

# Learning Theory

Nishant Mehta

Lecture 12 - Part II

# Generalization

Suppose a learning algorithm returns a hypothesis with low training error.

When can we guarantee that the hypothesis's true error is also low?

**Main question:** How can we use the training error of a learning algorithm to estimate the algorithm's true error?

# Generalization

Main question: How can we use the training error of a learning algorithm to estimate the algorithm's true error?

We will focus on *binary classification*

Recall:

Training error:  $\hat{R}(\hat{h}, D) = \frac{1}{n} \sum_{i=1}^n \mathbb{1} [\hat{h}(X_i) \neq Y_i]$

True error:  
(Risk)  $R(\hat{h}) = \mathbb{E}_{(X,Y) \sim P} [\mathbb{1} [\hat{h}(X) \neq Y]]$   
 $= \Pr_{(X,Y) \sim P} (\hat{h}(X) \neq Y)$

# Realizable setting

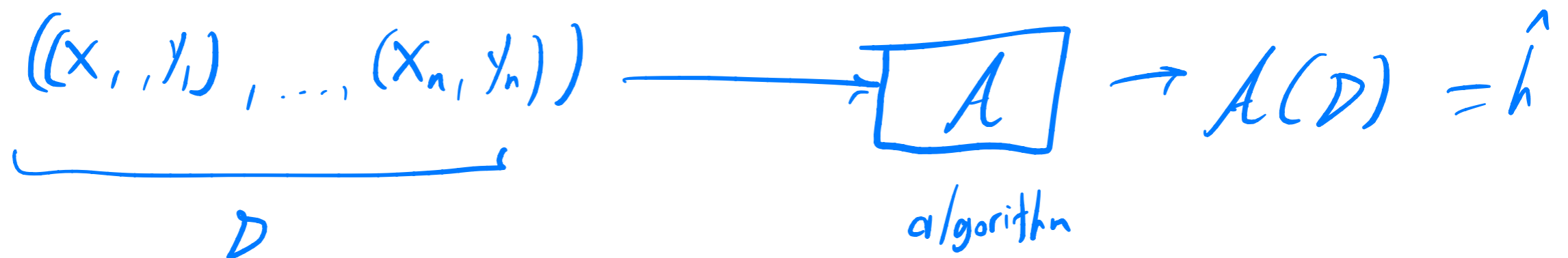
Suppose that each example is in input space  $\mathcal{X}$

In the *realizable setting*:

There is a known *concept class*  $\mathcal{C}$ , a set of concepts, where each concept  $c$  is a rule mapping from  $\mathcal{X}$  to  $\{0, 1\}$

There is a concept  $c \in \mathcal{C}$  such that, for any input  $X$ , the label is  $Y = c(X)$

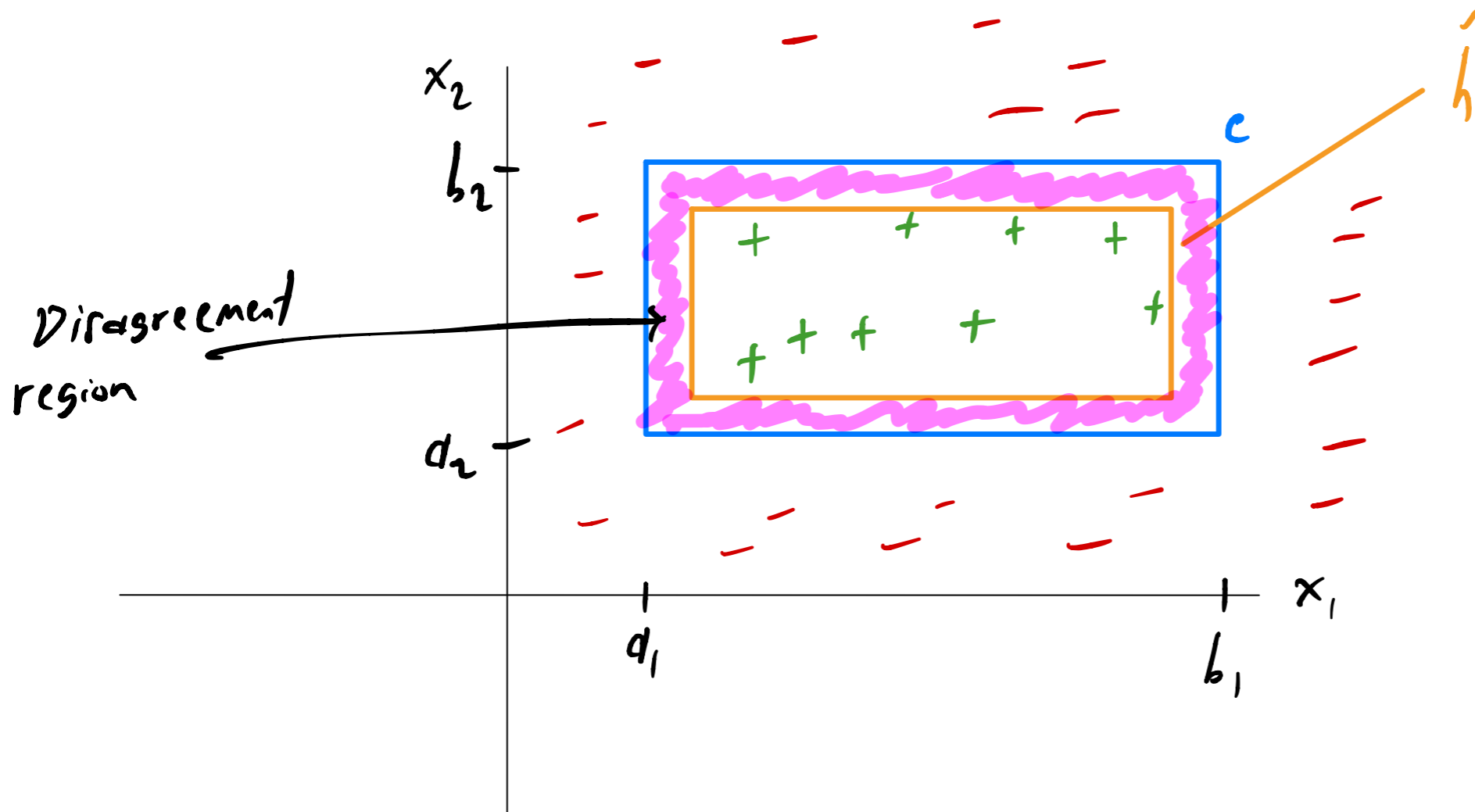
Given training data, learning algorithm selects hypothesis  $\hat{h} \in \mathcal{H}$



# Realizable setting: Example 1

Learning rectangles

$$\mathcal{X} = \mathbb{R}^2$$

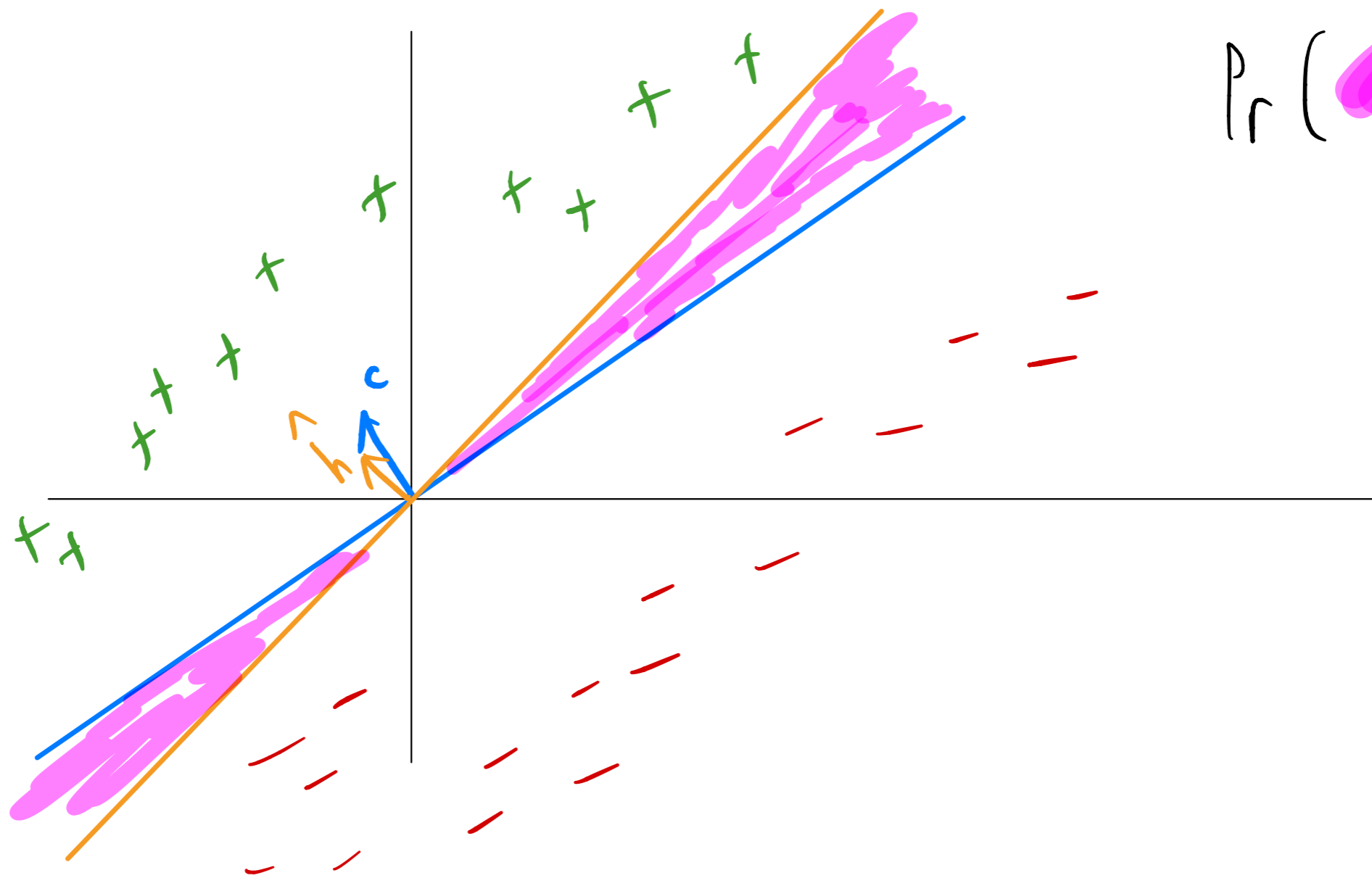


$$R(\hat{h}) = \Pr(x \in \square)$$

$$c(x) = \mathbb{1}[a_1 \leq x_1 \leq b_1] \cdot \mathbb{1}[a_2 \leq x_2 \leq b_2]$$

# Realizable setting: Example 2

Homogeneous linear separators in  $\mathbb{R}^2$



$$R(\hat{h}) = P_r(\text{shaded region})$$

# Back to generalization

In the realizable setting, we have:

$$\text{Training error: } \hat{R}(\hat{h}, D) = \frac{1}{n} \sum_{i=1}^n \mathbb{1} [\hat{h}(X_i) \neq c(X_i)]$$

$$\text{Risk: } R(\hat{h}) = \Pr_{X \sim P} (\hat{h}(X) \neq c(X))$$

Main question:

How can we use training error  $\hat{R}(\hat{h}, D)$  to upper bound risk  $R(\hat{h})$ ,  
***no matter what distribution the data comes from?***

**A bad learning algorithm**

memorization algorithm

$$\hat{h}(x) = \begin{cases} y_i & \text{if } \exists i \in [n] \ x = X_i \\ 0 & \text{otherwise} \end{cases}$$

$$R(\hat{h}) = ?$$