

2.8 General loss functions and random design

Goal: To quantify $L(\hat{\theta}) - L(\theta^*)$ for a general $L(\theta)$.

For an example $z = (x, y)$, $l(z, \theta)$ is the loss function and $\theta \in \mathbb{R}^d$.

Denote \mathcal{Z} as the set of all examples. Let $p^* = \Delta(\mathcal{Z})$.

Let $\theta^* \in \mathbb{R}^d$ be the minimal of expected risk:

$$\theta^* := \operatorname{argmax}_{\theta \in \mathbb{R}^d} L(\theta), \quad L(\theta) := E_{z \sim p^*}[l(z, \theta)]$$

Let $\hat{\theta} \in \mathbb{R}^d$ be the minimizer of the empirical risk:

$$\hat{\theta} := \operatorname{argmax}_{\theta \in \mathbb{R}^d} \hat{L}(\theta), \quad \hat{L}(\theta) := \frac{1}{n} \sum_{i=1}^n l(z^{(i)}, \theta)$$

where $z^{(i)}$ i.i.d. drawn from p^* .

Assumptions on $l(z, \theta)$:

- (i) $l(z, \theta)$ is twice differentiable in θ .
- (ii) $\nabla l(z, \theta) \in \mathbb{R}^d$ means the gradient at θ .
- (iii) $\nabla^2 l(z, \theta) \in \mathbb{R}^{d \times d}$ means the Hessian at θ .
- (iv) $E_{z \sim p^*}[\nabla^2 l(z, \theta)] > 0$ is positive definite for all θ .

[Definition 3] Well-specified model

- (i) $l(x, y; \theta) := -\log p_{\theta}(y|x)$
- (ii) $\{p_{\theta} : \theta \in \mathbb{R}^d\}$ is conditionally well-specified if $p^*(x, y) = p^*(x) p_{\theta^*}(y|x)$ for some $\theta^* \in \mathbb{R}^d$.
- (iii) Suppose each θ specifies a $p_{\theta}(x, y)$. $\{p_{\theta} : \theta \in \mathbb{R}^d\}$ is jointly well-specified if $p^*(x, y) = p_{\theta^*}(x, y)$ for some $\theta^* \in \mathbb{R}^d$.

[Theorem 3] Bartlett identity.

In the well-specified case (conditionally, thus jointly), the

following holds:

$$\nabla^2 L(\theta^*) = \text{Cov}[\nabla l(z, \theta^*)]$$

$$\begin{aligned} \text{Pf: } 1 &= \int p^*(z) dz = \int p^*(x) p_{\theta^*}(y|x) dz \\ \Rightarrow \int p^*(x) e^{-l(z, \theta^*)} dz &= 1 \end{aligned}$$

differentiate w.r.t. θ^* :

$$\int p^*(x) e^{-l(z, \theta^*)} (-\nabla l(z, \theta^*)) dz = 0$$

which implies $E[\nabla l(z, \theta^*)] = 0$. Differentiate again:

$$0 = \int p^*(x) [-e^{-l(z, \theta^*)} \nabla^2 l(z, \theta^*) + e^{-l(z, \theta^*)} \nabla l(z, \theta^*) \nabla l(z, \theta^*)^T] dz$$

$$= -E[\nabla^2 l(z, \theta^*)] + E[\nabla l(z, \theta^*) \nabla l(z, \theta^*)^T]$$

$$= -E[\nabla^2 l(z, \theta^*)] + \text{Cov}[\nabla l(z, \theta^*)] \quad \text{Since } E[\nabla l(z, \theta^*)] = 0.$$

$$\Rightarrow \nabla^2 L(\theta^*) = E[\nabla^2 l(z, \theta^*)] = \text{Cov}[\nabla l(z, \theta^*)] \quad \square$$

[Example 2] well-specified random design linear regression.

Model:

(i) $x \sim p^*(x)$ for some arbitrary $p^*(x)$

(ii) $y = \theta^* \cdot x + \varepsilon$ where $\varepsilon \sim N(0, 1)$

Loss function: $l(x, y; \theta) := \frac{1}{2} (\theta \cdot x - y)^2$

Property 1: $\nabla^2 L(\theta) = \text{Cov}[\nabla l(z, \theta^*)]$

$$\begin{aligned} \text{(i)} \quad \nabla^2 L(\theta) &= \nabla^2 E_{x \sim p^*, \varepsilon \sim N(0, 1)} \left[\frac{1}{2} (\theta \cdot x - \theta^* \cdot x - \varepsilon)^2 \right] \\ &= E[xx^T] \end{aligned}$$

$$\text{(ii)} \quad \text{Cov}[\nabla l(z, \theta^*)] = \text{Cov}[-\varepsilon x^T]$$

$$= E[\varepsilon x x^T \varepsilon] - E[\varepsilon x] E[\varepsilon x]^T$$

$$= E[xx^T]$$

Since x, ε are independent.

Now we study $\hat{\theta} - \theta^*$:

Step 1: perform a Taylor expansion of $\nabla \hat{L}$ around θ^* :

$$\nabla \hat{L}(\hat{\theta}) = \nabla \hat{L}(\theta^*) + \nabla^2 \hat{L}(\theta^*)(\hat{\theta} - \theta^*) + O_p(\|\hat{\theta} - \theta^*\|_2)$$

Using the fact $\nabla \hat{L}(\hat{\theta}) = 0$ since $\hat{\theta}$ is optimal:

$$\hat{\theta} - \theta^* = -\nabla^2 \hat{L}(\theta^*)^{-1} (\nabla \hat{L}(\theta^*) + O_p(\|\hat{\theta} - \theta^*\|_2)) \quad (I)$$

As $n \rightarrow \infty$, by the weak law of large numbers:

$$\nabla^2 \hat{L}(\theta^*) \xrightarrow{P} \nabla^2 L(\theta^*)$$

Since $\nabla^2 L(\theta^*) > 0$, $\nabla^2 L(\cdot)^{-1}$ is smooth around θ^* :

$$\nabla^2 \hat{L}(\theta^*)^{-1} \xrightarrow{P} \nabla^2 L(\theta^*)^{-1} \quad (CMT)$$

By central limit theorem:

$$\sqrt{n} \nabla \hat{L}(\theta^*) \xrightarrow{d} N(0, \text{Cov}[\nabla l(z, \theta^*)])$$

Suppose $\hat{\theta} - \theta^* = O_p(f(n))$, then by (I), $f(n)$ decays at a rate of $O(\frac{1}{\sqrt{n}})$ or $f^2(n)$, which implies $f^2(n) = \frac{1}{\sqrt{n}}$. Thus:

$$\sqrt{n} \cdot O_p(\|\hat{\theta} - \theta^*\|_2) \xrightarrow{P} 0.$$

By Slutsky's theorem, with (I):

$$\sqrt{n} \cdot (\hat{\theta} - \theta^*) \xrightarrow{d} N(0, \nabla^2 L(\theta^*)^{-1} \text{Cov}[\nabla l(z, \theta^*)] \nabla^2 L(\theta^*)^{-1})$$

Due to Property 1:

$$\sqrt{n} (\hat{\theta} - \theta^*) \xrightarrow{d} N(0, E[xx^T]^{-1})$$