DSAS202   Lec8 ⟶ Uncertainty Qualification

Confidence for model prediction ?

important in sensitive domain $\{$ medical field, finance

Today's Lecture:

1. **Model inference**   (Problem formulation)

⟶ We assume that our ground-truth $(x, y)$ comes from:

$$y = f_{\theta^*}(x) \qquad f_{\theta^*}(\cdot) \longrightarrow \text{oracle function}$$

$$\Rightarrow \mathcal{D} = \{(x_i, y_i)\}_{i=1}^{N} = \{(x_i, f_{\theta^*}(x_i))\}_{i=1}^{N}$$

⟶ our estimator of $\theta^*$ is $\hat{\theta}$, comes from:

ERM ⟵ $$\boxed{\hat{\theta} = \arg\min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^{N} \ell(f_\theta(x_i), \underbrace{f_{\theta^*}(x_i)}_{y_i \swarrow})}$$

Hypothesis Space $\mathcal{H} = \{f : f(x) = f_\theta(x), \theta \in \Theta\}$

⟶ parametric model like NN

⟶ We solve the ERM approximately via $\{$ GD, SGD, mGD

⇓

Optimization perspective

⇓

achieve $\hat{\theta} \approx \theta_n$

Question : when we achieve $\theta_n (\hat{\theta})$, we might be interested :

① the distribution of $\hat{\theta} / g(\hat{\theta})$

② the distribution of our prediction $f_{\hat{\theta}}(x)$

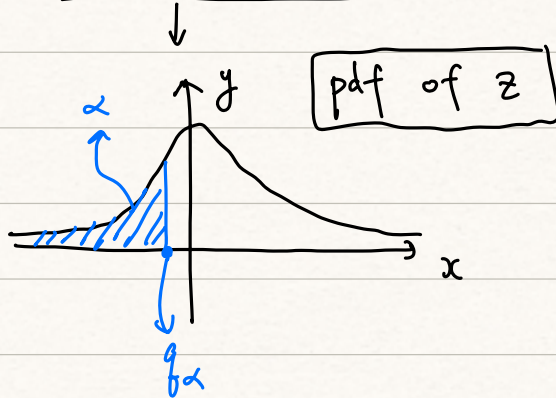↳ $\hat{\theta} = g(z_1 \ldots, z_n)$  $\widehat{z_i}$

→ Today's content

→ $\theta^*$ is constant

Solution : { Frequentist approach : confidence interval for $g(\theta^*)$

Bayesian approach :

our interest is :
$g(\theta) = f_\theta(z)$

conditional distribution

## 2. Definition of Quantile $q_\alpha$ for R.V. $z$

$\Longrightarrow q_\alpha \longrightarrow \underbrace{P(z \leq q_\alpha)}_{\text{cdf}} = \alpha$

(plug-in estimator)
estimator of $q_\alpha$.

$\dfrac{\sum_{i=1}^{n} \mathbb{1}\{z_i \leq \hat{q}_\alpha\}}{n} = \alpha$

[e.g.]. $\underline{z \sim N(0,1)}$

$\boxed{\text{ECDF}}$


pdf of $z$

→ plug-in estimator : use ECDF to replace CDF

$$\underline{\text{CDF}}$$
$$q_\alpha :\Rightarrow F(q_\alpha) = \alpha \longrightarrow$$

$$\underline{\text{ECDF}}$$
$$\hat{q}_\alpha :\Rightarrow \hat{F}(\hat{q}_\alpha) = \alpha$$

$$F(x) = P(X \leq x) \longrightarrow \hat{F}_n(x) = \dfrac{\sum_{i=1}^{n} \mathbb{1}\{x_i \leq x\}}{n}$$

$$\boxed{X_i \xrightarrow{\text{i.i.d}} F}$$

Definition of <mark>Confidence Interval ( C I )</mark>   $\underline{\underline{C_\alpha}}$

$\implies \theta \longrightarrow$ population parameter for $\underline{X \sim F(\cdot)}$

$\boxed{X_i \sim F(\cdot)}$ i.i.d sample from population

<mark>$C_\alpha \longrightarrow$   $\mathbb{P}( \theta \in \underline{C_\alpha(X_1,\ldots,X_n)}) = \alpha$</mark>
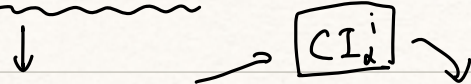
$\iff \mathbb{P}( L_\alpha(X_1,\ldots,X_n) \le \theta \le U_\alpha(X_1,\ldots,X_n)) = \alpha$

☆ <mark>Note:</mark> In statistical inference, the <mark>CI measures</mark> that:

$\Rightarrow$ if we repeat the estimation for $\underline{B \text{ times}}$, that is :

$\underline{X_{i1},\ldots, X_{in}} \sim F(\cdot)$   $i = 1, 2, 3,\ldots B$

$\downarrow$   $\boxed{CI_\alpha^i}$

$\underline{\text{Confidence Interval}}$   $[L_\alpha(X_{i1},\ldots X_{in}), U_\alpha(X_{i1},\ldots X_{in})]$   $i \in [B]$

<mark>$$\Rightarrow \quad \frac{\sum_{i=1}^{B} \mathbb{1}\{\theta \in CI_\alpha^i\}}{B} \xrightarrow[a.s.]{p} \mathbb{P}(\theta \in CI_\alpha) = \alpha$$</mark>

3. $\underline{\text{Frequentist Approach}}$ ( Confidence Interval )

$\longrightarrow \boxed{\text{idea}}$: for <mark>our model $f_{\hat\theta}(\cdot)$</mark> and <mark>$x_{new}$</mark>, we want an $\underline{\text{interval prediction}}$ CI such that $\underline{\mathbb{P}( f_{\theta^*}(x_{new}) \in CI) \ge 1 - \alpha}$

$\downarrow$

$\boxed{\alpha \longrightarrow \text{confidence level}}$

$\underline{\text{Assumption}}$ : 1. $g(\hat\theta) \xrightarrow{\text{unbiased estimator}} g(\theta^*) \iff \underline{\mathbb{E}[g(\hat\theta)] = g(\theta^*)}$

2. $g(\hat\theta) \implies \underline{\text{Gaussian Distributed}}$ R.V.

$$\Downarrow$$

$$g(\hat{\theta}) \sim \text{Gaussian} \left( \mathbb{E}[g(\hat{\theta})], \text{Var}[g(\hat{\theta})] \right)$$
$$= \text{Gaussian} \left( g(\theta^*), \mathbb{E}[\{g(\hat{\theta}) - g(\theta^*)\}^2] \right)$$

Therefore, we just need to figure out $\sigma^2 = \mathbb{E}[\{g(\hat{\theta}) - g(\theta^*)\}^2]$

$$\downarrow$$

Then, 95% CI of $\boxed{g(\theta^*)}$ is:

$$[ g(\hat{\theta}) - \boxed{1.96}\,\sigma \,, \, g(\hat{\theta}) + \boxed{1.96}\,\sigma ]$$

$$z_{\frac{\alpha}{2}} \qquad\qquad z_{\frac{\alpha}{2}}$$

Remark: we have $\mathbb{P}\left( g(\theta^*) \in [g(\hat{\theta}) - 1.96\,\sigma \,,\, g(\hat{\theta}) + 1.96\,\sigma] \right) = 0.95$

$$\downarrow$$

$$\boxed{\text{true value}} \qquad\qquad\qquad \boxed{\text{Sampling distribution}}$$

Issue: How to achieve $\sigma$ ?  $\boxed{\sigma^2 = \text{Var}[g(\hat{\theta})]}$

$$\downarrow$$

$$\boxed{\text{standard error of estimator } g(\hat{\theta})}$$

Solution:
$\begin{cases} 1. \text{ exact method for } \underline{\text{sampling distribution } g(\hat{\theta})} \\ 2. \text{ Bootstrap for } \underline{\text{sampling distribution } g(\hat{\theta})} \\ 3. \text{ exact method for } \underline{\text{standard error } \text{Var}[g(\hat{\theta})]} \\ 4. \text{ Bootstrap for } \underline{\text{standard error } \text{Var}[g(\hat{\theta})]} \end{cases}$

using asymptotic distribution of $g(\hat{\theta})$

$$\boxed{g(\hat{\theta}) \sim N\left( \mathbb{E}[g(\hat{\theta})], \text{Var}[g(\hat{\theta})] \right)}$$

4. Standard Error $\underline{\text{Var}[g(\hat{\theta})]}$     $\boxed{\text{Exact Method}}$

$$\rightarrow \text{Here, we focus on } g(\hat{\theta}) = f_{\hat{\theta}}(x_{new})$$

① Linear (Regression) case

consider $g(\hat{\theta}) = a^{T}\hat{\theta}$

$\longrightarrow$ $\text{var}[g(\hat{\theta})] = \text{var}[a^{T}\hat{\theta}]$

Denote $\Sigma_{\theta} = \text{cov}[\hat{\theta}]$

$\quad = a^{T} \text{cov}[\hat{\theta}] \cdot a$

$\quad = a^{T} \Sigma_{\hat{\theta}} \cdot a$

Assumption in Linear Regression :

$$y_i = f^{*}(x_i) + \varepsilon_i \longrightarrow \text{all the randomness comes from this term}$$

$\quad\quad x_i$ is fixed & non-random

$\longrightarrow$ in Linear Regression : $\quad f^{*}(x) = x^{T}\beta^{*}$

$$y_i = f^{*}(x_i) + \varepsilon_i$$

$f_{\hat{\beta}}(x) = x^{T}\hat{\beta} \longrightarrow$ if we have $\text{cov}[\hat{\beta}]$, then we are done!

measure the
uncertainty with respect to current dataset

$$\boxed{\text{cov}[\hat{\beta}]} = \text{cov}[(X^{T}X)^{+}X^{T}y]$$

$$\begin{cases} y = X\beta^{*} + \varepsilon \\ \mathbb{E}[y] = X\beta^{*} \end{cases}$$

$$= \text{cov}[(X^{T}X)^{-1}X^{T}\varepsilon]$$

$$= \sigma^{2}(X^{T}X)^{-1} \quad\quad \hat{\sigma}^{2} = \frac{SSE}{n-p+1}$$

$$\Longrightarrow \text{Var}[f_{\hat{\beta}}(x)] = \sigma^{2} x^{T}(X^{T}X)^{-1}x$$

$\quad\quad\quad\quad\quad \hookrightarrow$ exact variance for $f_{\hat{\beta}}(x)$

5. Standard Error $\text{Var}[g(\hat{\theta})]$ [ Bootstrap ]

$\quad\quad\quad\hookrightarrow$ for general function $g(\cdot)$, not just linear case

issue: ① do not know $\hat{\theta} \sim \underline{?} \longrightarrow$ no access to $\underline{\Sigma_{\hat{\theta}}}$

$\quad\quad$ ② even if we know $\Sigma_{\hat{\theta}}$, we cannot infer

$\quad\quad\quad\quad \text{Var}[g(\hat{\theta})] \longleftrightarrow \Sigma_{\hat{\theta}}$

$\quad\quad\quad$ due to the $\underline{\text{non-linearity of } g(\cdot) \text{ w.r.t } \hat{\theta}}$

$\quad\quad$ ✗✗✗

Surrogate : approximate $\text{Var}[g(\hat{\theta})]$

Pipeline : $S = \{(x_i, y_i)\}_{i=1}^{N}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \rightarrow$ draw with replacement $S_b$

$\quad\quad$ 1. draw $\underline{n_b}$ samples from $S$, $\boxed{b = 1, 2, \dots, B}$

$\quad\quad$ 2. Train model to achieve $\hat{\theta}$ from $S_b$

$\quad\quad\quad\quad\quad\quad\quad \Downarrow$

$\quad\quad\quad$ get model $\hat{f}_{sb}(x)$ $\quad \boxed{b = 2, 2, \dots, B}$

$\quad\quad$ 3. $\text{Var}[\underline{g(\hat{\theta})}] \approx \dfrac{1}{B} \sum\limits_{b=1}^{B} \left\{ \hat{f}_{sb}(x) - \overline{\hat{f}(x)} \right\}^2$

$\quad\quad\quad g(\theta) = f_{\hat{\theta}}(x) := f_S(x)$

$\quad\quad\quad$ since $\hat{\theta}$ is attained via $\underline{\text{dataset } S}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ bootstrap $\quad\quad \boxed{\overline{\hat{f}(x)} = \dfrac{1}{B} \sum\limits_{b=1}^{B} \hat{f}_{sb}(x)}$

# 6. Fisher Information for <u>MLE asymptotic distribution</u>

<u>Problem Setting</u> : $\qquad X_i \sim P_\theta \qquad i = 1, 2, \ldots, n.$

$\Rightarrow$ <u>likelihood</u> $\qquad \mathcal{L}(\theta ; \mathcal{D}) = \prod_{i=1}^{n} p(X_i \mid \theta)$

$\Rightarrow$ <u>log-likelihood</u> $\qquad \ell(\theta) = \sum_{i=1}^{n} \log p(X_i \mid \theta)$

$$:= \sum_i \ell_i(\theta)$$

$\Rightarrow$ MLE estimator $\quad \hat{\theta}_{MLE} := \underset{\theta \in \Theta}{\arg\max} \; \ell(\theta)$

$$= \underset{\theta \in \Theta}{\arg\max} \; \sum_{i=1}^{n} \log p(X_i \mid \theta)$$

## Notation

① Scoring Function :

$$S(\theta) = \nabla_\theta \ell(\theta)$$
$$= \sum_{i=1}^{n} S_i(\theta) := \sum_{i=1}^{n} \nabla_\theta \ell_i(\theta)$$

② Fisher Information:

$$I(\theta) = \mathbb{E}\left[ S(\theta) \cdot S(\theta)^T \right]$$

(independence between $X_i$) $\qquad = \underline{n \, \mathbb{E}\left[ S_i(\theta) \cdot S_i(\theta)^T \right]}$

$$= n \, \mathbb{E}_X\left[ \nabla_\theta \log p(x \mid \theta) \cdot \nabla_\theta \log p(x \mid \theta)^T \right]$$

$$\boxed{X \sim p(\cdot \mid \theta)}$$

Result 1: under some regularity condition,

$$I(\theta) = \mathbb{E}_{x \sim \theta} [S(\theta) \cdot S(\theta)^T]$$

$$= \mathbb{E}_{x \sim \theta} [\nabla \ell(\theta) \cdot \nabla \ell(\theta)^T]$$

$$= -\mathbb{E}_{x \sim \theta} [\nabla^2 \ell(\theta)]$$

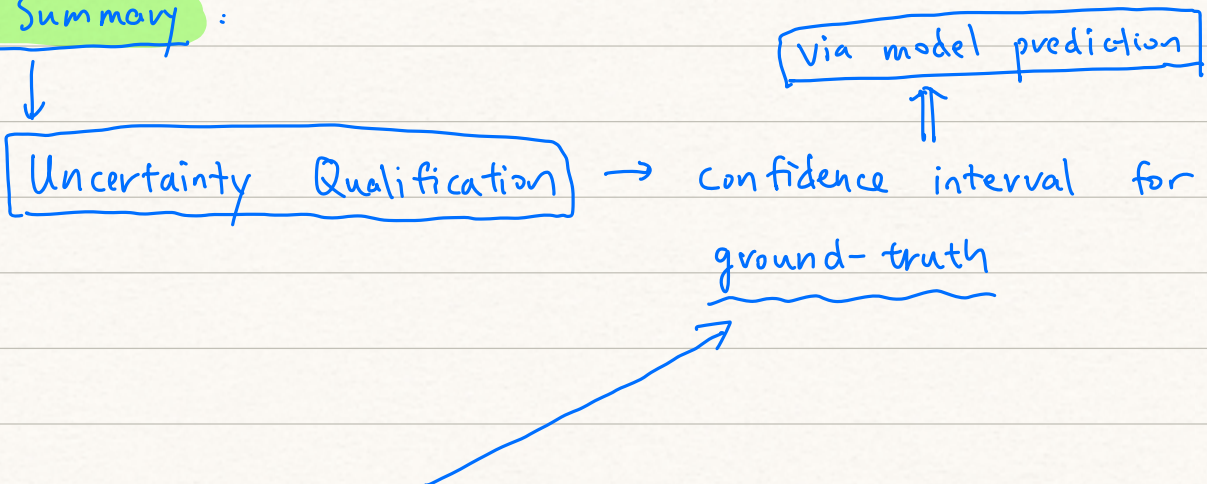Result 2: Asymptotic distribution for $\hat{\theta}_{MLE}$

$$\hat{\theta}_{MLE} \sim N(\theta^*, I(\theta^*)^{-1})$$

→ asymptotic normal

In practice, we use $\hat{\theta}_{MLE} - \theta^* \approx N(0, I(\hat{\theta}_{MLE})^{-1})$

In Summary:

↓

Uncertainty Qualification → confidence interval for

via model prediction
↑

ground-truth

From Frequentist Perspective: one general approach is:

Bootstrap → simulate the randomness via { empirical cdf

model