$\mathcal{X}$ :  Domain set

$\mathcal{Y}$ :  Label set

Let $D$ be a distribution over $\mathcal{X} \times \mathcal{Y}$. Let $D_x$ be the marginal distribution over $x$

We will assume that samples are drawn $(x,y) \sim D$

Let the risk be given by

$$L_D(h) := \underset{(x,y)\sim D}{\mathbb{P}}\left(h(x) \neq y\right) = D(\{(x,y): h(x) \neq y\})$$

Lemma (optimal Bayes predictor)

The classifier w/ minimum risk is

$$f_D(x) = \begin{cases} 1 & \text{if } \mathbb{P}(y=1|x) \geq \frac{1}{2} \\ 0 & \text{else} \end{cases}$$

pf

Note

$$L_D(f_D) = \underset{(x,y)\sim D}{\mathbb{E}} \mathbb{1}_{\{f_D(x) \neq y\}}$$

$$(*) \quad = \underset{x\sim D_x}{\mathbb{E}}\left(\underset{y\sim D_{y|x}}{\mathbb{E}}\left(\mathbb{1}_{\{f_D(\underline{x}) \neq y\}} \mid \underline{x}=x\right)\right)$$

Let $\alpha_x = \mathbb{P}(\underline{y}=1|\underline{x}=x)$    study this $\overset{\shortparallel}{\mathbb{P}}(f_D(\underline{x}) \neq y \mid \underline{x}=x)$

$$\mathbb{P}(f_D(\underline{x}) \neq y \mid \underline{x}=x) = \underbrace{\underbrace{\mathbb{1}_{\{\alpha_x \geq \frac{1}{2}\}}}_{\sout{\mathbb{P}(f(\underline{x})=0|\underline{x}=x)}}(1-\alpha_x)}_{\in \{0,1\}} + \underbrace{\underbrace{\mathbb{1}_{\{\alpha_x < \frac{1}{2}\}}}_{\sout{\mathbb{P}(f(\underline{x})=1|\underline{x}=x)}}\alpha_x}_{\in \{0,1\}}$$

$$= \min\left((1-\alpha_x), \alpha_x\right)$$

But also, suppose $g$ is any other classifier, possibly stochastic. Then

$$\mathbb{P}(g(\underline{x}) \neq \underline{Y} \mid \underline{X} = x) = \mathbb{P}(g(\underline{x}) = 1 \mid \underline{X} = x) \mathbb{P}(\underline{Y} = 0 \mid \underline{X} = x)$$

$$+ \mathbb{P}(g(\underline{x}) = 0 \mid \underline{X} = x) \mathbb{P}(\underline{Y} = 1 \mid \underline{X} = x)$$

$$= \mathbb{P}(g(\underline{x}) = 1 \mid \underline{X} = x)(1 - \alpha_x) + \mathbb{P}(g(\underline{x}) = 0 \mid \underline{X} = x) \alpha_x$$

$$\geq \min\left\{ (1 - \alpha_x), \alpha_x \right\}$$

Hence, ~~$\mathbb{P}(f_D(\underline{x}))$~~

going back to (*) we see that $\underset{\underline{X} \sim D_x}{\mathbb{E}}\left(1_{\{f_D(\underline{x}) = y\}} \mid \underline{X} = x\right) \leq \mathbb{E}\left(1_{\{g(\underline{x}) \neq y\}} \mid \underline{X} = x\right)$ $\forall x$

$$\cancel{(*)} = \underset{\underline{x} \sim D_x}{\cancel{\mathbb{E}}} \cancel{\min(\alpha_x, 1 - \alpha_x)}$$

$$(*) \leq \underset{\underline{x} \sim D_x}{\mathbb{E}}\left( \underset{\underline{Y} \sim D_{y\underline{x}}}{\mathbb{E}}\left(1_{\{g(\underline{x}) \neq y\}} \mid \underline{X} = x\right)\right)$$

$$= \underset{(x, y) \sim D}{\mathbb{E}}\left(1_{g(\underline{x}) \neq y}\right)$$

∎

Def ( Agnostic PAC Learnability )

Hypothesis class $\mathcal{H}$ is agnostic PAC learnable if $\exists \; m_{\mathcal{H}} : (0,1)^2 \to \mathbb{N}$ & learning algo $A$ s.t. : $\forall \; \varepsilon, \delta \in (0,1)$ & every distrib $D$ over $\mathcal{X}, \mathcal{Y}$, running $A$ on $m \geq m_{\mathcal{H}}(\varepsilon, \delta)$ i.i.d samples from $D$, $A$ returns a hypothesis $h$ s.t. w.p. at least $1-\delta$

$$L_D(h) \leq \min_{h' \in \mathcal{H}} L_D(h') + \varepsilon ,$$

where $L_D(h) = \underset{z \sim D}{\mathbb{E}} \; \lambda(h, z)$

Note: Above, we consider general loss funct $\lambda$. e.g.

$$\lambda_{01}(h, (x,y)) := \begin{cases} 0 & \text{if } h(x) = y \\ 1 & h(x) \neq y \end{cases}$$

$$\lambda_{sq}(h, (x,y)) := (h(x) - y)^2$$

Uniform Convergence

Def ( $\varepsilon$-rep. sample) A training set $S$ is called $\varepsilon$-representative if

$$|L_S(h) - L_D(h)| < \varepsilon \qquad \forall h \in \mathcal{H}$$

Lemma Assume $S$ is $\varepsilon/2$-rep. Then any $h_S \in ERM_{\mathcal{H}}(S)$ satisfies

$$L_D(h_S) \leq \min_{h \in \mathcal{H}} L_D(h) + \varepsilon$$

pf

$$\forall h, \quad L_D(h_S) \overset{\varepsilon\text{-rep}}{\leq} L_S(h_S) + \frac{\varepsilon}{2} \overset{h_S \in ERM}{\leq} L_S(h) + \frac{\varepsilon}{2} \leq L_D(h) + \varepsilon.$$

Minimizing over $h$ completes the proof.

**Def** (uniform convergence) We say a hypo class $\mathcal{H}$ has the uniform conv property if $\exists \; m_{\mathcal{H}} : (0,1)^2 \to \mathbb{N}$ s.t. $\forall \varepsilon, \delta \in (0,1)$ & $\forall$ prob distribution $D$ over $Z$, if $S$ is a sample of $m \geq m_{\mathcal{H}}(\varepsilon, \delta)$ iid draws, then w.p at least $(1-\delta)$ $S$ is $\varepsilon$-rep.

**Note** If $\mathcal{H}$ has uniform conv property w/ funck $m_{\mathcal{H}}^{uc}$, then $\mathcal{H}$ is agnostic PAC learnable w/ sample complexity $m_{\mathcal{H}}(\varepsilon, \delta) \leq m_{\mathcal{H}}^{uc}(\frac{\varepsilon}{2}, \delta)$. Also, the ERM paradigm is a successful agnostic PAC learner for $\mathcal{H}$.

**Claim** Let $\mathcal{H}$ be a finite hypothesis class, $Z$ be a domain, & let ~~$\ell: \mathcal{H} \times Z$ all~~ $\ell: \mathcal{H} \times Z \to [0,1]$ be a loss function. Then $\mathcal{H}$ has the uniform conv property w/ sample complexity

$$m_{\mathcal{H}}^{uc}(\varepsilon, \delta) \leq \frac{\log(2|\mathcal{H}|/\delta)}{2\varepsilon^2}$$

Fix $\varepsilon, \delta$. W.t.s

$$D^m\left(\{S : \forall h \in \mathcal{H}, \; |L_S(h) - L_D(h)| \leq \varepsilon\}\right) \geq 1-\delta$$

equivalently,

$$D^m\left(\{S : \exists h \in \mathcal{H}, \; |L_S(h) - L_D(h)| > \varepsilon\}\right) < \delta$$

write

$$\{S : \exists h \in \mathcal{H}, \; |L_S(h) - L_D(h)| > \varepsilon\} = \bigcup_{h \in \mathcal{H}} \{S : |L_S(h) - L_D(h)| > \varepsilon\}$$

union bd gives

$$D^m\left( \overset{\nearrow}{\phantom{xxxxxxx}} \right) \leq \sum_{h \in \mathcal{H}} \{S : |L_S(h) - L_D(h)| > \varepsilon\} \qquad (*)$$

Now, Just want to bound terms inside sum. Note that

$$ L_S(h) = \frac{1}{m} \sum_{i=1}^{m} \underbrace{l(h, z_i)}_{\text{has}} $$

$$ \underset{z \sim D}{\text{mean}} \; \mathbb{E}(l(h,z)) =: L_D(h) $$

So,

$$ \mathbb{E}(L_S(h)) = L_D(h). $$

Recall Hoeffding ineq.: $(\Theta_i)$ iid w/ $\mathbb{E}\,\Theta_1 = \mu$ &

$P(a \le \Theta_1 \le b) = 1.$ Th $\forall \varepsilon > 0,$

$$ \mathbb{P}\left( \left| \frac{1}{m} \sum_{i=1}^{m} \Theta_i - \mu \right| > \varepsilon \right) \le 2 \exp\left( \frac{-2m\varepsilon^2}{(b-a)^2} \right). $$

Applying Hoeffding, we see

$$ D^m\left( \{ S : |L_S(h) - L_D(h)| > \varepsilon \} \right) \le 2 \exp\left( \frac{-2m\varepsilon^2}{1} \right) $$

Hence, w/ (*) this implies

$$ D^m\left( \{ S : \exists h \in \mathcal{H}, |L_S(h) - L_D(h)| > \varepsilon \} \right) \le \sum_{h \in \mathcal{H}} 2 e^{-2m\varepsilon^2} $$

$$ \le 2|\mathcal{H}| e^{-2m\varepsilon^2} $$