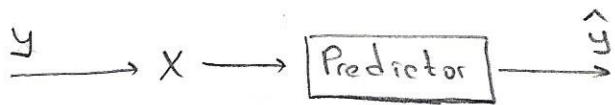


# Lecture 6:



Data:

$\{(x_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d}}{\sim} P$  ( $P$  is the true unknown distribution from which nature generates these samples)

Hypothesis class:  $\mathcal{H} = \{h : X \rightarrow Y\}$   
the class of functions we can use as predictors

Loss function:  $l(y, \hat{y}) : Y \times Y \rightarrow \mathbb{R}^+$   
quantifies the error associated with predicting  $y$  as  $\hat{y}$ .

Empirical loss of  $h \in \mathcal{H}$ :

$$L_n(h) = \frac{1}{n} \sum_{i=1}^n l(h(x_i), y_i)$$

Population loss of  $h \in \mathcal{H}$ :

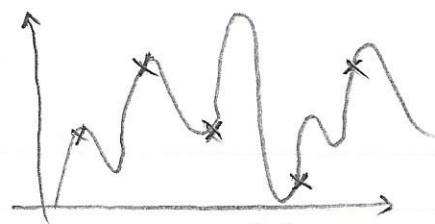
$$L(h) = \mathbb{E}_P[l(h(x), y)]$$

Empirical Risk Minimization (ERM):

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} L_n(h)$$

The ERM solution or the empirical solution.

Key question: How can we ensure that  $L(\hat{h})$  is close to  $L(h)$ ?



$$L_n(h) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i) \xrightarrow{n \rightarrow \infty} \mathbb{E}_P[\ell(h(x), y)] = L(h)$$

$(x_i, y_i) \stackrel{i.i.d.}{\sim} P$       because of LLN.

$L_n(\hat{h}) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{h}(x_i), y_i)$

ERM: We approximate the true expectation of  $z_i = \ell(h(x), y)$  with the empirical or sample mean.

Uniform Convergence: Show that if take  $n$  iid samples from  $P$  and compute  $L_n(h)$  for all  $h \in \mathcal{H}$ , then  $L_n(h) \approx L(h)$  for all  $h \in \mathcal{H}$ .

Hoeffding's Inequality: Let  $X_1, X_2, \dots, X_n$  i.i.d.  $\mathbb{E}[X_i] = \mu$  &  $a \leq X_i \leq b$ , for any  $t \geq 0$

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| \geq t\right) \leq 2e^{-2nt^2/(b-a)^2}$$

Assume  $0 \leq \ell(\hat{y}, \tilde{y}) \leq 1$  (e.g. 0-1 loss) ✓  
 or  $\ell(y, \hat{y}) = (y - \hat{y})^2$   $0 \leq y, \hat{y} \leq 1$ .

Then for any  $h \in \mathcal{H}$

$$P(|L_n(h) - L(h)| \geq \epsilon) \leq 2e^{-2n\epsilon^2}$$

We want  $|L_n(h) - L(h)| \leq \epsilon \quad \forall h \in \mathcal{H}$

$$P(\text{There exists } h \in \mathcal{H} \text{ s.t. } |L_n(h) - L(h)| \geq \epsilon)$$

$$= P\left(\bigcup_{h \in \mathcal{H}} \{|L_n(\hat{h}) - L(h)| \geq \epsilon\}\right)$$

$$\leq \sum_{h=1}^{|\mathcal{H}|} P(|L_n(\hat{h}) - L(h)| \geq \epsilon)$$

$$\leq 2|\mathcal{H}| 2e^{-2n\epsilon^2} = 2e^{\log|\mathcal{H}| - 2n\epsilon^2} \quad (*)$$

# of functions  
in the class  
 $\mathcal{H}$

If we want  $(*)$  to be  $\leq \delta$ , then

$$e^{\log|\mathcal{H}| - 2n\epsilon^2} \leq e^{\delta/2} \Rightarrow n > \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2}$$

If  $\mathcal{H}$  is a function class with  $d$  parameters (e.g. a neural network with  $d$  coefficient's) and each parameter is represented by a fixed precision, say  $k$  bits

$$|\mathcal{H}| \approx (2^k)^d$$



Biggest open problem in ML: Deep network generalize well even when  $n \ll d$ .

Population

Solution

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} L(h) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathbb{E}_P[l(h(x), y)]$$

$$L(\hat{h}) = \underbrace{L(\hat{h}) - L(h^*)}_{\text{excess risk / estimation error}} + \underbrace{L(h^*)}_{\text{approximation error}}$$

population loss  
of ERM solution

excess risk /  
estimation error

Typically  
larger if  $\mathcal{H}$  is large

approximation error

Typically smaller if  
 $\mathcal{H}$  is large

$$\underbrace{L(\hat{h}) - L(h^*)}_{\text{excess risk}} = L(\hat{h}) - L_n(\hat{h}) + L_n(\hat{h}) - L_n(h^*) + L_n(h^*) - L(h^*)$$

$$\leq \underbrace{|L(\hat{h}) - L_n(\hat{h})|}_{\leq \varepsilon} + \underbrace{|L_n(h^*) - L(h^*)|}_{\leq \varepsilon}$$

$$\leq 2\varepsilon \quad \text{if uniform convergence holds}$$

$$\leq 2\varepsilon \quad \text{with probability } 1 - \delta$$

if  $e^{\log|\mathcal{H}| - 2n\varepsilon^2} \leq e^{\log\delta/2}$

$$2n\varepsilon^2 \geq \log 2|\mathcal{H}|/\delta$$

$$n \geq \frac{\log 2|\mathcal{H}|/\delta}{2\varepsilon^2}$$