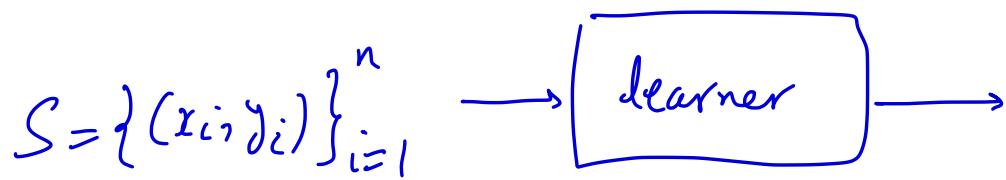


Lecture 23:



ultimate goal:

$$\min_h \Pr_{(x,y) \sim D} \{ h(x) \neq y \}$$



the first relaxation:

$$\min_h \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\} \rightarrow \text{bad} = 0$$



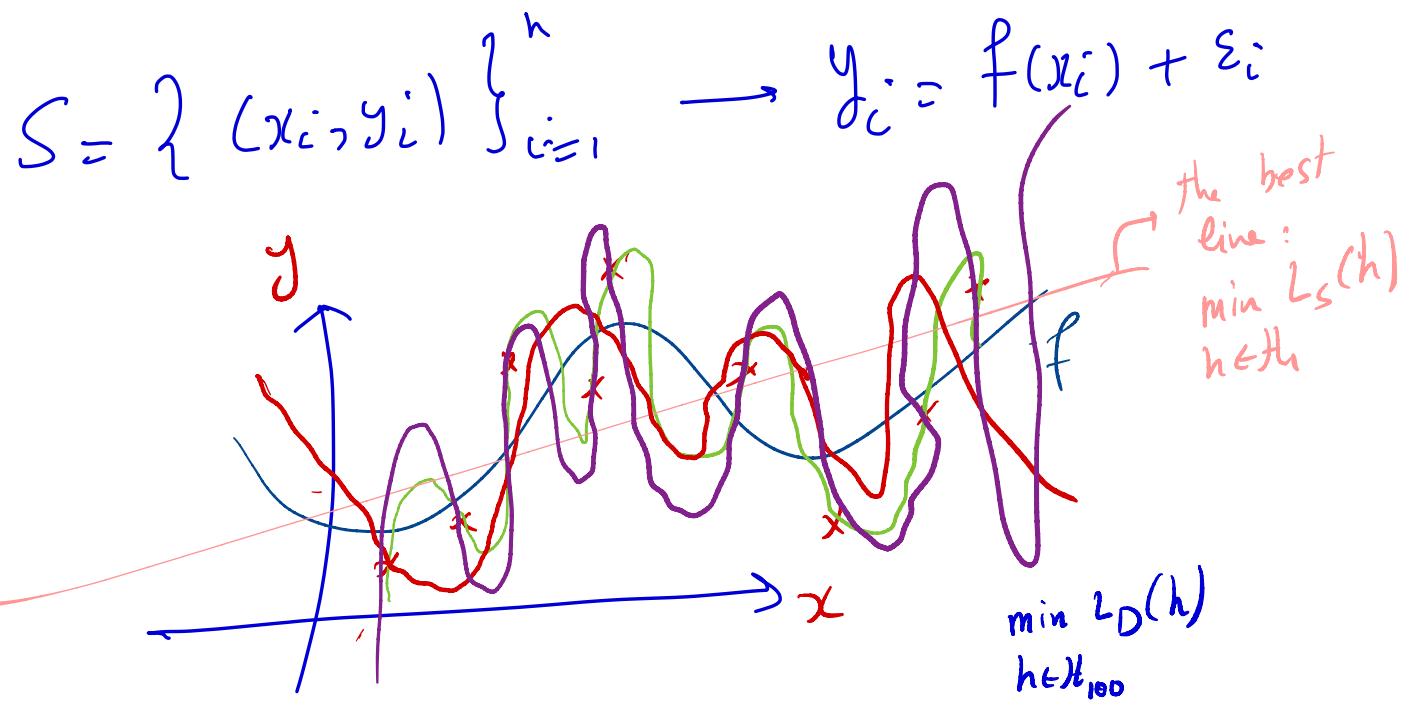
restricted ERM_H

$$\min_{h \in H} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\}$$

the predictor is inside the class H.

Example:

$$y = f(x) + \varepsilon \quad E[\varepsilon] = 0, \quad \text{var}[\varepsilon] = 6^2$$



$$L_D(h) = E_{(x,y) \sim D} \left[(h(x) - y)^2 \right]$$

$\min_h L_D(h) \Rightarrow$ the minimizer is f

$$\begin{aligned} \min_{h \in H} L_D(h) &= E[(y - f(x))^2] \\ &= E[\varepsilon^2] = 6^2 \end{aligned}$$

training error

$$\min_h L_S(h) = 0$$

$$\left\{ \begin{array}{l} \min_h L_D(h) = \underline{\epsilon}^2 \\ \min_h L_S(h) = 0 \end{array} \right.$$

$$\min_{\substack{h \in \mathcal{H}}} L_S(h)$$

examples for \mathcal{H} :

$\mathcal{H}_1 = \left\{ \begin{array}{l} \text{all the linear functions} \\ \text{i.e. polynomials of degree 1} \end{array} \right\}$

↳ since \mathcal{H}_1 is not rich enough, we will have underfitting

\Rightarrow Hence the function class \mathcal{H} that we're going to use, should be sufficiently rich/complex.

$$H_6 = \{ \text{all the polynomials of degree 6} \}$$

⇒ the problem with using rich/complex function classes is: overfitting

In summary: we have 2 problems:

Problem 1:

if H is not rich enough, it'll underfit:

$$\min_{h \in H} L_D(h) \gg \min_h L_D(h)$$

Solution: use a class H that's sufficiently rich.

Problem 2:

if H is complex, then it may suffer from overfitting:

$$\min_{h \in \mathcal{H}} L_D(h) \approx \min_{h} L_D(h)$$

training error

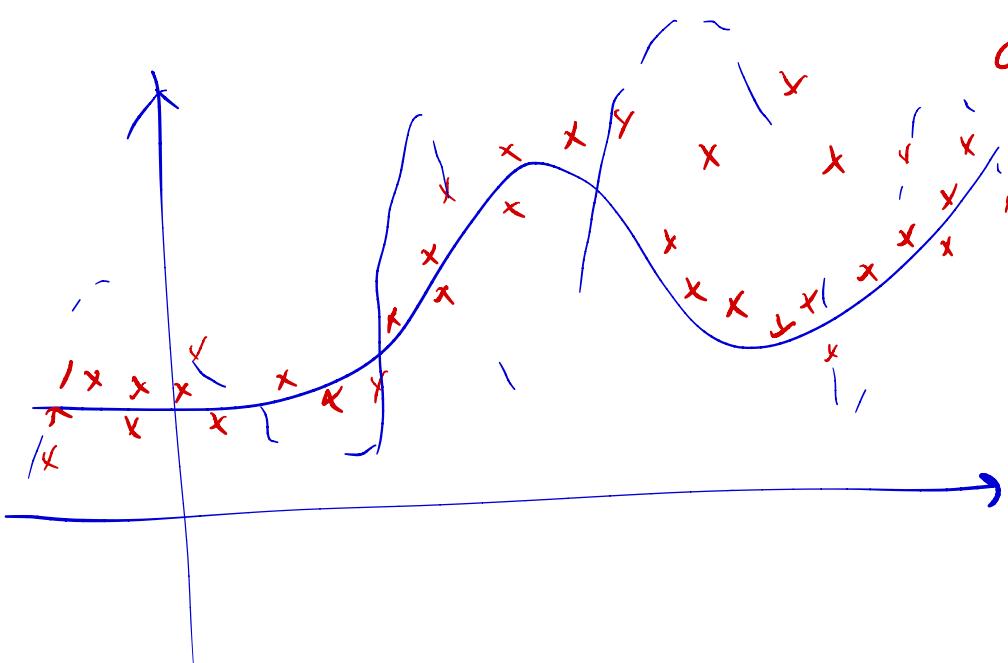
$$\min_{h \in \mathcal{H}} L_S(h) \rightarrow h_S^* \rightarrow L_S(h_S^*) \approx 0$$

very rich/complex

$$\min_h L_D(h) \ll L_D(h_S^*)$$

*true error of
what we learned
from fitting the
function class \mathcal{H}*

*to the training
data S.*



$$L_S(h) \rightarrow L_D(h)$$

$$\min_{h \in \mathcal{H}} L_S(h) \longrightarrow \min_{h \in \mathcal{H}} L_D(h)$$

\Rightarrow to avoid overfitting we should have a sufficient number of training data points.

$$\text{for every } h \in \mathcal{H} \rightarrow L_S(h) \approx L_D(h) \pm \epsilon$$

$$\min_{h \in \mathcal{H}} L_S(h) \approx \min_{h \in \mathcal{H}} L_D(h) \pm \epsilon$$

It seems that to avoid overfitting the number of training data points should be sufficiently large

\hookrightarrow - depends on $\cdot \mathcal{H}$

- depends on the accuracy ϵ

PAC learning (PAC = Probably Approximately Correct)

We say that a function class \mathcal{H} is PAC-learnable if for any $\epsilon, \delta > 0$, there exists a number $n_0(\epsilon, \delta)$, such that for any distribution of data, D , the following holds: If $n \geq n_0(\epsilon, \delta)$, then with probability $1 - \delta$ we have that:

$$\forall h \in \mathcal{H} : |L_D(h) - L_S(h)| \leq \epsilon.$$

$$\underbrace{\min_{h \in \mathcal{H}} L_S(h)}_{\approx} \approx \underbrace{\min_{h \in \mathcal{H}} L_D(h)}_{\approx}$$