## Motivation

CDF

$$\begin{cases} \text{ECDF} & \widehat{F}_n(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\{X_i \leq x\} \\ \text{CDF} & F(x) = \mathbb{P}(X \leq x) = \mathbb{E}[\mathbb{1}\{X \leq x\}] \end{cases}$$

$$\begin{cases} ① & \text{LLN} \Rightarrow \widehat{F}_n(x) \xrightarrow{P/a.s.} F(x) \\ ② & \text{DKW} \Rightarrow \mathbb{P}(\|F_n - F\|_\infty \geq \varepsilon) \leq 2\exp(-2n\varepsilon^2) \\ ③ & \text{G-C thm} \Rightarrow \sup_x |\widehat{F}_n(x) - F(x)| \xrightarrow{P} 0 \rightarrow \text{uniform converge} \end{cases}$$

interest: $\gamma(F)$

plug-in estimator: $\gamma(F_n)$

$\|\widehat{F}_n - F\|_\infty \xrightarrow{P} 0$

Hoeffding can only give **point-wise** bound $\mathbb{P}(|\widehat{F}_n(x) - F(x)| \geq \varepsilon) \leq C$

uniform

$\Rightarrow$ Empirical Process $\qquad X_1, \dots, X_n \sim \mathbb{P}(\cdot)$

$$\mathbb{P}_n(f) = \frac{1}{n} \sum_{i=1}^{n} f(X_i) = \int f \, d\mathbb{P}_n$$

$$\mathbb{P}(f) = \mathbb{E}_\mathbb{P}[f(X)] = \int f \, d\mathbb{P}$$

Answer this question

Our interest:
$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} f(X_i) - \mathbb{E}_F[f(X)] \right| \xrightarrow{P} 0 \quad ?$$

$$\Leftrightarrow \sup_{f \in \mathcal{F}} |\mathbb{P}_n(f) - \mathbb{P}(f)| \xrightarrow{P} 0$$

$$\Leftrightarrow \|\mathbb{P}_n - \mathbb{P}\|_{f \in \mathcal{F}} \xrightarrow{P} 0 \quad ?$$

Application: ERM

$$\begin{cases} \widehat{R}_n(\theta, \theta^*) = \frac{1}{n} \sum_{i=1}^{n} \ell_\theta(X_i, Y_i) \\ \\ R(\theta, \theta^*) = \mathbb{E}_{(X,Y) \sim \theta^*}[\ell_\theta(X,Y)] \end{cases}$$

$(X_i, Y_i) \rightsquigarrow \theta^*$

$$\begin{cases} \text{MLE} \Leftrightarrow \ell_\theta(x,y) = -\log \frac{P_{\theta^*}(x,y)}{P_\theta(x,y)} \\ \text{classify} \Rightarrow \ell_\theta(x,y) = \mathbb{1}\{f_\theta(x) \neq y\} \end{cases}$$

$$\begin{cases} \longrightarrow \hat{\theta} = \underset{\theta \in \Theta}{\text{argmin}} \ \hat{R}_n(\theta, \theta^*) \iff \hat{f} = \underset{f \in \mathcal{H}}{\text{argmin}} \ \hat{R}_n(f, f^*) \\[2em] \longrightarrow \tilde{\theta} = \underset{\theta \in \Theta}{\text{argmin}} \ R(\theta, \theta^*) \iff \tilde{f} = \underset{f \in \mathcal{H}}{\text{argmin}} \ R(f, f^*) \end{cases}$$

$\longrightarrow$ Excess Risk

Consider 
$$E(\hat{\theta}, \tilde{\theta}) = R(\hat{\theta}, \theta^*) - R(\tilde{\theta}, \theta^*)$$

$$= \boxed{(R(\hat{\theta}, \theta^*) - R_n(\hat{\theta}, \theta^*))} \rightarrow \text{r.v.}$$

$$+ \boxed{(R_n(\hat{\theta}, \theta^*) - R_n(\tilde{\theta}, \theta^*))} \rightarrow \text{optimization} \leq 0$$

$$+ \boxed{(R_n(\tilde{\theta}, \theta^*) - R(\tilde{\theta}, \theta^*))} \rightarrow \text{fixed Bound}$$

$\searrow$ capacity of Hypothesis

$\boxed{\text{need a uniform bound}}$

$$\underset{\theta \in \Theta}{\sup} \ | \ \mathbb{P}_n(\ell_\theta) - \mathbb{P}(\ell_\theta) \ |$$

$$= \boxed{\| \mathbb{P}_n - \mathbb{P} \|_{\mathcal{F}}} \qquad \text{where} \quad \mathcal{F} = \{ f: \ f = \ell(\theta) \quad \theta \in \Theta \}$$

$\boxed{\text{to solve this, VC-dim comes in}}$

$\downarrow$

$\boxed{\text{result like D-K-W}}$

$\boxed{\text{VC-dim attention}}$

it focus on the question:

$$\begin{cases} \mathbb{P}_n(A) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\{X_i \in A\} \\ \mathbb{P}(A) = \mathbb{P}(X \in A) \end{cases} \rightarrow \boxed{\underset{A \in \mathcal{A}}{\sup} \ | \ \mathbb{P}_n(A) - \mathbb{P}(A) |}$$

$\boxed{\text{apply to Binary Classification ERM framework}}$

'Toy' is because: linearity of $(x,y)$ in estimator $\boxed{\hat{F}(x) = F(x) \cdot y}$

**A Toy Example to explain** :
$\left\{\begin{array}{l} \text{Training error} \\ \text{Testing error} \end{array}\right.$

our problem consider the simple case that:

① fix $x$

② randomness comes from $\varepsilon$

3. Consider the fixed design nonparametric regression set up with observations $(x_1, Y_1), \ldots, (x_n, Y_n)$. Suppose further that the noise has constant variance, i.e. $\text{Var}\{\epsilon_i\} = \sigma^2$. Both regressograms and Nadaraya-Watson kernel estimators are examples of *linear smoothers*. This means that the vector of predicted values $\hat{\mathbf{r}} = (\hat{r}_n(x_1), \ldots \hat{r}_n(x_n))$ is given by $\hat{\mathbf{r}} = \mathbf{L}\mathbf{y}$ where $\mathbf{y} = (Y_1, \ldots, Y_n)$. The training error can therefore be written as $\frac{1}{n} \|\mathbf{L}\mathbf{y} - \mathbf{y}\|_2^2$. What is the expected training error (your answer may contain Trace($\mathbf{L}$))? What is the average predictive risk, i.e. $\mathbb{E}\left\{\frac{1}{n} \|\mathbf{L}\mathbf{y} - \mathbf{y}^*\|_2^2\right\}$ where $\mathbf{y}^* = (Y_1^*, \ldots, Y_n^*)$ are new observations at each $x_i$ (your answer may contain Trace($\mathbf{L}$))? The difference in the two values should be $2\sigma^2\text{Trace}(\mathbf{L})/n$. In other words, this is the amount by which the training error is overly optimistic. *Hint: If $M$ is any matrix, and $\epsilon$ a vector of independent random variables each with mean 0 and variance $\sigma^2$, then $\mathbb{E}\{\epsilon^T M \epsilon\} = Trace(M)\sigma^2$.*

Model: $\quad \hat{Y}_i = \hat{r}(x_i) + \varepsilon_i$

$\left\{\begin{array}{l} \text{expected training error,} \quad \mathbb{E}_\varepsilon\left[\frac{1}{n} \left\| \left(L(\hat{r}+\varepsilon) - (\hat{r}+\varepsilon)\right) \right\|_2^2\right] \\[4em] \text{expected testing error,} \quad \mathbb{E}_{\varepsilon, \varepsilon^*}\left[\frac{1}{n} \left\| L\cdot(\hat{r}+\varepsilon) - (\hat{r}+\varepsilon^*)\right\|_2^2\right] \end{array}\right.$