

Machine Learning from Data

Lecture 5: Spring 2021

Today's Lecture

- Training Vs Testing
 - The Two Questions of Learning ✓
 - Theory of Generalization ($E_{in} \approx E_{out}$)
 - An Effective Number of Hypotheses

Feasibility

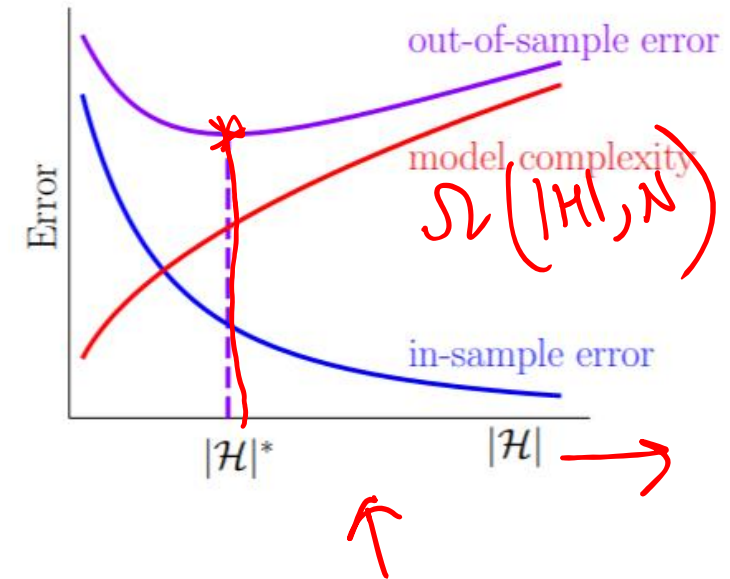
1. Can we make sure that $E_{\text{out}}(g)$ is close enough to $E_{\text{in}}(g)$? $E_{\text{in}}(g) \approx E_{\text{out}}(g)$
2. Can we make $E_{\text{in}}(g)$ small enough?

The Hoeffding generalization bound:

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \underbrace{\sqrt{\frac{1}{2N} \ln \frac{2|\mathcal{H}|}{\delta}}}_{\text{generalization error bar}}$$

E_{in} : training (eg. the practice exam)

E_{out} : testing (eg. the real exam)



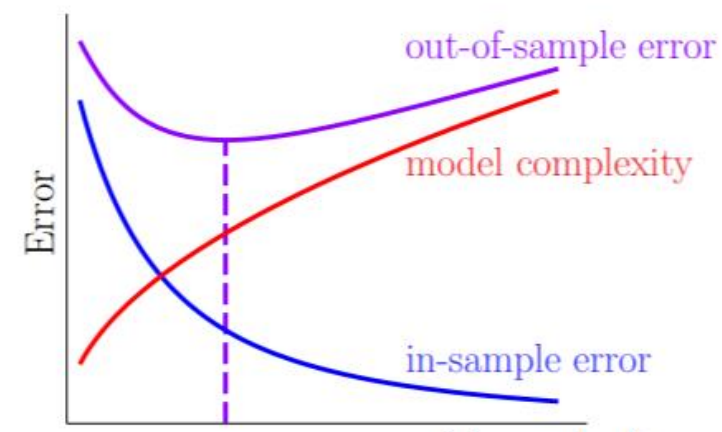
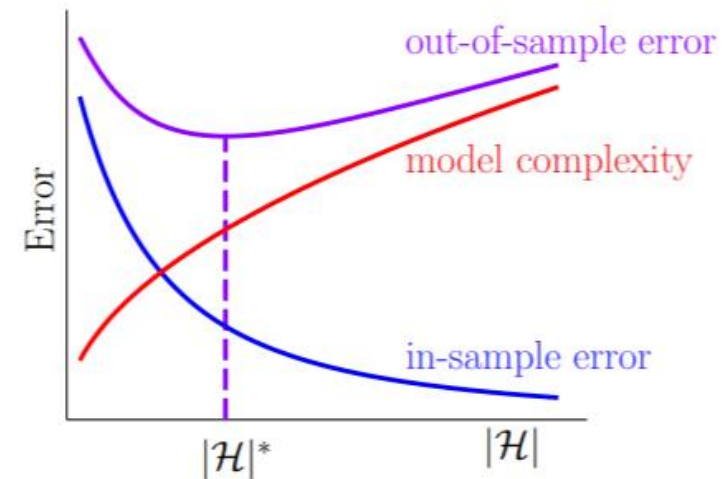
There is a tradeoff when picking $|\mathcal{H}|$.

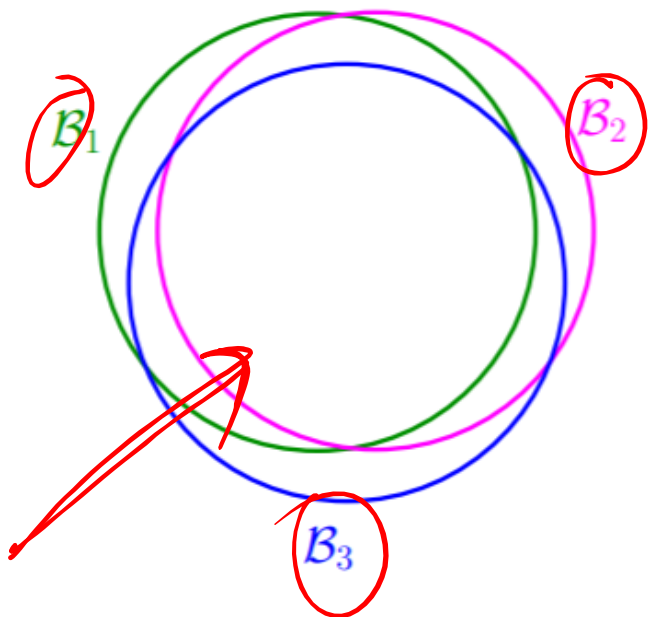
$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \underbrace{\left(\sqrt{\frac{1}{2N} \ln \frac{2|\mathcal{H}|}{\delta}} \right)}_{\text{error bar}} \quad \begin{matrix} |\mathcal{H}| \\ H \rightarrow \infty \end{matrix}$$

Most models



$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \frac{4m_{\mathcal{H}}}{\delta}}$$





- \mathcal{B}_m are events (sets of outcomes); they can overlap.
- If the \mathcal{B}_m overlap, the union bound is loose.
- If many h_m are similar, the \mathcal{B}_m overlap.
- There are “effectively” fewer than $|\mathcal{H}|$ hypotheses.
- We can replace $|\mathcal{H}|$ by something smaller.

Meaning of cardinality of \mathcal{H} ?

$|\mathcal{H}|$

How did \mathcal{H} come in?

$$|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon, g \in \mathcal{H}$$

\mathcal{B}_m is a bad event

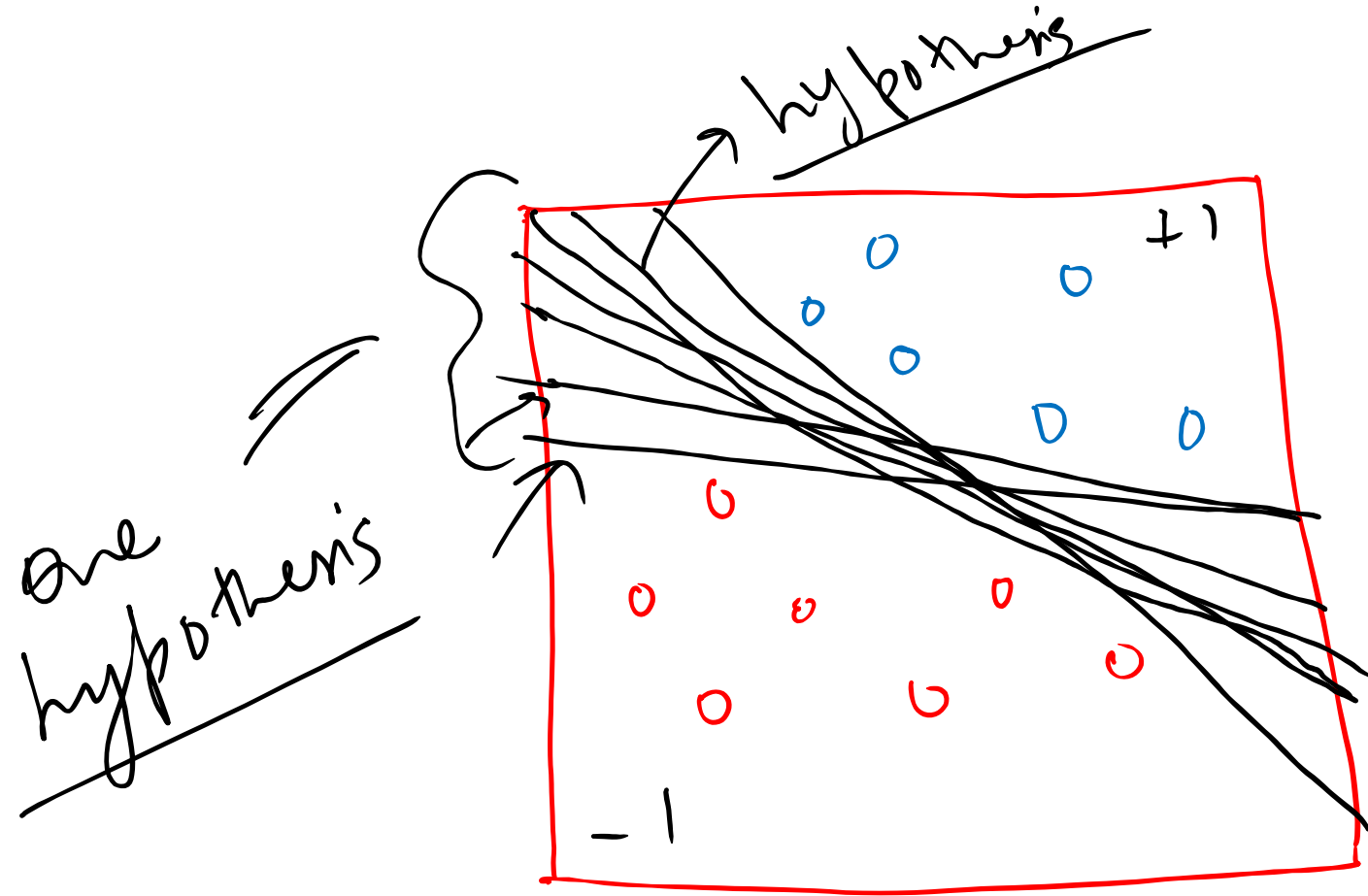
$$|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon$$

$$P[\mathcal{B}_g] \leq P[\text{any } \mathcal{B}_m] \leq \sum_{m=1}^{|\mathcal{H}|} P[\mathcal{B}_m]$$

Effective Number of Hypothesis

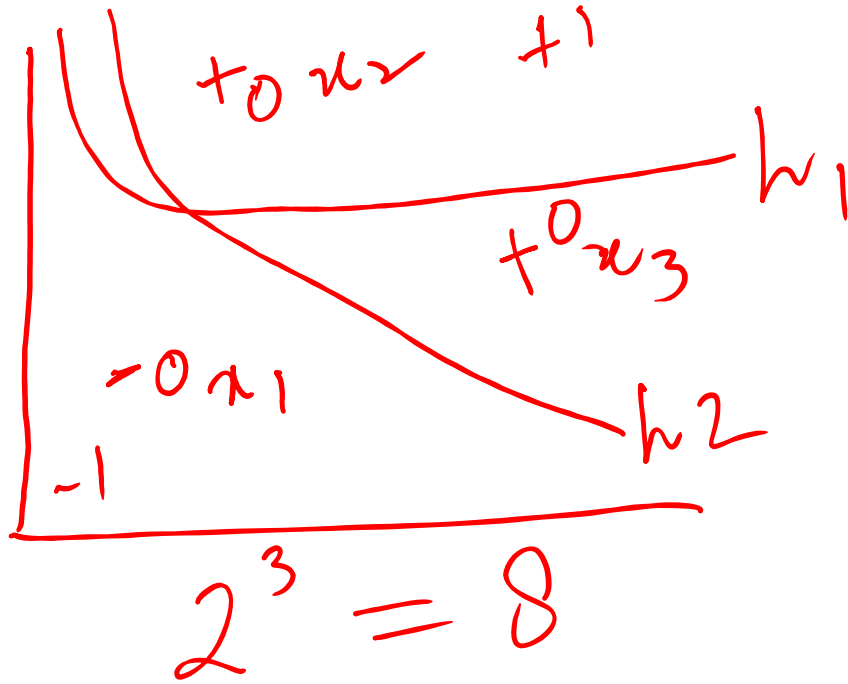
- We need a way to measure the diversity of H .
- A simple idea:
 - Fix any set of N data points.
 - If H is diverse it should be able to implement all functions . . . on these N points.

Similarity
 $|H|$
infinite



The Growth Function

•



dictating →

	x_1	x_2	x_3
h_1	$\bar{+}$	$\bar{+}$	$\bar{+}$
h_2	$-$	$+$	$+$
\vdots			
h_q	$-$	$+$	$+$

Another Example

(N -points)

$$2^N$$

$\delta \rightarrow$ dichotomy

	x_1	x_2	...	x_N
h_1	\ominus	$-$		$+$
h_2	$+$	$-$		$+$
\vdots				

\mathcal{H}

Restriction of the hypothesis set \mathcal{H} to the points x_1, x_2, \dots, x_N

$$\mathcal{H}(x_1, x_2, \dots, x_N) = \{ \delta \mid \delta \text{ is implemented by } h \in \mathcal{H} \}$$

Quantify the complexity of K on x_1, x_2, \dots, x_N
using $\left| K(x_1, \dots, x_N) \right| \leq \underline{\underline{2^N}}$?

Q1: What x_1, x_2, \dots, x_N should we use to
compute the complexity of a hypothesis set?

Q2: Is this bound for complexity useful?
f.

Ans] Yes,
Worst case Analysis

$$m_H(N) = \max_{x_1, x_2, \dots, x_N} |H(x_1, x_2, \dots, x_N)|$$

growth function!

$$m_H(N) \leq 2^N$$

Ans 2 Error bar:

$$\sqrt{\frac{1}{2N} \frac{\ln 2(H)}{\delta}} \rightarrow \underline{\underline{m_H(N)}}$$

$$\sqrt{\frac{1}{2N} \frac{\ln 2^{2^N}}{\delta}}$$

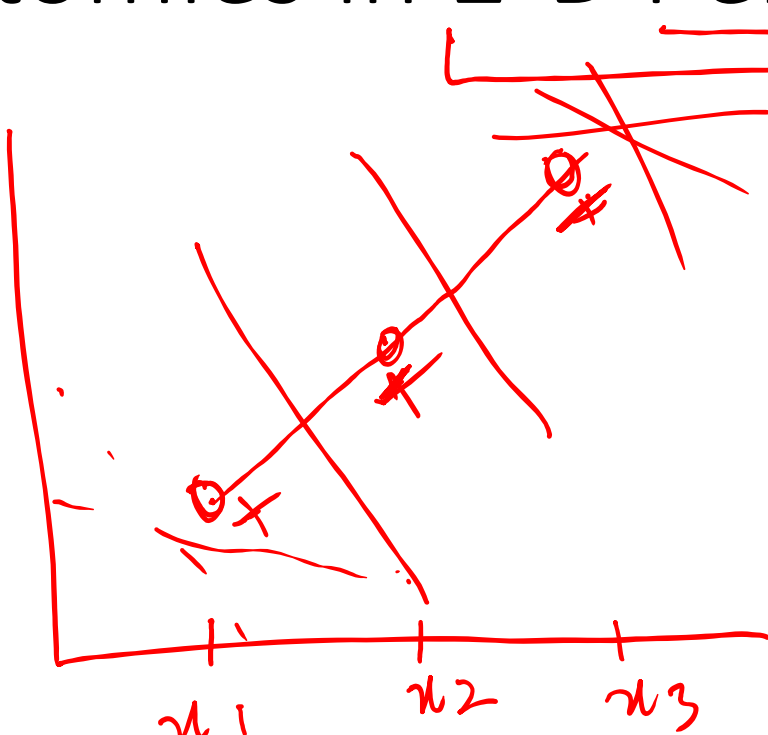
$$\sqrt{\frac{1}{2N} \frac{d \ln 2^{N^2}}{\delta}}$$

$$\frac{\ln N}{N}$$

$$N \rightarrow \infty \left(\text{limit } \frac{\ln N}{N} = 0 \text{ as } N \rightarrow \infty \right)$$

Dichotomies in 2-D Perceptron

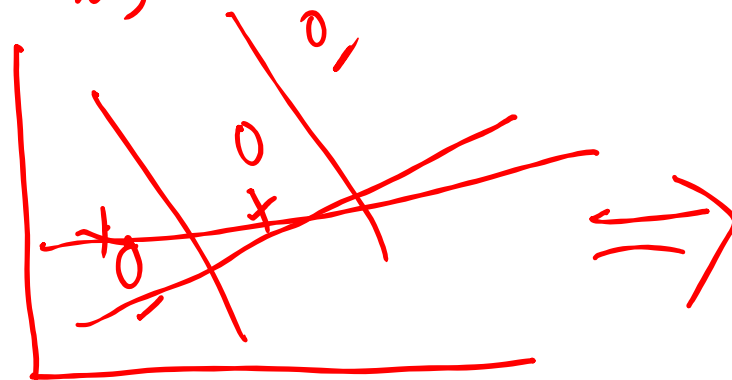
•



x_1	x_2	x_3
+	+	+
-	+	+
-	-	+
-	-	-
+	-	-
+	+	-

3 points = 2^N

$$|H(x_1, x_2, x_3)| = 6$$

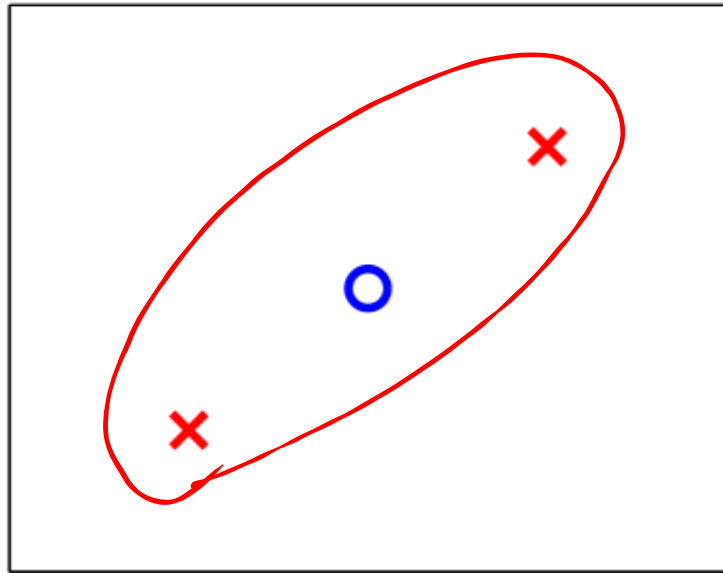


-	+	-
+	-	+

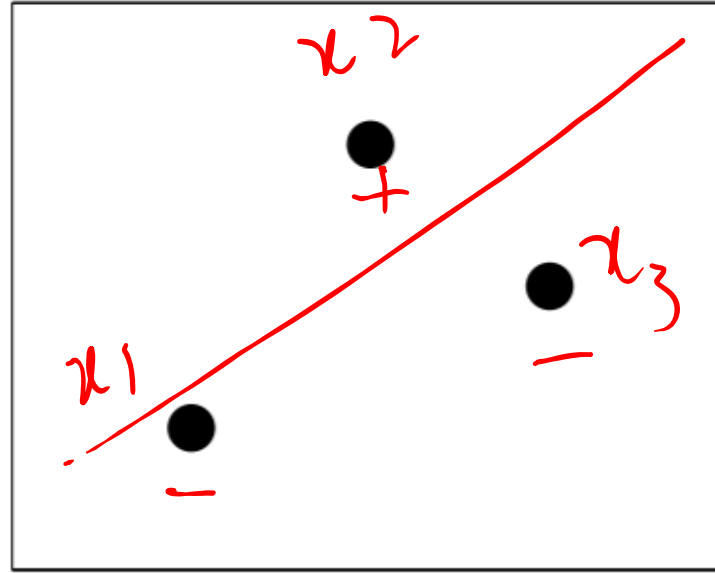
Example - Growth Functions

2-D Perceptron Model

For 3 points $= 2^3 = 8$

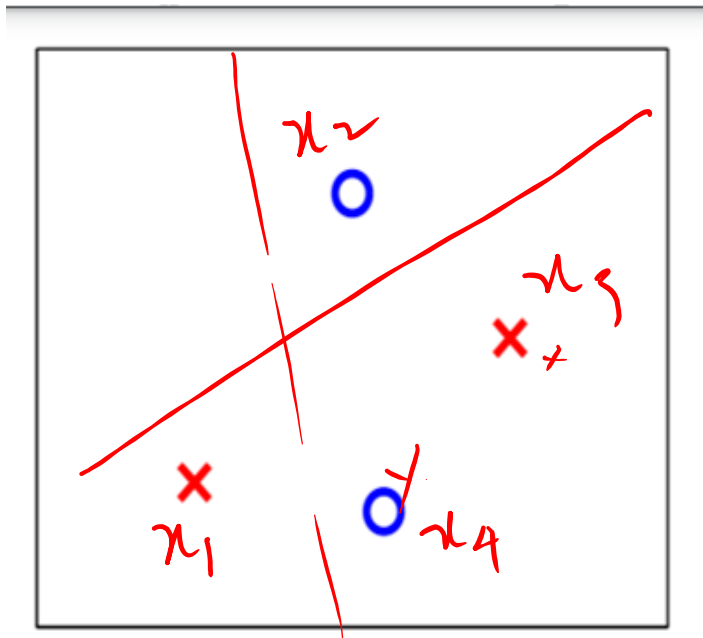


x



←
✓

$$m_H(N) = 8 \\ = 2^3$$



HL \rightarrow Perception

$$\underline{\underline{2^4 = 16}}$$

$$\begin{matrix} - & + & - & + \\ + & - & + & - \end{matrix} \} m_H(4) = 14 < 2^N$$

\uparrow \nearrow
 $14 < 16$

3 \rightarrow 4 points

\rightarrow deficiency

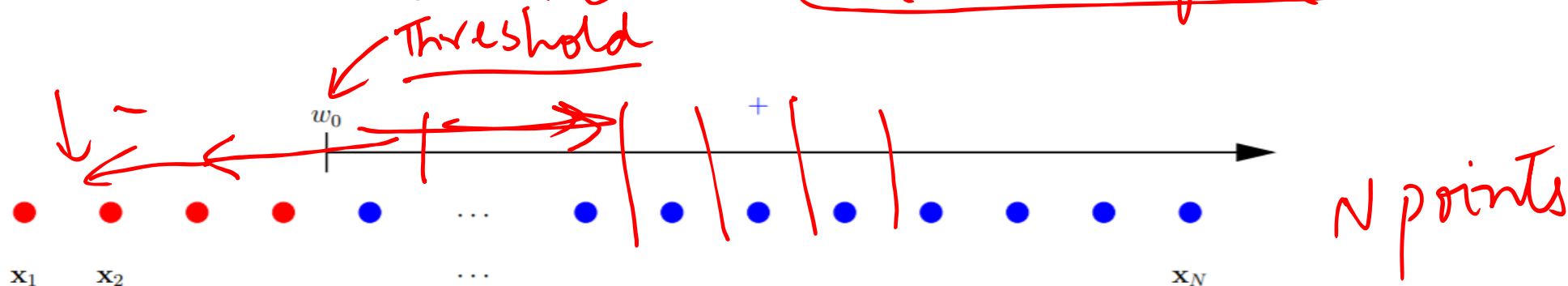
Value = 4

is important

$$\underline{\underline{2^4}}$$

1-D Positive Ray \mathcal{H}

$$h(x) = \text{sign}(x - w_0)$$



N points

$N+1$ regions

$$m_{\mathcal{H}}(N) = N+1$$

N	1	2	3	4
$m_H(N)$	2	3	4	5
2^N	2	4	8	16

Broken at $N=2$

Worst Case 2^N

$$3 < 4$$

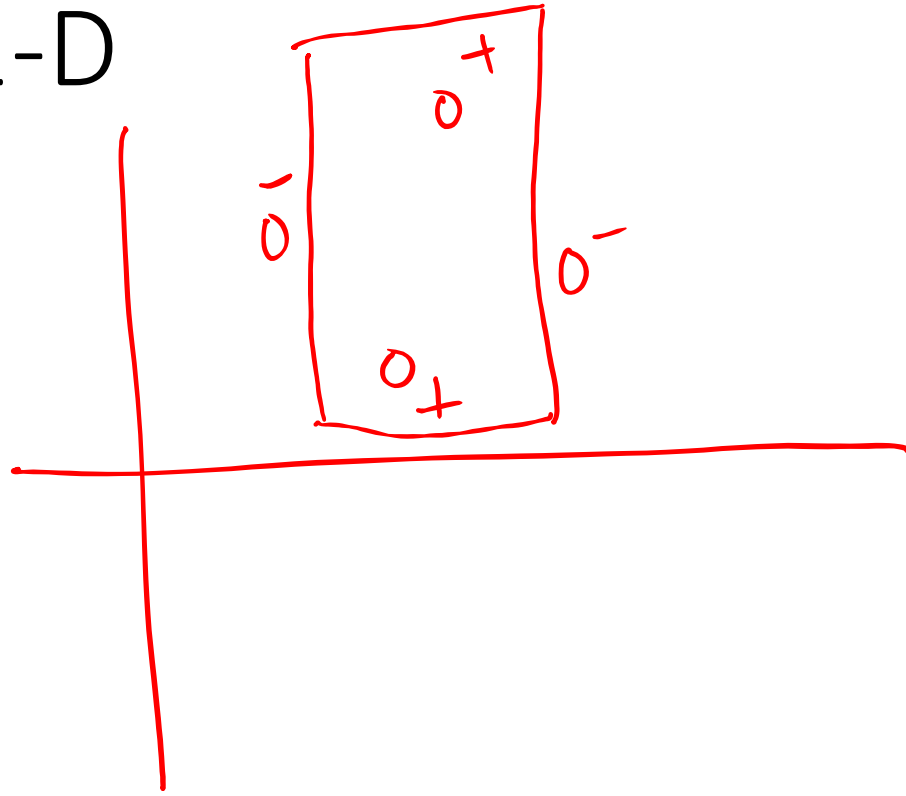
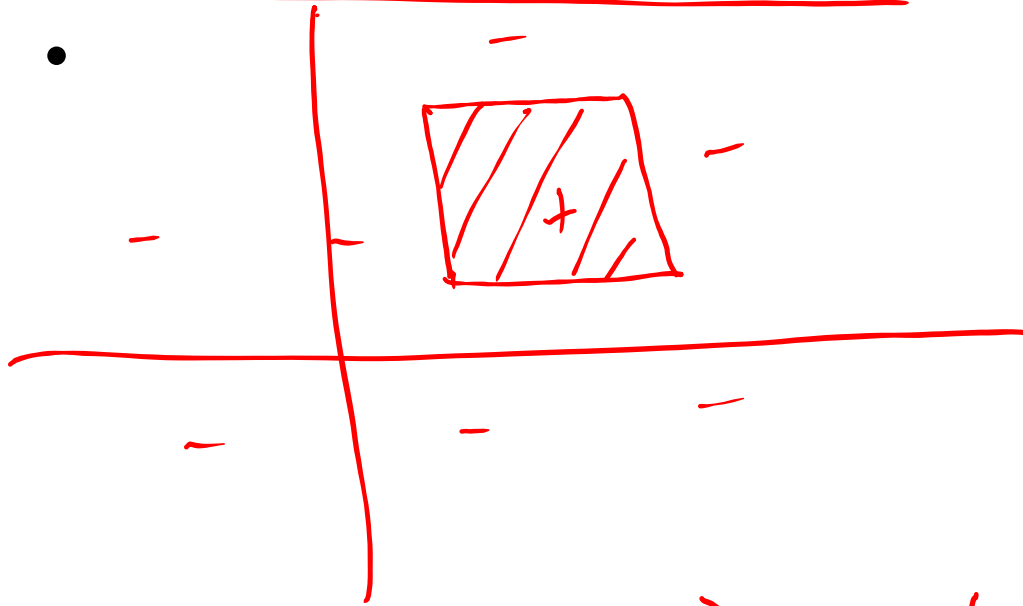
Perceptions

N	1	2	3	4	5
$m_H(N)$	2	4	8	14	
2^N	2	4	8	16	

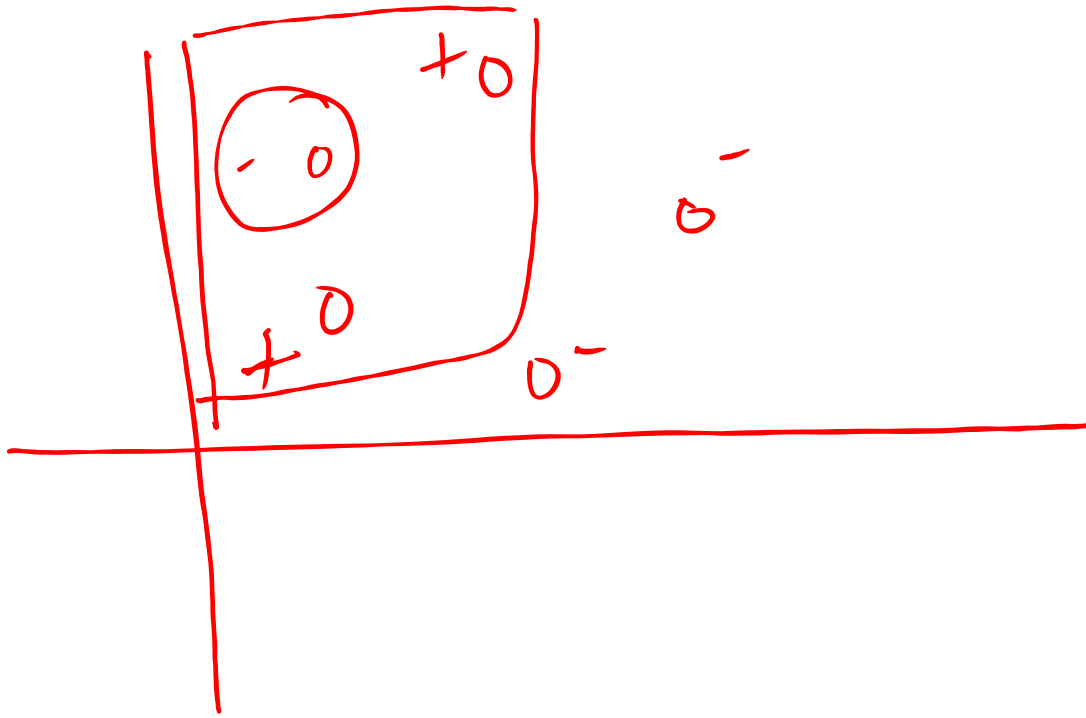
*

Positive Rectangles in 2-D

-



$$m_H(N) = m_n(4) = 16(2^N)$$



$$\underbrace{m_h(5)} < 2^5$$

Summarize

Effective no. $\rightarrow K, N$
 2^N

K

	1	2	3	$N=4$ 4	5	...
<u>2-D perceptron</u>	2	4	8	14	...	
<u>1-D pos. ray</u>	2	3	4	5	...	
2-D pos. rectangles	2	4	8	16	$< 2^5$...

- $m_{\mathcal{H}}(N)$ drops below 2^N — there is hope for the generalization bound.
- A break point is any n for which $m_{\mathcal{H}}(n) < 2^n$.

Definition: Shatter a Data Set

\mathcal{H}

- x_1, x_2, \dots, x_N

2^N

$\mathcal{H} \xrightarrow[\text{implement}]{\text{can}}$ All possible dichotomies (2^N)

Finding
4 dichotomies
that cannot
shatter

$$2^N = 8$$

Combinatorial Puzzle

any pair of
2 points

0 -
● +

\checkmark \mathbf{X}_1	\checkmark \mathbf{X}_2	\checkmark \mathbf{X}_3
○	○	○
○	○	●
○	●	○
○	●	●

2 points are
shattered

	X_1	X_2	X_3
→	○	○	○
→	○	○	● +
	○	● +	○ -
	○	● +	● +

No Pair of Points
is Shattered

X_1	X_2	X_3
$\overset{-}{\circ}$	$\overset{-}{\circ}$	$\overset{-}{\circ}$
\circ	\circ	\bullet
$\overset{\checkmark}{\circ}$	\bullet	$\overset{\checkmark}{\circ}$
\bullet	$\overset{-}{\circ}$	$\overset{-}{\circ}$

If $N = 4$ how many possible dichotomys with no 2 points shattered?

$$2^4 = 16$$

\mathbf{X}_1	\mathbf{X}_2	\mathbf{X}_3
○	○	○
○	○	●
○	●	○
●	○	○

\mathbf{X}_1	\mathbf{X}_2	\mathbf{X}_3	\mathbf{X}_4
○	○	○	○
○	○	○	●
⋮			

Thanks!