \subseteq \mathbb{R}^m)$
of size d such that
 C is shattered by h .

② \nexists set C' of size $d+1$,
 C' cannot be shattered by h

Resist 1D

- ① \exists a set of 2 points which can be shattered by h .

	○	○
+	+	
+	-	
-	+	
-	-	

- ② If set of 3 points, h cannot shatter these 3 points

	○	○	○	
+	+	-	✓	
-	+	-	✗ ↗	
+	-	+	✗ ↘	
-	-	+	✓	
+	+	-	✓ ✓	
-	-	-	✓ ✓ ✓	

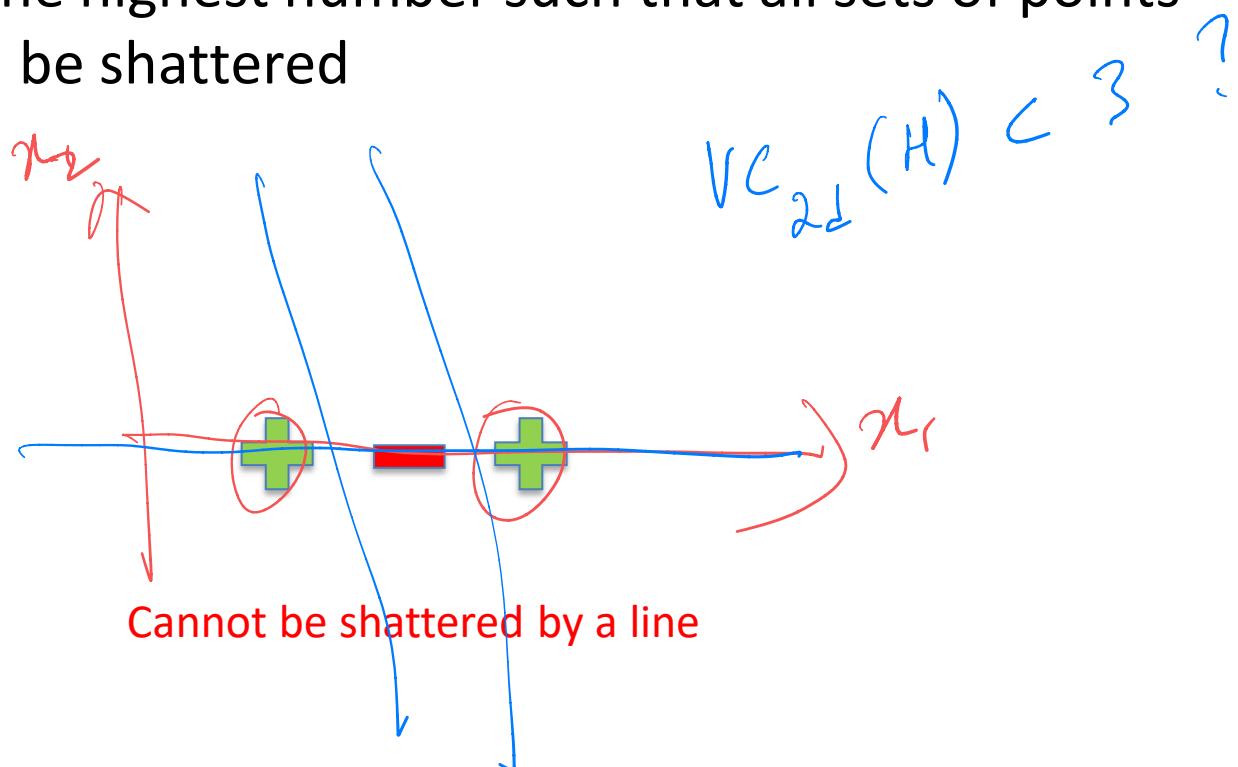
VC Dimension

Separation

- A set of points is shattered by a hypothesis space H if and only if for every partition of the set of points into positive and negative examples, there exists some consistent $h \in H$
- The Vapnik–Chervonenkis (VC) dimension of H over inputs from X is the size of the largest finite subset of X shattered by H

VC Dimension

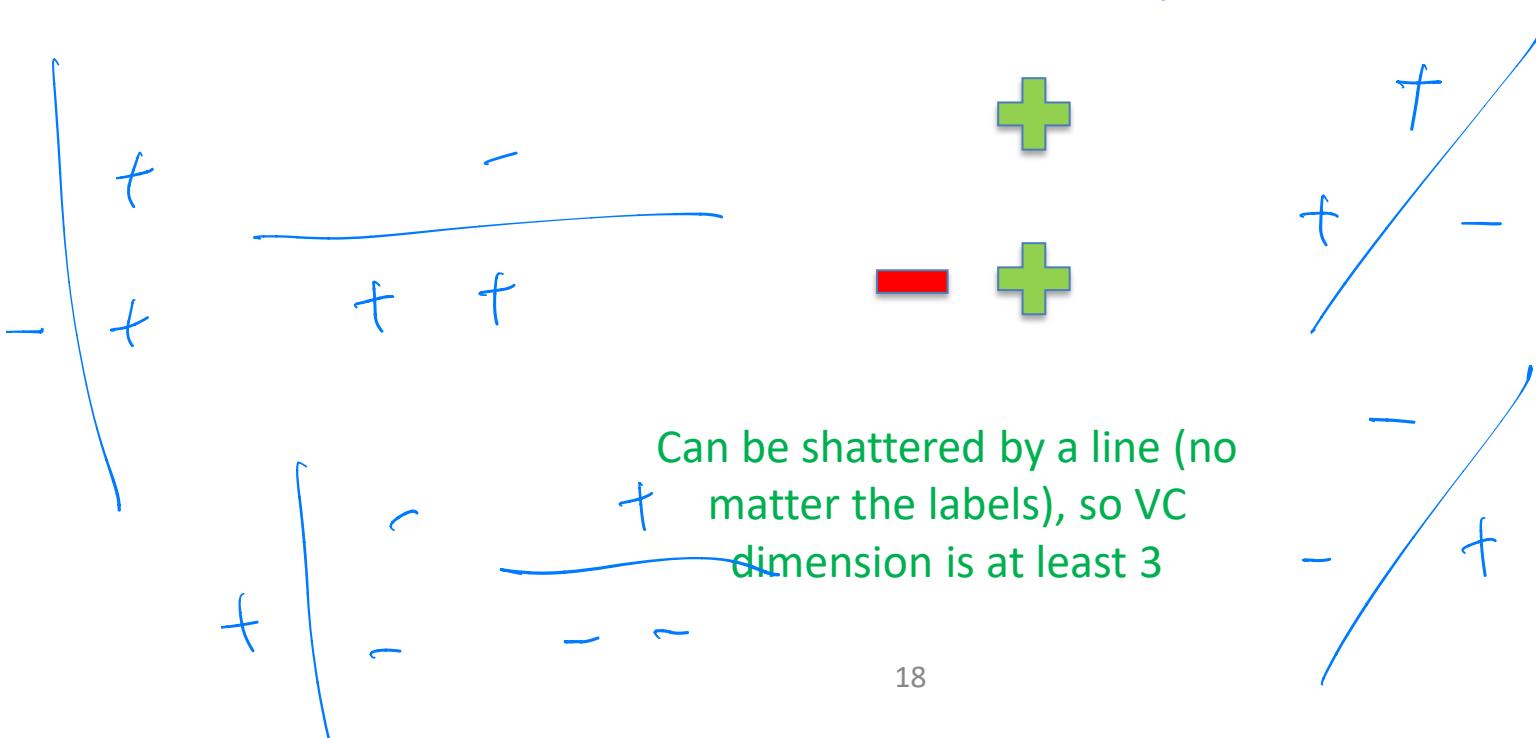
- Common misconception:
 - VC dimension is determined by the largest shattered set of points, not the highest number such that all sets of points that size can be shattered



VC Dimension

- Common misconception:
 - VC dimension is determined by the largest shattered set of points, not the highest number such that all sets of points that size can be shattered

$$VC_{2d}(\mathcal{H}) \geq 3$$

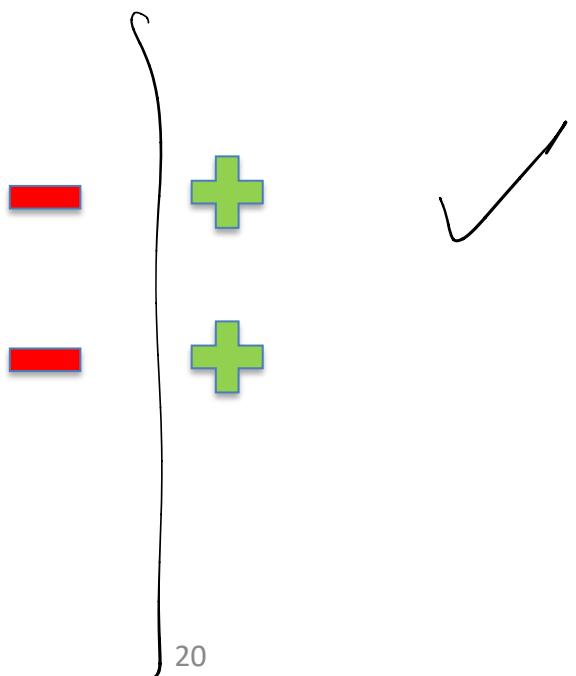


VC Dimension

- What is the VC dimension of 2-D space under linear separators?
 - It is at least three from the last slide
 - Can some set of four points be shattered?

VC Dimension

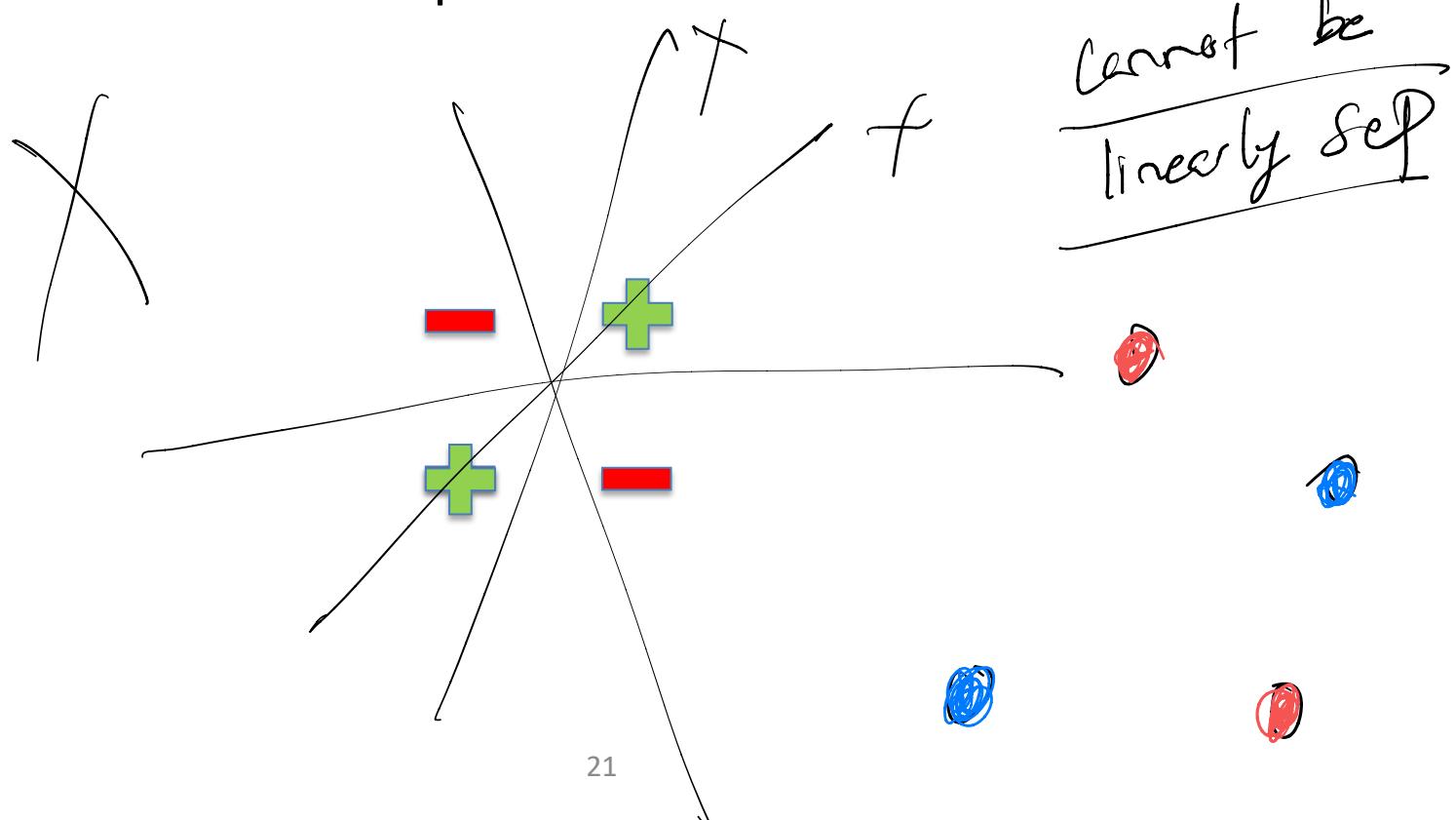
- What is the VC dimension of 2-D space under linear separators?
 - It is at least three from the last slide
 - Can some set of four points be shattered?



VC Dimension



- What is the VC dimension of 2-D space under linear separators?
 - It is at least three from the last slide
 - Can some set of four points be shattered?



VC Dimension

- What is the VC dimension of 2-D space under linear separators?
 - It is at least three from the last slide
 - Can some set of four points be shattered?



NO! This means that
the VC dimension is at
most 3

VC Dimension

- There exists a set of size $d + 1$ in a $d - \text{dimensional}$ space that can be shattered by a linear separator, but not a set of size $d + 2$
- The larger the subset of X that can be shattered, the more expressive the hypothesis space is
- If arbitrarily large finite subsets of X can be shattered, then $VC(H) = \infty$

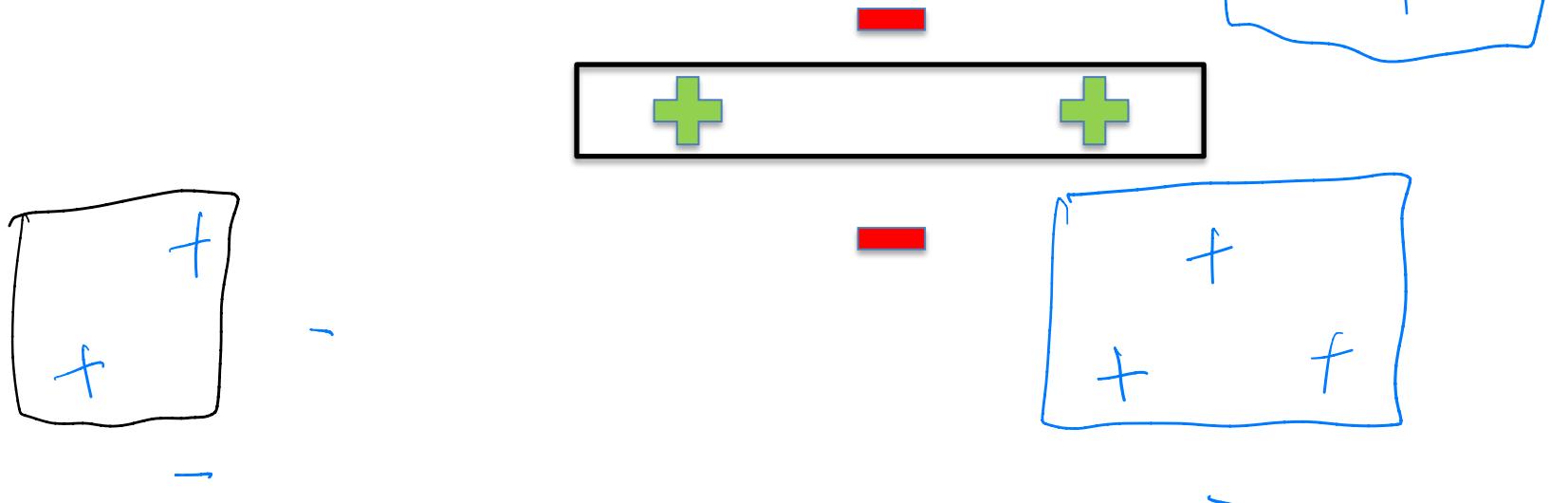
$$VC_d(h_{\text{Lin}}) = d+1$$

Axis Parallel Rectangles

- Let X be the set of all points in \mathbb{R}^2
- Let H be the set of all axis parallel rectangles in 2-D (inside + outside -)
 - What is $VC(H)$?

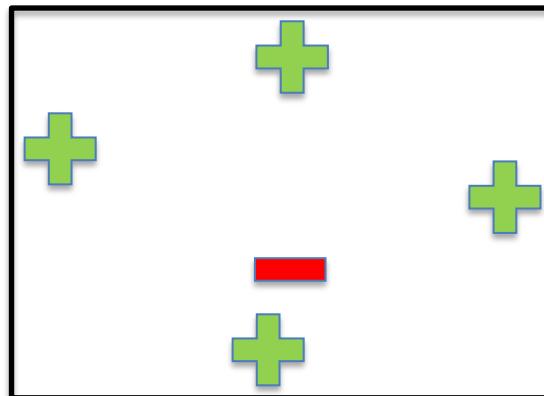
Axis Parallel Rectangles

- Let X be the set of all points in \mathbb{R}^2
- Let H be the set of all axis parallel rectangles in 2-D (inside + outside -)
 - $VC(H) \geq 4$



Axis Parallel Rectangles

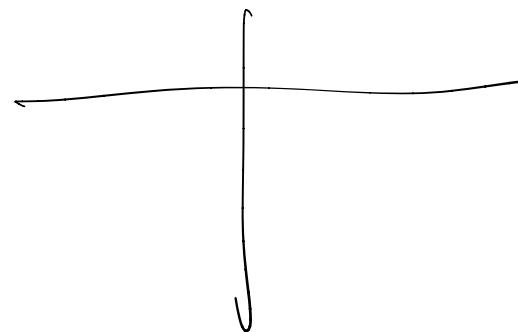
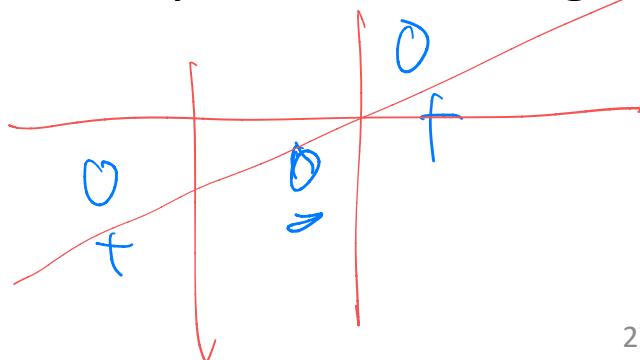
- Let X be the set of all points in \mathbb{R}^2
- Let H be the set of all axis parallel rectangles in 2-D
 - $VC(H) = 4$
 - A rectangle can contain at most 4 extreme points, the fifth point must be contained within the rectangle defined by these points



Examples

$d \geq 2$

- VC dimension of one-level decision trees over real vectors of length 2?
- VC dimension of linear separators through the origin?
- VC dimension of a hypothesis space with exactly one hypothesis in it for binary vectors of length $n \geq 1$?



Examples

$$n \in X \subseteq \mathbb{R}^d$$



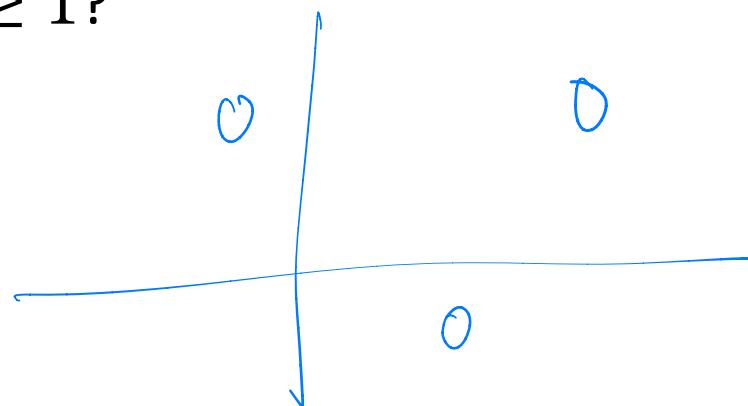
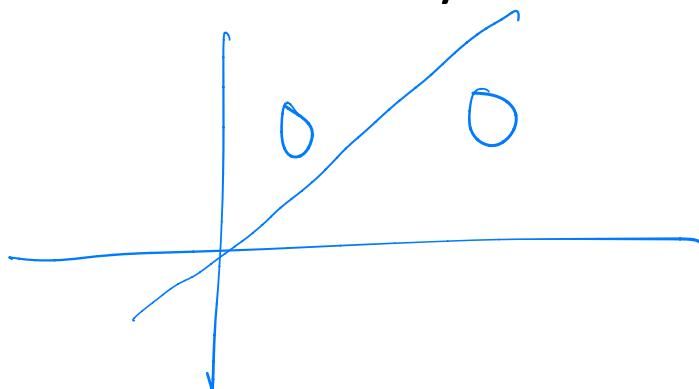
- VC dimension of one-level decision trees over real vectors of length 2?



- Three

- VC dimension of linear separators through the origin?

- VC dimension of a hypothesis space with exactly one hypothesis in it for binary vectors of length $n \geq 1$?



Examples

$\text{VC} \rightarrow "d+1"$, d -dim.

- VC dimension of one-level decision trees over real vectors of length 2?
 - Three
- VC dimension of linear separators through the origin?
 - Two
- VC dimension of a hypothesis space with exactly one hypothesis in it for binary vectors of length $n \geq 1$?

$\text{VC} \rightarrow "d"$ d -dim

" h " $h(n) = 1 \nexists^n$

⊕

Examples

- VC dimension of one-level decision trees over real vectors of length 2?
 - Three
- VC dimension of linear separators through the origin?
 - Two
- VC dimension of a hypothesis space with exactly one hypothesis in it for binary vectors of length $n \geq 1$?
 - Zero

$$\mathcal{E}^{\text{test}} \leq \mathcal{E}^{\text{true}} + \sqrt{\frac{1}{2m} \log_2 (|H| + 1)}$$

$$\text{VC}(H) \leq \log_2 |H|$$

PAC Bounds with VC Dimension



- VC dimension can be used to construct PAC bounds

$$M \geq \frac{1}{\epsilon} \left(4 \ln \frac{2}{\delta} + 8 \cdot VC(H) \ln \frac{13}{\epsilon} \right)$$

Sample Complexity
Bound

- Then, with probability at least $(1 - \delta)$ every $h \in H$ satisfies

$$\epsilon_h \leq \underbrace{\epsilon_h^{train}}_{T\epsilon} + \sqrt{\frac{1}{M} \left(VC(H) \left(\ln \left(\frac{2M}{VC(H)} \right) + 1 \right) + \ln \frac{4}{\delta} \right)}$$

test bound

- These bounds (and the preceding discussion) only work for binary classification, but there are generalizations

$\log_2 |H|$

PAC Learning



- Given:
 - Set of data X possible input / fech
 - Hypothesis space H model
 - Set of target concepts C Label-/target function
 - Training instances from unknown probability distribution over X of the form $(x, c(x))$ labels
- Goal:
 - Learn the target concept $c \in C$

PAC Learning

- Given:
 - A concept class C over n instances from the set X
 - A learner L with hypothesis space H
 - Two constants, $\epsilon, \delta \in (0, \frac{1}{2})$
- C is said to be PAC learnable by L using H iff for all distributions over X , learner L by sampling n instances will with probability at least $1 - \delta$ output a hypothesis $\underline{h \in H}$ such that
 - $\epsilon_h \leq \epsilon$
 - Running time is polynomial in $\frac{1}{\epsilon}, \frac{1}{\delta}, n, \underline{\text{size}(c)}$

PAC Learning

- Given:
 - A concept class C over n instances from the set X
 - A learner L with hypothesis space H
 - Two constants, $\epsilon, \delta \in (0, \frac{1}{2})$
- C is said to be PAC learnable by L using H iff for all distributions over X , learner L by sampling n instances will with probability at least $1 - \delta$ output a hypothesis $h \in H$ such that
 - $\epsilon_h \leq \epsilon$
 - Running time is polynomial in $\frac{1}{\epsilon}, \frac{1}{\delta}, n, \text{size}(c)$