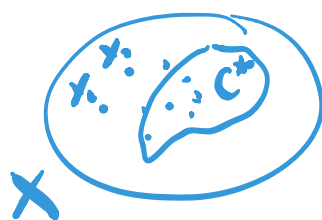


Chapter 5 continued -  
 formulation of learning problem,  
 ERM, and uniform  
 convergence of empirical  
 distributions (UCEM)

We saw  $(X, P_X, \mathcal{C})$  - realizable concept  
 learning



Observe  $Z^n = ((x_1, 1_{x_1 \in C^*}), \dots, (x_n, 1_{x_n \in C^*}))$

Learning algorithm  $A = (A_n)_{n \geq 1}$

$$A_n(Z^n) = \hat{C}$$

$$P(\hat{C} \neq C^*) = \text{error probability of using } \hat{C} \text{ on a fresh sample.}$$

$A$  is PAC if for any  $\epsilon, \delta > 0$  there  
 exist  $n(\epsilon, \delta)$  so for any  $P \in P_X$ , and  
 $C^* \in \mathcal{C}$ , if  $n \geq n(\epsilon, \delta)$  then with prob.  
 at least  $1 - \delta$ ,  $P(\hat{C}_n \neq C^*) \leq \epsilon$ , if  $n \geq n(\epsilon, \delta)$ .

Sect. 5.1.2)  $(X, P_X, \mathcal{F})$  with mean square error.

"realizable function learning"

$\mathcal{F}$  is a set of functions  $f: X \rightarrow [0, 1]$ .

$$Z^n = ((x_1, f^*(x_1)), (x_2, f^*(x_2)), \dots, (x_n, f^*(x_n)))$$

where  $f^* \in \mathcal{F}$ ,  $f^*$  is the true function.

$$A_n(Z^n) = \hat{f} \in \mathcal{F} \quad l(f(x), f^*(x)) = (f(x) - f^*(x))^2$$

$x_1, \dots, x_n$  iid  $P_X$

$$L_P(\hat{f}_n, f^*) = E_P[|\hat{f}_n(x) - f^*(x)|^2] \\ \stackrel{\uparrow}{=} A(Z^n)$$

$\mathcal{A}$  is PAC if for any  $\epsilon, \delta > 0$  there exists  $n(\epsilon, \delta)$  so for any  $P \in \mathcal{P}_X$ , any  $f^* \in \mathcal{F}$  if  $\hat{f}_n = A(Z^n)$  then with prob. at least  $1 - \delta$

$$L_P(\hat{f}_n, f^*) \leq \epsilon \quad \text{if } n \geq n(\epsilon, \delta)$$

### 5.3 Function learning agnostic (model free) case $(X, Y, U, P, \mathcal{F}, l)$

- $X$  = space of feature vectors
- $Y$  = space of labels
- $U$  = space of output labels (usually  $U=Y$ )
- $P$  = set of probability dist<sup>n</sup> on  $Z = X \times Y$ .
- $\mathcal{F}$  = set of functions  $f: X \rightarrow U$
- $l$  : loss function  $l(y, u) = \begin{matrix} \text{loss for} \\ \text{output } u \\ \text{when label is } y. \end{matrix}$

Training samples  $Z^n = (Z_1, \dots, Z_n)$   
 $Z_i = (X_i, Y_i)$  independent, dist<sup>n</sup>  $P \in \mathcal{P}$

$$A_n(Z^n) = \hat{f}_n \in \mathcal{F}.$$

Performance of  $\hat{f}_n$  on a fresh sample

$$\begin{aligned} \text{is } L_P(\hat{f}_n) &= E_P[l(Y, \hat{f}_n(X))] \\ &= \int_{X \times Y} l(y, \hat{f}_n(x)) P(dx, dy) \end{aligned}$$

$\mathcal{A}$  is PAC if for any  $\epsilon, \delta > 0$  there exists  $n(\epsilon, \delta)$  so that for any  $P \in \mathcal{P}$  if  $\hat{f}_n \in \mathcal{A}(Z^n)$  and  $n \geq n(\epsilon, \delta)$  then with probability at least  $1 - \delta$ :

$$L_P(\hat{f}_n) \leq \underbrace{L_P^*(\mathcal{F})}_{\uparrow = \min_{f \in \mathcal{F}} L_P(f)} + \epsilon$$

Section 5.3.2 Learning to classify with noisy labels.

Start with realizable model  $(X, P_X, \mathcal{C})$   
Let  $\eta \in (0, 1/2)$

$\mathcal{P}$ : distributions on  $X \times Y$  of the following form:

- select  $P \in \mathcal{P}_X$
- select  $C^* \in \mathcal{C}$

$W$  is  $\text{Ber}(\eta)$   
 $P(W=1) = \eta$   
 $P(W=0) = 1 - \eta$

$X$  has dist<sup>n</sup>  $P$ ,  $Y = 1_{\{x \in C^*\}} \oplus W$

$$A_n(\mathbb{Z}^n) = \hat{C}_n \in \mathcal{C}.$$

How good is a given  $C \in \mathcal{C}$   
 (for  $P \in \mathcal{P}_X$ ,  $C^* \in \mathcal{C}$ ,  $\eta$ )

$$\begin{aligned} \underline{L_{P_X, \mathcal{C}^*}(C)} &= P_X(\underline{1_{\{x \in C\}} \neq 1_{\{x \in C^*\}} \oplus W}) \\ &= P_X(W \neq 1_{\{x \in C^* \Delta C\}}) \\ &= (1-\eta)P_X(C^* \Delta C) \\ &\quad \eta(1 - P_X(C^* \Delta C)) \\ &= \eta + \underbrace{(1-2\eta)}_{>0} P_X(C^* \Delta C) \end{aligned}$$

