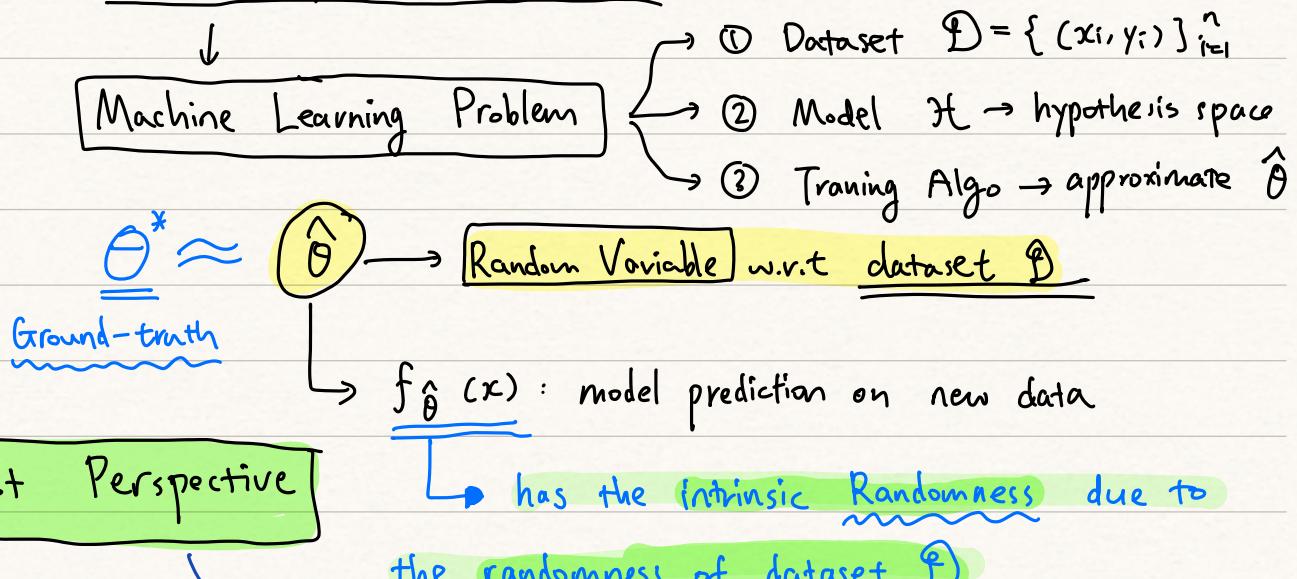


Recap : → Uncertainty Quantification



Example : ① $E[X] = \theta^*$

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i$$



asymptotic

⇒ from CLT, we have $\hat{\theta} \sim \text{Gaussian}(\theta^*, \frac{1}{n} \text{Var}[X])$

→ uncertainty quantification

$$\Rightarrow P(\theta^* \in [\text{std}(\hat{\theta} - g_{\frac{1-\alpha}{2}}), \text{std}(\hat{\theta} - g_{\frac{\alpha}{2}})]) = 1-\alpha$$

→ Here, randomness comes from $\mathcal{D} = \{X_i\}_{i=1}^n$

Example : ② $y = X\beta^* + \varepsilon$ $f_{\hat{\beta}}(x) = x^T \hat{\beta}$ (Linear Reg)

Here, $\hat{\beta} = (X^T X)^{-1} X^T Y$ → randomness comes from $Y (\varepsilon)$

$$\hat{\beta} \sim \text{Gaussian}(\beta^*, \sigma^2 (X^T X)^{-1})$$

$$\Rightarrow f_{\hat{\beta}}(x) = x^T \hat{\beta} \sim \text{Gaussian}(x^T \beta^*, \sigma^2 x^T (X^T X)^{-1} x)$$

$$\Rightarrow P(f^*(x) = x^T \beta^* \in [\dots, \dots]) = 1-\alpha$$

→ uncertainty quantification

Example : ③ \rightarrow Bootstrap

we have $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n + \text{Hypothesis Space } \mathcal{H}$

optimal hyper-parameter $\hat{\theta} \approx \theta_n$

Basically, $f_{\hat{\theta}}(x)$ can be viewed as the random variable of $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

model prediction $f_{\hat{\theta}}(x)$

use Bootstrap to estimate

{ variance
confidence interval }

draw with replacement

$\{(X_{i1}, Y_{i1}), \dots, (X_{in}, Y_{in})\}_{i=1}^B$

$\hookrightarrow \{f_{\hat{\theta}_i}(x)\}_{i=1}^B$

★ this can also construct C.I under normal assumption

$$\Rightarrow \textcircled{1} \underbrace{\text{variance}}_{\text{Var}[f_{\hat{\theta}}(x)]} \approx \sum_{i=1}^B (f_{\hat{\theta}_i}(x) - \overline{f_{\hat{\theta}}(x)})^2$$

$$\overline{f_{\hat{\theta}}(x)} = \frac{1}{B} \sum_{i=1}^B f_{\hat{\theta}_i}(x)$$

$$\Rightarrow \textcircled{2} \underbrace{\text{confidence interval}}_{\{f_{\hat{\theta}_i}(x) - f_{\hat{\theta}}(x)\}_{i=1}^B} \rightarrow$$

$$\downarrow \quad \left\{ \begin{array}{l} q_{\frac{\alpha}{2}} = f_{\hat{\theta}_{\frac{\alpha}{2}}} - f_{\hat{\theta}} \\ q_{1-\frac{\alpha}{2}} = f_{\hat{\theta}_{1-\frac{\alpha}{2}}} - f_{\hat{\theta}} \end{array} \right. \text{quantile}$$

$$\Pr(f_{\hat{\theta}}(x) - f_{\theta^*}(x) \leq t) \approx \frac{1}{B} \sum_{i=1}^B \mathbb{1}\{f_{\hat{\theta}_i}(x) - f_{\hat{\theta}}(x) \leq t\}$$

$$\Rightarrow \Pr(f_{\theta^*}(x) \in [2f_{\hat{\theta}}(x) - f_{\hat{\theta}_{1-\frac{\alpha}{2}}}(x), 2f_{\hat{\theta}}(x) - f_{\hat{\theta}_{\frac{\alpha}{2}}}(x)]) = 1 - \alpha$$

To conclude: from Frequentist Perspective:

- ① all randomness comes from the data \mathcal{D} itself!
- ② Ground-truth parameters are fixed!

Today's lecture: Bayesian Approach $\rightarrow \beta$ is random!

1. Bayes Theorem

$$\begin{aligned}
 P(\beta | \mathcal{D}) &= \frac{\text{likelihood} \cdot \text{prior distribution}}{P(\mathcal{D})} \\
 &= \frac{P(\mathcal{D}|\beta) \cdot P(\beta)}{\int_{\beta} P(\mathcal{D}|\beta) \cdot P(\beta) d\beta} \\
 &\propto P(\mathcal{D}|\beta) \cdot P(\beta)
 \end{aligned}$$

↳ our interest to determine the posterior distribution $P(\beta | \mathcal{D})$

Advantages:

① easy to combine prior knowledge

② robust against over-fitting \leftarrow prior distribution

③ give the whole distribution of β

encode knowledge to distribution

\uparrow regularization

$$\boxed{P(\beta | \mathcal{D})} \Rightarrow \underline{\text{easy to construct}} \quad (1)$$

Disadvantages:

① Expensive to achieve closed-form sol²

need approximation via MCMC

② sensitive to prior choices

2. Example: Bayesian Linear Regression

$$\text{Setting: } \begin{cases} \beta \sim N(\mu_0, \Sigma_0) \\ y | X, \beta \sim N(X\beta, \Sigma_\epsilon) \end{cases}$$

$$① \text{ likelihood: } P(\mathcal{D} | \beta) = P(y | X, \beta)$$

$$\propto \exp \left\{ - (y - X\beta)^T \Sigma_\epsilon^{-1} (y - X\beta) \right\}$$

$$② \text{ prior: } P(\beta) = \exp \left\{ - (\beta - \mu_0)^T \Sigma_0^{-1} (\beta - \mu_0) \right\}$$

$$\Rightarrow ③ \text{ posterior: } P(\beta | \mathcal{D}) \propto P(\mathcal{D} | \beta) \cdot P(\beta)$$

$$= \underbrace{\exp \left\{ - (y - X\beta)^T \Sigma_\epsilon^{-1} (y - X\beta) \right\}}_{\times}$$

$$\underbrace{\exp \left\{ - (\beta - \mu_0)^T \Sigma_0^{-1} (\beta - \mu_0) \right\}}_{\times}$$

Simplification: $\exp \{ -(\gamma - X\beta)^T \Sigma_{\varepsilon}^{-1} (\gamma - X\beta) - (\beta - \mu_0)^T \Sigma_0^{-1} (\beta - \mu_0) \}$

$$\propto \exp \{ -\beta^T X^T \Sigma_{\varepsilon}^{-1} X \beta - \beta^T \Sigma_0^{-1} \beta$$

$$+ 2 \gamma^T \Sigma_{\varepsilon}^{-1} X \beta + 2 \mu_0^T \Sigma_0^{-1} \beta \}$$

$$= \exp \{ -\beta^T (X^T \Sigma_{\varepsilon}^{-1} X + \Sigma_0^{-1}) \beta + 2 (\gamma^T \Sigma_{\varepsilon}^{-1} X + \mu_0^T \Sigma_0^{-1}) \beta \}$$

if $X \sim N(\mu, \Sigma)$, $p(x) \propto \exp \{ -x^T \Sigma^{-1} x + 2\mu^T \Sigma^{-1} x \}$

$$\Rightarrow \begin{cases} \Sigma^{-1} = (X^T \Sigma_{\varepsilon}^{-1} X + \Sigma_0^{-1}) \\ \mu^T \Sigma^{-1} = \gamma^T \Sigma_{\varepsilon}^{-1} X + \mu_0^T \Sigma_0^{-1} \end{cases}$$

$$\Rightarrow \begin{cases} \Sigma = (X^T \Sigma_{\varepsilon}^{-1} X + \Sigma_0^{-1})^{-1} \\ \mu = (X^T \Sigma_{\varepsilon}^{-1} X + \Sigma_0^{-1})^{-1} (X^T \Sigma_{\varepsilon}^{-1} \gamma + \Sigma_0^{-1} \mu_0) \end{cases}$$

Therefore, we can conclude that:

for Bayesian LR. under the setting that:

$$\begin{cases} \beta \sim N(\mu_0, \Sigma_0) \\ y | X, \beta \sim N(X\beta, \Sigma_{\varepsilon}) \end{cases}$$

(special case)

posterior dist. is exact Gaussian for Bayesian

we have :

$$\boxed{\beta | \mathcal{D} \sim N(\mu_p, \Sigma_p)}$$

LR

$$\text{Here, } \begin{cases} \mu_p = (X^T \Sigma_{\varepsilon}^{-1} X + \Sigma_0^{-1})^{-1} (X^T \Sigma_{\varepsilon}^{-1} y + \Sigma_0^{-1} \mu_0) \\ \Sigma_p = (X^T \Sigma_{\varepsilon}^{-1} X + \Sigma_0^{-1})^{-1} \end{cases}$$

Issue: ① Generally speaking, we cannot attain the exact solution for posterior distribution \rightarrow disadvantage 1

① If our choice of β (prior distribution) is far from reality, then the achieved $\beta | \mathcal{D}$ (posterior distribution) will be twisted.

disadvantage 2

sensitive to prior choices

3. Monte Carlo Simulation Outline

Recap: In Bayesian LR, based on $\begin{cases} \beta \sim N(\mu_0, \Sigma_0) \\ y|x, \beta \sim N(X\beta, \Sigma_e) \end{cases}$

we can achieve $\beta | \mathcal{D} \sim N(\mu_p, \Sigma_p)$

exact posterior distribution

→ ① we can attain exact $\begin{cases} \text{expectation} & \mu_p \\ \text{variance} & \Sigma_p \end{cases}$

(if our interest is $\begin{cases} \mathbb{E}[\beta] \\ \text{Var}[\beta] \end{cases}$ (posterior))

→ ② we can sample from it easily

* sometimes even if we have $g(\beta)$, it is expensive to do integral

(if our interest is $\begin{cases} \mathbb{E}[g(\beta)] \\ \text{Var}[g(\beta)] \end{cases}$ (posterior), now we have no access to the posterior distribution for $g(\beta)$ explicitly)

$\Rightarrow \beta_1, \dots, \beta_n \sim \beta | \mathcal{D}$ (since $\beta | \mathcal{D}$ is good enough)

approximate

$\begin{cases} \mathbb{E}[g(\beta)] \\ \text{Var}[g(\beta)] \end{cases}$

via Monte Carlo

exact Gaussian

→ ③ Sometimes it is difficult to sample from $\beta | \mathcal{D}$

directly \Rightarrow use some Markov Chain Monte Carlo



sample from $\beta | \mathcal{D}$ approximately

[E.g.] Metropolis-Hastings Algo [M-H Algo]
Gibbs Sampling

4. Monte Carlo Approach



Our interest is like $\begin{cases} \mathbb{E}[f(\beta) | \mathcal{D}] \\ \text{Var}[f(\beta) | \mathcal{D}] \end{cases}$

Assume: ① it is intractable to compute $\int_{\beta} f(\beta) \cdot p(\beta | \mathcal{D}) d\beta$

directly

② We can sample from $\beta | \mathcal{D}$ directly

Solution: $\mathbb{E}[f(\beta) | \mathcal{D}] = \int_{\beta} f(\beta) p(\beta | \mathcal{D}) d\beta$

e.g. $f(\beta) = h_{\beta}(x) \Rightarrow h_{\beta}(x)$: neural network prediction

$$\mathbb{E}_{\beta | \mathcal{D}}[f(\beta)] = \mathbb{E}_{\beta | \mathcal{D}}[h_{\beta}(x)]$$

no closed-form solution

\Rightarrow approximate through $\mathbb{E}_{\beta}[f(\beta) | \mathcal{D}]$

$$\approx \hat{f}_m = \frac{1}{m} \sum_{i=1}^m f(\beta^i) \quad \beta^i \sim p(\beta | \mathcal{D})$$

① From LLN, we have $\hat{f}_m \xrightarrow[\text{a.s.}]{P} \mathbb{E}_{\beta} [f(\beta) | \mathcal{D}] \quad \underline{m \rightarrow \infty}$

② if $\beta^i \sim p(\beta | \mathcal{D})$ i.i.d. then $\text{var}[\hat{f}_m] = \frac{1}{m} \text{var}_{\beta} [f(\beta) | \mathcal{D}]$

From CLT, we have :

$$\hat{f}_m - \mathbb{E}_{\beta} [f(\beta) | \mathcal{D}] \approx N(0, \frac{1}{m} \text{var}_{\beta} [f(\beta) | \mathcal{D}])$$

$\frac{1}{\sqrt{m}}$ convergence rate

Remark: for MCMC method, generally we do not have this property, which means the variance will be slightly bigger.

From autocorrelation to ESS

$$\begin{aligned} \text{Var}\left\{\frac{1}{T} \sum_{t=1}^T f(\theta^{(t)})\right\} &= \frac{1}{T^2} \sum_{t=1}^T \left(\text{Var}\{f(\theta^{(t)})\} + \sum_{s \neq t} \text{Cov}\{f(\theta^{(t)}), f(\theta^{(s)})\} \right) \\ &\stackrel{\text{if independent, then } ESS = T!}{=} \frac{\text{Var}\{f(\theta^{(0)})\}}{T^2} \sum_{t=1}^T \left(1 + \sum_{s \neq t} \text{Corr}\{f(\theta^{(t)}), f(\theta^{(s)})\} \right) \\ &= \frac{\text{Var}\{f(\theta^{(0)})\}}{T^2} \sum_{t=1}^T \left(1 + \sum_{s \neq t} \rho_{s-t} \right) \\ &= \frac{\text{Var}\{f(\theta^{(0)})\}}{T^2} \sum_{t=1}^T \left(1 + \sum_{s \neq t} \rho_{|s-t|} \right) \\ &\approx \frac{\text{Var}\{f(\theta^{(0)})\}}{T^2} \sum_{t=1}^T \left(1 + 2 \sum_{L=1}^{\infty} \rho_L \right) \\ &= \frac{\text{Var}\{f(\theta^{(0)})\}}{T(1 + 2 \sum_{L=1}^{\infty} \rho_L)^{-1}} \quad \text{ESS} \end{aligned}$$

Idea: use Autocorrelation to calculate this term!

can estimate with plugin estimators



Remark: when dimensionality is high, we prefer Frequentist Approach

$\beta \in \mathbb{R}^p$ and p is large

since we need : $\mathbb{E}_{\beta} [f(\beta) | \mathcal{D}] \approx \frac{1}{m} \sum_{i=1}^m f(\beta^i)$ for Bayesian

↳ works well in low-dimension case

similarly : $\text{var}_{\beta} [f(\beta) | \mathcal{D}] = \int_{\beta} (f(\beta) - \mathbb{E}_{\beta} [f(\beta) | \mathcal{D}])^2 P(\beta | \mathcal{D}) d\beta$

$$\approx \frac{1}{m} \sum_{i=1}^m (f(\beta^i) - \hat{f}_m)^2$$

$$\begin{cases} \hat{f}_m = \frac{1}{m} \sum_{i=1}^m f(\beta^i) \approx \mathbb{E}_{\beta}[f(\beta) | \mathcal{D}] \\ \beta^i \sim p(\beta | \mathcal{D}) \end{cases}$$

5. Application of Monte Carlo method

Model : $\begin{cases} \beta \sim N(0, \tau I) \\ y | x, \beta \sim N(H\beta, \sigma^2 I) \end{cases}$

$$H := \begin{pmatrix} h(x_1)^T \\ \vdots \\ h(x_n)^T \end{pmatrix}$$

→ General Linear Model

$$\Rightarrow \beta | \mathcal{D} \sim N(\mu_p, \Sigma_p)$$

where $\begin{cases} \mu_p = (H^T H + \sigma^{-2} I)^{-1} H^T y \\ \Sigma_p = (\sigma^{-2} H^T H + I)^{-1} \end{cases}$

Ridge Regression estimator

→ if apply Maximize A Posterior (MAP) framework.

then $\hat{\beta}_{MAP} = (H^T H + \sigma^{-2} I)^{-1} H^T y$

↳ Ridge Regularization

⇒ Moreover, we have :

$$p(y | x_{new}, \mathcal{D})$$

$$= \int_w p(y, w | x_{new}, \mathcal{D}) dw$$

$$\xrightarrow{N(w^T h(x_{new}) - \beta^T)} \xrightarrow{N(\mu_p, \Sigma_p)}$$

$$= \int_w p(y | w, x_{new}) \cdot p(w | \mathcal{D}) dw$$

$$\Rightarrow p(y | x_{new}, \mathcal{D}) \sim \underline{N(\mu_p^T h(x_{new}), \beta^{-1} + h(x_{new})^T \Sigma_p h(x_{new}))}$$

$(\text{std})^2$

construct "confidence interval"

$$[\mu_p^\top h(x_{\text{new}}) - \text{std}, \mu_p^\top h(x_{\text{new}}) + \text{std}]$$

Generally when our model is $y = f_w(x)$

and we have $w | \mathcal{D} \sim \text{some distribution}$

Issue: ① it is difficult to determine

$y | x_{\text{new}}, \mathcal{D}$ distribution

② cannot compute $\begin{cases} \mathbb{E}[y | x_{\text{new}}, \mathcal{D}] \\ \text{Var}[y | x_{\text{new}}, \mathcal{D}] \end{cases}$ explicitly

\Rightarrow Monte Carlo to approximate $\begin{cases} \mathbb{E}[y | x_{\text{new}}, \mathcal{D}] \rightarrow \hat{\mu} \\ \text{Var}[y | x_{\text{new}}, \mathcal{D}] \rightarrow \hat{\text{std}} \end{cases}$

then, the approximate "Confidence Interval" is:

$$[\hat{\mu} - \hat{\text{std}}, \hat{\mu} + \hat{\text{std}}]$$

two case

① $w | \mathcal{D} \rightarrow$ sample directly
like Bayesian LR case

$$w | \mathcal{D} \sim N(\mu_p, \Sigma_p)$$

② $w | \mathcal{D} \rightarrow$ cannot sample directly
 \Rightarrow use MC MC to sample approximately

Bayesian inference vs frequentist inference

- Two fundamentally different ways of thinking about inference

Frequentist

- Parameter is fixed
- Data is random
- Randomness represents sampling

$X \sim \text{random}$

Bayesian

- Parameter is random
- Data is fixed
- Randomness represents uncertainty

$\Theta \sim \text{random}$

super-rare event
not reasonable in frequentist framework

Bayesian perspective allows us to reason about questions like "what is the probability of an earthquake?" (after some observation)