

3 - Asymptotic Statistics

- Supervised learning setup:

- (y, x) pairs $\sim p(y, x)$ joint distribution

$y \in \mathcal{Y}$ is response, $x \in \mathcal{X}$ is inputs (feature)
 $\in \mathbb{R}$ or $\{0, 1\}$ or $\{\pm 1\}$, $x \in \mathbb{R}^d$

- Observe data: $(y_i, x_i) \stackrel{\text{iid}}{\sim} p(y, x)$ for $i=1, 2, \dots, n$.

Goal: Find a function $f: \mathcal{X} \rightarrow \mathcal{Y}$ s.t.

$$y \approx f(x)$$

using data.

! This is the focus of next six lectures!

- Need to measure the quality of the function f .

* Loss function: $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

- measures distance between y and $f(x)$

- Ex: $l(y, f(x)) = (y - f(x))^2$ or $|y - f(x)|$

* Risk function: $R: \mathcal{F} \rightarrow \mathbb{R}_+$ for $\mathcal{F} = \{f\}$.

- measures expected loss of function f .

$$R(f) = \mathbb{E}[l(y, f(x))]$$

$(y, x) \sim p(y, x)$

- Depends on l, p, f .

Goal (revised): Find $f \in \mathcal{F}$ s.t. $R(f)$ is small.

* Bias - Variance Decomposition

$$\mathbb{E}[\mathbb{E}[y|x]] = \mathbb{E}[y]$$

- Squared error loss: $l(a, b) = (a - b)^2$

$$l(y, f(x)) = (y - f(x))^2$$

- Risk: $R(f) = \mathbb{E}[(y - f(x))^2] = \mathbb{E}[\mathbb{E}[(y - f(x))^2 | x]] + \mathbb{E}[y|x]$

$$\begin{aligned}
R(f) &= \mathbb{E} \left[\mathbb{E} \left[\underbrace{(y - \mathbb{E}[y|x])^2}_{\text{I}} + \underbrace{(\mathbb{E}[y|x] - f(x))^2}_{\text{II}} \mid x \right] \right] \\
&= \mathbb{E} \left[\mathbb{E} \left[\underbrace{(y - \mathbb{E}[y|x])^2}_{\text{I}^2} + \underbrace{(\mathbb{E}[y|x] - f(x))^2}_{\text{II}^2} \right. \right. \\
&\quad \left. \left. + 2 \cdot \underbrace{(y - \mathbb{E}[y|x]) (\mathbb{E}[y|x] - f(x))}_{2 \text{I} \cdot \text{II}} \mid x \right] \right] \\
&= \mathbb{E} \left[\mathbb{E} \left[(y - \mathbb{E}[y|x])^2 \mid x \right] + \mathbb{E} \left[(\mathbb{E}[y|x] - f(x))^2 \mid x \right] \right. \\
&\quad \left. + 2 \underbrace{\mathbb{E} \left[(y - \mathbb{E}[y|x]) (\mathbb{E}[y|x] - f(x)) \mid x \right]}_{\substack{\text{are} \\ \text{fncs of } x}} \right] \\
&= 2 \underbrace{(\mathbb{E}[\mathbb{E}[y|x]] - \mathbb{E}[y])}_{=0} (\mathbb{E}[\mathbb{E}[y|x]] - f(x)) \\
&= 0.
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left[\mathbb{E} \left[(y - \mathbb{E}[y|x])^2 \mid x \right] \right] + \mathbb{E} \left[(\mathbb{E}[y|x] - f(x))^2 \right] \\
&= \underbrace{\mathbb{E} \left[\text{Var}(y|x) \right]}_{\text{irreducible}} + \underbrace{\mathbb{E} \left[(\mathbb{E}[y|x] - f(x))^2 \right]}_{\text{Bias}^2} \\
&\quad \text{minimized at } f(x) = \mathbb{E}[y|x]
\end{aligned}$$

* Parametric Models

- Need to restrict \mathcal{F}
 e.g. 2-layer NNs, linear funcs, etc.

- Assume $\mathcal{F} = \{ f_\theta : \theta \in \mathcal{U} \}$ where \mathcal{U} is the parameter set.

e.g. $f_\theta(x) = \langle x, \theta \rangle$

$$\ell(y, f_\theta(x)) = (y - \langle x, \theta \rangle)^2 \triangleq \ell(\underbrace{(y, x)}_{\text{data}}, \underbrace{\theta}_{\text{parameter}})$$

- Recall: Want $\theta_* = \operatorname{argmin}_{\theta \in \Theta} R(\theta) = \mathbb{E}[l(y, x), \theta]$

- we can't do this, since we don't know P .

- Can do: Observe data $(y_i, x_i) \stackrel{iid}{\sim} P \quad i=1, \dots, n$

Estimate $R(\theta)$ w/ $\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n l(y_i, x_i), \theta$

Find $\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \hat{R}(\theta)$

Hope $\hat{\theta} \approx \theta_*$

$\hat{R}(\hat{\theta})$: training error

$R(\hat{\theta})$: test error

Def (Excess risk): $R(\hat{\theta}) - R(\theta_*)$

Goal (revised+): Quantify how small excess risk is.

* MLE

- Assume $y|x \sim P_{\theta_*}(y|x)$ where $\theta_* \in \Theta$

↑ unknown true parameter.

but we know the distributional form, that is $P_{\theta}(y|x)$.

- Loss: $l(y, x), \theta = -\log P_{\theta}(y|x)$

Ex: $y|x \sim \mathcal{N}(\langle \theta_*, x \rangle, 1)$ $P_{\theta}(y|x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(y - \langle x, \theta \rangle)^2\right\}$

$l(y, x), \theta = \frac{1}{2}(y - \langle x, \theta \rangle)^2 + \text{constant}$.

- MLE problem: - Observe data $(x_i, y_i) \sim p(x, y)$
 $= p_{\theta_*}(y|x) \cdot p(x)$

- Want θ_*

$$\text{or } \theta_{**} = \underset{\theta}{\operatorname{argmin}} R(\theta)$$

- Can do $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \hat{R}(\theta)$

$$= \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n -\log p_{\theta}(y_i|x_i)$$

- θ_* or θ_{**} ?

$$* \nabla R(\theta_*) = -\mathbb{E} \left[\nabla \log p_{\theta_*}(y|x) \right] = 0 \quad (\text{hw1})$$

$$* \nabla^2 R(\theta_*) = -\mathbb{E} \left[\nabla^2 \log p_{\theta_*}(y|x) \right] = I_{\theta_*} \geq 0$$

- Can assume $I_{\theta_*} > 0 \Rightarrow \theta_*$ is a min of $R(\theta)$.

- Identifiability condition:

$$\theta \neq \theta', p_{\theta} \neq p_{\theta'} \Rightarrow \theta_* \text{ is global minimizer } R(\theta).$$

(exercise)

- Asymptotics of MLE

Def (Conv. in prob): $X_n \xrightarrow{P} X$ if $\forall \epsilon > 0,$

$$\mathbb{P}(|X_n - X| > \epsilon) \rightarrow 0.$$

Def (Conv. in dist.): Let F_n and F be the CDFs of X_n and X , respectively.

$$X_n \xrightarrow{d} X \text{ if } \lim_n F_n(x) = F(x) \text{ for } \forall x \text{ s.t.}$$

F is cont.

Slutsky's thm: Let $X_n \xrightarrow{d} X$, $Y_n \xrightarrow{P} a$, $Z_n \xrightarrow{P} b$,
 then $X_n Y_n + Z_n \xrightarrow{d} aX + b$.

Def (Consistency): $\{\hat{\theta} \equiv \hat{\theta}_n\}_{n \in \mathbb{N}}$ is consistent if $\hat{\theta} \xrightarrow{P} \theta_*$.

Theorem (Asymp of MLE): Assume $\{\hat{\theta}\}_n$ is consistent,
 and $I_{\theta_*} > 0$ and $\|\nabla^3 \log p_{\theta_*}\|_{op} < B, \forall \theta$.

Then,

$$1. \sqrt{n} (\hat{\theta} - \theta_*) \xrightarrow{d} \mathcal{N}(0, I_{\theta_*}^{-1})$$

$$2. n (R(\hat{\theta}) - R(\theta_*)) \xrightarrow{d} \frac{1}{2} \chi_d^2$$

Remarks: 1- $I_{\theta_*} \uparrow \Rightarrow$ less variance \Rightarrow more info on θ_* .

2- Asymptotic efficiency

$$3- \mathbb{E} \chi_d^2 = d \Rightarrow R(\hat{\theta}) - R(\theta_*) \approx O(d/n)$$

generalization error
(excess risk)

dimension
 \uparrow

\downarrow
num of samples.

Tensor notation:

$$A \in \mathbb{R}^{n \times n}, u \in \mathbb{R}^n$$

$$Au \in \mathbb{R}^n \rightarrow A[u]$$

$$u^T A u \in \mathbb{R} \rightarrow A[u, u]$$

$$A \in \mathbb{R}^{n \times n \times n}, u \in \mathbb{R}^n$$

$$A[u] \in \mathbb{R}^{n \times n}$$

$$A[u, u] \in \mathbb{R}^n$$

$$A[u, u, u] \in \mathbb{R}$$

\hookrightarrow tensor

$$\text{Ex: } f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\nabla f \in \mathbb{R}^d$$

$$\nabla^2 f \in \mathbb{R}^{d \times d}$$

$$\nabla^3 f \in \mathbb{R}^{d \times d \times d}$$

* proof: $\nabla \hat{R}(\hat{\theta}) = 0 = \nabla R(\theta_*)$

By Taylor's thm,

1.

$$0 = \nabla \hat{R}(\hat{\theta}) = \nabla \hat{R}(\theta_*) + \nabla^2 \hat{R}(\theta_*) (\hat{\theta} - \theta_*) + \frac{1}{2} \nabla^3 \hat{R}(\bar{\theta}) [\hat{\theta} - \theta_*, \hat{\theta} - \theta_*]$$

$\bar{\theta} \in (\theta_*, \hat{\theta})$
↑

$$\Rightarrow -\nabla \hat{R}(\theta_*) = \left\{ \nabla^2 \hat{R}(\theta_*) + \frac{1}{2} \nabla^3 \hat{R}(\bar{\theta}) [\hat{\theta} - \theta_*] \right\} (\hat{\theta} - \theta_*)$$

\mathbb{R}^d $\mathbb{R}^{d \times d}$ $\mathbb{R}^{d \times d \times d}$ \mathbb{R}^d \mathbb{R}^d

$$\Rightarrow \sqrt{n} \frac{1}{n} \sum_{i=1}^n \nabla \log p_{\theta_*}(y_i | x_i) = \left\{ \frac{1}{n} \sum_{i=1}^n -\nabla^2 \log p_{\theta_*}(y_i | x_i) + \frac{1}{2} \nabla^3 R(\bar{\theta}) [\hat{\theta} - \theta_*] \right\} \sqrt{n} (\hat{\theta} - \theta_*)$$

\downarrow by CLT \downarrow LLN $\leq B$ \downarrow p \downarrow Want

$$\mathcal{N}(0, \text{Cov}(\nabla \log p_{\theta_*}(y_i | x_i)))$$

$= I_{\theta_*}$

$$= \mathbb{E}[-\nabla^2 \log p_{\theta_*}(y_i | x_i)]$$

$= I_{\theta_*}$

by Slutsky's thm

$$\Rightarrow \mathcal{N}(0, I_{\theta_*}) = I_{\theta_*} \cdot \sqrt{n} (\hat{\theta} - \theta_*)$$

$$\Rightarrow \sqrt{n} (\hat{\theta} - \theta_*) \xrightarrow{d} \mathcal{N}(0, I_{\theta_*}^{-1})$$

$$x \sim \mathcal{N}(\mu, \Sigma)$$

$$Ax \sim \mathcal{N}(A\mu, A\Sigma A^T)$$

2. For excess risk:

$$R(\hat{\theta}) = R(\theta_*) + \underbrace{\langle \nabla R(\theta_*), \hat{\theta} - \theta_* \rangle}_{=0} + \frac{1}{2} \langle \hat{\theta} - \theta_*, \nabla^2 R(\theta_*) (\hat{\theta} - \theta_*) \rangle + \frac{1}{6} \nabla^3 R(\bar{\theta}) [\hat{\theta} - \theta_*, \hat{\theta} - \theta_*, \hat{\theta} - \theta_*]$$

$$n (R(\hat{\theta}) - R(\theta_*)) = \frac{1}{2} \left\langle \underbrace{\sqrt{n}(\hat{\theta} - \theta_*)}_{\downarrow N(0, I_{\theta_*}^{-1})}, \left\{ \underbrace{\nabla^2 R(\theta_*)}_{I_{\theta_*}} + \frac{1}{3} \nabla^3 R(\bar{\theta}) [\hat{\theta} - \theta_*] \right\} \underbrace{\sqrt{n}(\hat{\theta} - \theta_*)}_{\downarrow N(0, I_{\theta_*}^{-1})} \right\rangle$$

\leftarrow same n.v. \rightarrow

by Slutsky's thm: $= \frac{1}{2} \left\langle N(0, I_{\theta_*}^{-1}), I_{\theta_*} \cdot N(0, I_{\theta_*}^{-1}) \right\rangle$

$$= \frac{1}{2} \left\langle I_{\theta_*}^{1/2} N(0, I_{\theta_*}^{-1}), I_{\theta_*}^{1/2} N(0, I_{\theta_*}^{-1}) \right\rangle$$

$= I_{\theta_*}^{1/2} I_{\theta_*}^{1/2}$

$$= \frac{1}{2} \left\langle N(0, I), N(0, I) \right\rangle$$

\uparrow same n.v. \uparrow

If $\{z_i\} \stackrel{i.i.d.}{\sim} N(0, 1)$
 $\sum_{i=1}^d z_i^2 \sim \chi_d^2$

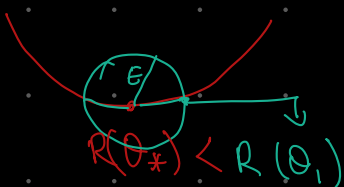
$$\sim \frac{1}{2} \chi_d^2$$

□

Consistency of MLE

Theorem: Assume 1- uniform convergence $\sup_{\theta \in \Theta} |\hat{R}(\theta) - R(\theta)| \xrightarrow{P} 0$

2- identifiability $\inf_{\theta: \|\theta - \theta_*\| \geq \epsilon > 0} R(\theta) > R(\theta_*)$



$\theta: \|\theta - \theta_*\| \geq \epsilon > 0$

\mathcal{U} is compact, non-empty and R is cts.

Then, $\hat{\theta} = \underset{\theta \in \mathcal{U}}{\operatorname{argmin}} \hat{R}(\theta)$ is consistent for θ_* .

proof: - Item 3 $\Rightarrow \hat{\theta} \in \mathcal{U}$ and $\theta_* \in \mathcal{U}$

- $\hat{\theta}$ minimizes $\hat{R}(\theta)$ in \mathcal{U}

$$\Rightarrow \hat{R}(\hat{\theta}) \leq \hat{R}(\theta_*) - R(\theta_*) + R(\theta_*)$$

$$\leq \sup_{\theta \in \mathcal{U}} |\hat{R}(\theta) - R(\theta)| + R(\theta_*) \xrightarrow{P} R(\theta_*) \quad (\text{by item 1})$$

$\xrightarrow{P} 0$

$$\Rightarrow \hat{R}(\hat{\theta}) \approx R(\theta_*) \text{ as } n \rightarrow \infty.$$

$$(*) - 0 \leq R(\hat{\theta}) - R(\theta_*) \leq R(\hat{\theta}) - \hat{R}(\hat{\theta}) \quad (\text{when } n \rightarrow \infty)$$

$$\leq \sup_{\theta \in \mathcal{U}} |\hat{R}(\theta) - R(\theta)| \xrightarrow{P} 0 \quad (\text{by item 1})$$

- Then, compare the events

$$\forall \epsilon, \delta: \left\{ \|\hat{\theta} - \theta_*\| \geq \epsilon \right\} \stackrel{\text{by item 2}}{\subseteq} \left\{ R(\hat{\theta}) - R(\theta_*) > \delta \right\}$$

$$\text{but } \mathbb{P}(\quad) \xrightarrow{P} 0$$

by * above. \square