

## §5 Estimation

### §5.1 Estimator

5.1.1  $\mathbf{X}$ : observed dataset

$\theta$ : unknown parameter

$\psi(\theta)$ : target for estimation, for some specified function  $\psi$  of  $\theta$

*Point estimation* of  $\psi(\theta)$ : selection of a “good” value to *estimate*  $\psi(\theta)$  based on  $\mathbf{X}$ .

5.1.2 **Definition.** An *estimator* of  $\psi(\theta)$  is a statistic  $T = T(\mathbf{X})$  for estimating  $\psi(\theta)$ .

5.1.3  $T(\mathbf{X})$  is random (varies from sample to sample) and has its own sampling distribution, based on which the quality of  $T(\mathbf{X})$  is assessed in the frequentist sense.

5.1.4 Treating  $T(\mathbf{X})$  as a “decision rule”, we may assess its quality by its risk function

$$R(\theta, T) = \mathbb{E}_{\theta}[L(\theta, T(\mathbf{X}))],$$

where  $L(\theta, a)$  defines the “loss” resulting from estimating  $\psi(\theta)$  by the “action” (estimate)  $a$ .

### §5.2 Bias and mean squared error

5.2.1 **Definition.** Let  $T$  be an estimator of  $\psi(\theta)$ . Then the *bias* of  $T$  is:

$$\text{bias}_{\theta}(T) = \mathbb{E}_{\theta}[T] - \psi(\theta).$$

If  $T$  has zero bias, it is *unbiased*.

5.2.2 Bias of  $T$  measures its *accuracy*.

For  $T \in \mathbb{R}$ ,  $\text{Var}_{\theta}(T)$  (or *standard deviation* of  $T$ ,  $\sqrt{\text{Var}_{\theta}(T)}$  ) measures its *precision*.

A good estimator should be both *accurate* and *precise*.

5.2.3 If we set the loss function to be  $L(\theta, a) = \|\psi(\theta) - a\|_2^2$ , the risk function  $R(\theta, T)$  reduces to the *mean squared error*.

**Definition.** An estimator  $T = (T_1, \dots, T_d)^\top$  of  $\psi(\theta) \in \mathbb{R}^d$  has the *mean squared error* (MSE)

$$\begin{aligned} \text{MSE}_\theta(T) &= \mathbb{E}_\theta[\|T - \psi(\theta)\|_2^2] \\ &= \mathbb{E}_\theta[\|T - \mathbb{E}_\theta[T]\|_2^2] + 2\mathbb{E}_\theta[(T - \mathbb{E}_\theta[T])^\top (\mathbb{E}_\theta[T] - \psi(\theta))] + \|\mathbb{E}_\theta[T] - \psi(\theta)\|_2^2 \\ &= \sum_{j=1}^d \text{Var}_\theta(T_j) + \sum_{j=1}^d \{\text{bias}_\theta(T_j)\}^2 = \sum_{j=1}^d \text{MSE}_\theta(T_j). \end{aligned}$$

MSE provides a measure of the quality of an estimator by taking into account both accuracy (bias) and precision (variance).

Small MSE  $\Rightarrow$  sampling distribution of  $T$  highly concentrated near  $\psi(\theta)$ .

### §5.3 Rao-Blackwell Theorem

5.3.1 **Lemma.** Let  $T = T(\mathbf{X})$  be complete sufficient for  $\theta$ . Then, for any function  $\psi(\theta)$  of  $\theta$ , there exists at most one unique function of  $T$  which is unbiased for  $\psi(\theta)$ .

.....

*Proof:*

Suppose  $g(T)$ ,  $h(T)$  are unbiased estimators of  $\psi(\theta)$  such that

$$\mathbb{E}_\theta[g(T)] = \mathbb{E}_\theta[h(T)] = \psi(\theta) \quad \forall \theta.$$

Then  $\mathbb{E}_\theta[g(T) - h(T)] = 0 \quad \forall \theta \implies \mathbb{P}_\theta(g(T) = h(T)) = 1.$

■

5.3.2  $\mathbf{X} \sim f(\cdot|\theta)$

Target for estimation:  $\psi(\theta)$

Loss function:  $L(\theta, a)$  [ = loss incurred when  $\psi(\theta)$  is estimated by  $a$  ]

**Theorem.** (*Rao-Blackwell*)

Suppose that  $T(\mathbf{X})$  is sufficient for  $\theta$  and  $L(\theta, a)$  is convex in  $a$ .

Let  $\rho(\mathbf{X})$  be an estimator of  $\psi(\theta)$ . Define

$$\rho^*(t) = \mathbb{E}[\rho(\mathbf{X}) | T(\mathbf{X}) = t].$$

Then, if  $\rho(\mathbf{X})$  has finite expectation and risk, i.e. both  $\mathbb{E}_\theta[\rho(\mathbf{X})]$  and  $\mathbb{E}_\theta[L(\theta, \rho(\mathbf{X}))]$  are finite, then  $\rho^*(T(\mathbf{X}))$  is an estimator of  $\psi(\theta)$  such that

$$\mathbb{E}_\theta[L(\theta, \rho^*(T(\mathbf{X})))] \leq \mathbb{E}_\theta[L(\theta, \rho(\mathbf{X}))].$$

Moreover, if  $\rho(\mathbf{X})$  is unbiased and  $T(\mathbf{X})$  is complete, then

- (i)  $\rho^*(T(\mathbf{X}))$  is the **unique** unbiased estimator of  $\psi(\theta)$  which is a function of  $T(\mathbf{X})$ ;
- (ii)  $\mathbb{E}_\theta[L(\theta, \rho^*(T(\mathbf{X})))] \leq \mathbb{E}_\theta[L(\theta, S(\mathbf{X}))]$  for any unbiased estimator  $S(\mathbf{X})$  of  $\psi(\theta)$ .  
[i.e.  $\rho^*(T(\mathbf{X}))$  either strictly dominates, or has at least the same risk function as, any unbiased estimator  $S(\mathbf{X})$ .]
- .....

*Proof:*

*Sufficiency of  $T(\mathbf{X})$  for  $\theta$  implies that the conditional distribution of  $\rho(\mathbf{X})$  given  $T(\mathbf{X})$  is free of  $\theta$ , so  $\rho^*(T(\mathbf{X}))$  is a **legitimate** estimator of  $\psi(\theta)$ .*

*Applying Jensen's inequality,*

$$L(\theta, \rho^*(T(\mathbf{X}))) \leq \mathbb{E}[L(\theta, \rho(\mathbf{X}))|T(\mathbf{X})] \quad \Rightarrow \quad \mathbb{E}_\theta[L(\theta, \rho^*(T(\mathbf{X})))] \leq \mathbb{E}_\theta[L(\theta, \rho(\mathbf{X}))].$$

*Next suppose  $\rho(\mathbf{X})$  is unbiased and  $T(\mathbf{X})$  is complete. Then*

$$\mathbb{E}_\theta[\rho^*(T(\mathbf{X}))] = \mathbb{E}_\theta[\mathbb{E}[\rho(\mathbf{X})|T(\mathbf{X})]] = \mathbb{E}_\theta[\rho(\mathbf{X})] = \psi(\theta),$$

*i.e.  $\rho^*(T(\mathbf{X}))$  is also unbiased for  $\psi(\theta)$ .*

*Using Lemma §5.3.1,  $\rho^*(T(\mathbf{X}))$  is the unique function of  $T(\mathbf{X})$  which is unbiased for  $\psi(\theta)$ . This proves part (i).*

*To prove (ii), for any unbiased estimator  $S(\mathbf{X})$ , define  $S^*(T(\mathbf{X})) = \mathbb{E}[S(\mathbf{X})|T(\mathbf{X})]$ . Then  $S^*(T(\mathbf{X}))$  is unbiased for  $\psi(\theta)$  and is a function of  $T(\mathbf{X})$ . By part (i), we have  $S^* = \rho^*$  with probability one for all  $\theta$ . Thus*

$$\mathbb{E}_\theta[L(\theta, \rho^*(T(\mathbf{X})))] = \mathbb{E}_\theta[L(\theta, S^*(T(\mathbf{X})))] \leq \mathbb{E}_\theta[L(\theta, S(\mathbf{X}))]. \quad \blacksquare$$

### 5.3.3 Important applications of Rao-Blackwell Theorem —

- We may reduce, or at least not increase, the risk of an arbitrary estimator by evaluating its expected value conditional on a sufficient statistic.
- If a complete sufficient statistic  $T$  exists, we can obtain an unbiased estimator which has the smallest risk among all unbiased estimators. This can be done by  
*either*

finding the expected value of any unbiased estimator conditional on  $T$ ,

*or*

finding an unbiased estimator which is a function of  $T$ .

5.3.4 Consider a scalar target  $\psi(\theta) \in \mathbb{R}$ . Setting  $L(\theta, a) = \{\psi(\theta) - a\}^2$  (squared loss), we obtain immediately the following.

**Corollary.** Let  $\rho(\mathbf{X})$  be an estimator of  $\psi(\theta)$  with finite second moment (i.e.  $\mathbb{E}_\theta[\rho(\mathbf{X})^2] < \infty$ ). Then

$$\text{MSE of } \rho^*(T(\mathbf{X})) \leq \text{MSE of } \rho(\mathbf{X}).$$

If we assume further that  $\rho(\mathbf{X})$  is unbiased and  $T(\mathbf{X})$  is complete, then

- (i)  $\rho^*(T(\mathbf{X}))$  is the unique unbiased estimator of  $\psi(\theta)$  which is a function of  $T(\mathbf{X})$ ;
- (ii)  $\text{Var}_\theta(\rho^*(T(\mathbf{X}))) \leq \text{Var}_\theta(S(\mathbf{X}))$  for any unbiased estimator  $S(\mathbf{X})$  of  $\psi(\theta)$ , i.e.

$\rho^*(T(\mathbf{X}))$  is a *uniformly minimum variance unbiased* (UMVU) estimator, and is also the unique UMVU estimator which is a function of  $T(\mathbf{X})$ .

### 5.3.5 Examples.

- (i) Example §4.5.2(1) has proved that  $\max_i X_i$  is complete sufficient for  $\theta$  for  $X_i \text{ iid } \sim U[0, \theta]$ . Since  $\mathbb{E}_\theta[2X_1] = 2(\theta/2) = \theta$ ,  $2X_1$  is unbiased for  $\theta$ .

Note that

$$\mathbb{P}_\theta(X_1 = \max_i X_i) = 1/n, \quad \mathbb{P}_\theta(X_1 < \max_i X_i) = 1 - 1/n.$$

Given  $\max_i X_i = t$  and  $X_1 < \max_i X_i$ ,  $X_1 \sim U[0, t]$ . Thus

$$\mathbb{E}[2X_1 | \max_i X_i = t] = (2t)(1/n) + t(1 - 1/n) = (1 + 1/n)t,$$

so that

$(1 + 1/n)\max_i X_i$  has the smallest risk among all unbiased estimators of  $\theta$ , for any convex loss function. In particular, it is a UMVU estimator of  $\theta$ .

- (ii)  $\mathbf{X} = (X_1, \dots, X_n) \text{ iid } \sim N(\mu, \sigma^2) \leftarrow$  here  $\theta = (\mu, \sigma)$

By exponential family properties,  $T = (\sum_{i=1}^n X_i^2, \sum_{i=1}^n X_i)$  is complete sufficient for  $\theta = (\mu, \sigma)$ .

- The sample mean  $\bar{X}$  is unbiased for  $\psi(\theta) = \mu$  and is a function of  $T$ . Thus

$\bar{X}$  has the smallest risk among all unbiased estimators of  $\mu$ , for any convex loss function. In particular, it is a UMVU estimator of  $\mu$ .

- The sample variance  $S^2 = \sum_{i=1}^n (X_i - \bar{X})^2 / (n - 1)$  is unbiased for  $\psi(\theta) = \sigma^2$  and is a function of  $T$ . Thus

$S^2$  has the smallest risk among all unbiased estimators of  $\sigma^2$ , for any convex loss function. In particular, it is a UMVU estimator of  $\sigma^2$ .

## §5.4 Information inequality

### 5.4.1 Regularity assumptions.

- $\mathbf{X} \sim f(\cdot|\theta)$ ,  $\theta = (\theta_1, \dots, \theta_k)^\top \in \Theta \subset \mathbb{R}^k$ , where  $\Theta$  contains an **open rectangle**,
- sample space  $\{\mathbf{x} : f(\mathbf{x}|\theta) > 0\} = \mathcal{S}$  is common to all  $\theta \in \Theta$ ,
- for any  $\mathbf{x} \in \mathcal{S}$ ,  $\theta \in \Theta$  and  $i = 1, \dots, k$ ,  $\frac{\partial f(\mathbf{x}|\theta)}{\partial \theta_i}$  exists and is finite.

5.4.2 **Definitions.** Under regularity assumptions, the *score function* is defined to be

$$\mathbf{U}(\theta) = [U_1(\theta), \dots, U_k(\theta)]^\top = \left[ \frac{\partial \ln f(\mathbf{X}|\theta)}{\partial \theta_1}, \dots, \frac{\partial \ln f(\mathbf{X}|\theta)}{\partial \theta_k} \right]^\top,$$

and the (*Fisher*) *information matrix* is the  $k \times k$  matrix

$$I(\theta) = \mathbb{E}_\theta[\mathbf{U}(\theta)\mathbf{U}(\theta)^\top],$$

whose  $(i, j)$ th entry is given by

$$I_{ij}(\theta) = \mathbb{E}_\theta[U_i(\theta)U_j(\theta)] = \mathbb{E}_\theta \left[ \left( \frac{\partial \ln f(\mathbf{X}|\theta)}{\partial \theta_i} \right) \left( \frac{\partial \ln f(\mathbf{X}|\theta)}{\partial \theta_j} \right) \right].$$

Note: Equivalently, we may define

$$\mathbf{U}(\theta) = \left[ \frac{\partial \ln \ell_{\mathbf{X}}(\theta)}{\partial \theta_1}, \dots, \frac{\partial \ln \ell_{\mathbf{X}}(\theta)}{\partial \theta_k} \right]^\top = \left[ \frac{\partial S_{\mathbf{X}}(\theta)}{\partial \theta_1}, \dots, \frac{\partial S_{\mathbf{X}}(\theta)}{\partial \theta_k} \right]^\top,$$

where  $\ell_{\mathbf{X}}(\theta) \propto f(\mathbf{X}|\theta)$  denotes the likelihood function and  $S_{\mathbf{X}}(\theta)$  the loglikelihood function.

5.4.3 **Lemma.** Under regularity assumptions and assuming that  $\frac{\partial}{\partial \theta_i} \int f(\mathbf{x}|\theta) d\mathbf{x} = \int \frac{\partial f(\mathbf{x}|\theta)}{\partial \theta_i} d\mathbf{x}$ , we have

$$\mathbb{E}_\theta[\mathbf{U}(\theta)] = \mathbf{0} \quad \text{and} \quad I(\theta) = \text{Var}_\theta(\mathbf{U}(\theta)).$$

If, in addition,  $f(\mathbf{x}|\theta)$  has second derivatives  $\frac{\partial^2}{\partial \theta_i \partial \theta_j} f(\mathbf{x}|\theta)$  for all  $i, j$ , then

$$I(\theta) = -\mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln f(\mathbf{X}|\theta) \right].$$

Note: the lemma says that componentwise,

$$I_{ij}(\theta) = \text{Cov}_\theta(U_i(\theta), U_j(\theta)) = -\mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln f(\mathbf{X}|\theta) \right].$$

.....

*Proof:*

$$\mathbb{E}_\theta[U_i(\theta)] = \mathbb{E}_\theta \left[ \frac{1}{f(\mathbf{X}|\theta)} \frac{\partial f(\mathbf{X}|\theta)}{\partial \theta_i} \right] = \int \frac{\partial f(\mathbf{x}|\theta)}{\partial \theta_i} d\mathbf{x} = \frac{\partial}{\partial \theta_i} \int f(\mathbf{x}|\theta) d\mathbf{x} = \frac{\partial}{\partial \theta_i}(1) = 0,$$

which implies immediately  $I(\theta) = \mathbb{E}_\theta[\mathbf{U}(\theta)\mathbf{U}(\theta)^\top] = \text{Var}_\theta(\mathbf{U}(\theta))$ .

To prove the last assertion, consider

$$\begin{aligned} \mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln f(\mathbf{X}|\theta) \right] &= -\mathbb{E}_\theta \left[ \frac{1}{f(\mathbf{X}|\theta)^2} \frac{\partial f(\mathbf{X}|\theta)}{\partial \theta_i} \frac{\partial f(\mathbf{X}|\theta)}{\partial \theta_j} \right] + \mathbb{E}_\theta \left[ \frac{1}{f(\mathbf{X}|\theta)} \frac{\partial^2 f(\mathbf{X}|\theta)}{\partial \theta_i \partial \theta_j} \right] \\ &= -\mathbb{E}_\theta[U_i(\theta)U_j(\theta)] + \int \frac{\partial^2 f(\mathbf{x}|\theta)}{\partial \theta_i \partial \theta_j} d\mathbf{x} \\ &= -I_{ij}(\theta) + \frac{\partial^2}{\partial \theta_i \partial \theta_j} \int f(\mathbf{x}|\theta) d\mathbf{x} = -I_{ij}(\theta). \end{aligned} \quad \blacksquare$$

#### 5.4.4 **Theorem.** (Information Inequality)

Let  $T = T(\mathbf{X})$  be a statistic with  $\mathbb{E}_\theta[T^2] < \infty$ . Assume that  $\boldsymbol{\alpha}(\theta) = \left[ \frac{\partial}{\partial \theta_1} \mathbb{E}_\theta[T], \dots, \frac{\partial}{\partial \theta_k} \mathbb{E}_\theta[T] \right]^\top$  exists, and can be obtained by differentiating under the integral sign.

Then, under regularity assumptions and assuming that  $I(\theta)$  is positive definite,

$$\text{Var}_\theta(T) \geq \boldsymbol{\alpha}(\theta)^\top I(\theta)^{-1} \boldsymbol{\alpha}(\theta).$$

.....

*Proof:*

Note that

$$\mathbb{E}_\theta[TU_i(\theta)] = \int T(\mathbf{x}) \left\{ \frac{1}{f(\mathbf{x}|\theta)} \frac{\partial f(\mathbf{x}|\theta)}{\partial \theta_i} \right\} f(\mathbf{x}|\theta) d\mathbf{x} = \frac{\partial}{\partial \theta_i} \int T(\mathbf{x}) f(\mathbf{x}|\theta) d\mathbf{x} = \frac{\partial}{\partial \theta_i} \mathbb{E}_\theta[T],$$

so that  $\boldsymbol{\alpha}(\theta) = \mathbb{E}_\theta[\mathbf{U}(\theta)T] = \mathbb{E}_\theta[\mathbf{U}(\theta)(T - \mathbb{E}_\theta[T])]$ .

Consider

$$\begin{aligned} 0 &\leq \mathbb{E}_\theta \left[ \{T - \mathbb{E}_\theta[T] - \boldsymbol{\alpha}(\theta)^\top I(\theta)^{-1} \mathbf{U}(\theta)\}^2 \right] \\ &= \text{Var}_\theta(T) - 2\boldsymbol{\alpha}(\theta)^\top I(\theta)^{-1} \mathbb{E}_\theta[\mathbf{U}(\theta)(T - \mathbb{E}_\theta[T])] + \boldsymbol{\alpha}(\theta)^\top I(\theta)^{-1} \mathbb{E}_\theta[\mathbf{U}(\theta)\mathbf{U}(\theta)^\top] I(\theta)^{-1} \boldsymbol{\alpha}(\theta) \\ &= \text{Var}_\theta(T) - \boldsymbol{\alpha}(\theta)^\top I(\theta)^{-1} \boldsymbol{\alpha}(\theta), \end{aligned}$$

which implies the Information Inequality.  $\blacksquare$

5.4.5 **Corollary.** Let  $\psi(\theta)$  be a differentiable function of  $\theta$ , and  $T$  be an unbiased estimator of  $\psi(\theta)$ . Then

$$\text{Var}_\theta(T) \geq \left[ \frac{\partial \psi(\theta)}{\partial \theta_1}, \dots, \frac{\partial \psi(\theta)}{\partial \theta_k} \right] I(\theta)^{-1} \left[ \frac{\partial \psi(\theta)}{\partial \theta_1}, \dots, \frac{\partial \psi(\theta)}{\partial \theta_k} \right]^\top.$$

.....  
*Proof: Immediate from Theorem §5.4.4, noting that  $\mathbb{E}_\theta[T] = \psi(\theta)$ .* ■

Practical implication:

*The above corollary provides a lower bound for the variance of **any** unbiased estimator of  $\psi(\theta)$ . If an unbiased estimator has variance **equal** to this lower bound, it must be a UMVU estimator of  $\psi(\theta)$ .*

5.4.6 For the special case where  $k = 1$  and  $\psi(\theta) = \theta$ , the lower bound given in Corollary §5.4.5, i.e.  $I(\theta)^{-1} = 1/I(\theta)$ , is known as the *Cramér-Rao Lower Bound* (CRLB).

5.4.7 **Example §5.4.1** (*Normal mean regression*)

$(u_1, \dots, u_n)$ :  $u_i$  a non-random covariate vector associated with  $i$ th response

$\mu_\beta(u)$ : (possibly non-linear) regression function parameterised by  $\beta \in \mathbb{R}^q$

$\mathbf{X} = (X_1, \dots, X_n)$ : independent responses, with  $X_i \sim N(\mu_\beta(u_i), \nu)$

Joint pdf:  $f(\mathbf{x}|\beta, \nu) = (2\pi\nu)^{-n/2} \exp \left[ -\frac{\sum_{i=1}^n \{x_i - \mu_\beta(u_i)\}^2}{2\nu} \right]$ . Then

$$\frac{\partial}{\partial \beta} \ln f(\mathbf{X}|\beta, \nu) = \nu^{-1} \sum_{i=1}^n \{X_i - \mu_\beta(u_i)\} \frac{\partial \mu_\beta(u_i)}{\partial \beta},$$

$$\frac{\partial}{\partial \nu} \ln f(\mathbf{X}|\beta, \nu) = -\frac{n}{2\nu} + \frac{1}{2\nu^2} \sum_{i=1}^n \{X_i - \mu_\beta(u_i)\}^2,$$

$$\frac{\partial^2}{\partial \beta \partial \beta^\top} \ln f(\mathbf{X}|\beta, \nu) = \nu^{-1} \sum_{i=1}^n \left[ -\frac{\partial \mu_\beta(u_i)}{\partial \beta} \frac{\partial \mu_\beta(u_i)}{\partial \beta^\top} + \{X_i - \mu_\beta(u_i)\} \frac{\partial^2 \mu_\beta(u_i)}{\partial \beta \partial \beta^\top} \right],$$

$$\frac{\partial^2}{\partial \beta \partial \nu} \ln f(\mathbf{X}|\beta, \nu) = -\frac{1}{\nu^2} \sum_{i=1}^n \{X_i - \mu_\beta(u_i)\} \frac{\partial \mu_\beta(u_i)}{\partial \beta},$$

$$\frac{\partial^2}{\partial \nu^2} \ln f(\mathbf{X}|\beta, \nu) = \frac{n}{2\nu^2} - \frac{1}{\nu^3} \sum_{i=1}^n \{X_i - \mu_\beta(u_i)\}^2,$$

so that

$$I(\beta, \nu) = \begin{bmatrix} \nu^{-1} \sum_{i=1}^n \frac{\partial \mu_\beta(u_i)}{\partial \beta} \frac{\partial \mu_\beta(u_i)}{\partial \beta^\top} & \mathbf{0} \\ \mathbf{0}^\top & n/(2\nu^2) \end{bmatrix}$$

and

$$I(\beta, \nu)^{-1} = \begin{bmatrix} \nu \left\{ \sum_{i=1}^n \frac{\partial \mu_\beta(u_i)}{\partial \beta} \frac{\partial \mu_\beta(u_i)}{\partial \beta^\top} \right\}^{-1} & \mathbf{0} \\ \mathbf{0}^\top & 2\nu^2/n \end{bmatrix}.$$

Let  $T$  be any unbiased estimator of  $\mathbf{c}^\top \beta$  for a given vector  $\mathbf{c} \in \mathbb{R}^q$ , i.e.  $\mathbb{E}_{\beta, \nu}[T] = \mathbf{c}^\top \beta$  for all  $\beta, \nu$ . Note

$$\boldsymbol{\alpha}(\beta, \nu) = \begin{bmatrix} \frac{\partial \mathbf{c}^\top \beta}{\partial \beta} \\ \frac{\partial \mathbf{c}^\top \beta}{\partial \nu} \end{bmatrix} = \begin{bmatrix} \mathbf{c} \\ 0 \end{bmatrix}, \quad \boldsymbol{\alpha}(\beta, \nu)^\top I(\beta, \nu)^{-1} \boldsymbol{\alpha}(\beta, \nu) = \nu \mathbf{c}^\top \left\{ \sum_{i=1}^n \frac{\partial \mu_\beta(u_i)}{\partial \beta} \frac{\partial \mu_\beta(u_i)}{\partial \beta^\top} \right\}^{-1} \mathbf{c}.$$

Thus

$$\text{Var}_{\beta, \nu}(T) \geq \nu \mathbf{c}^\top \left\{ \sum_{i=1}^n \frac{\partial \mu_\beta(u_i)}{\partial \beta} \frac{\partial \mu_\beta(u_i)}{\partial \beta^\top} \right\}^{-1} \mathbf{c}.$$

Special cases:

(i) (*Multiple linear regression*)  $\mu_\beta(u) = u^\top \beta$

$$\rightarrow \frac{\partial \mu_\beta(u_i)}{\partial \beta} = u_i \quad \text{and} \quad \text{Var}_{\beta, \nu}(T) \geq \nu \mathbf{c}^\top \left( \sum_{i=1}^n u_i u_i^\top \right)^{-1} \mathbf{c}.$$

Take  $T = \mathbf{c}^\top \left( \sum_{i=1}^n u_i u_i^\top \right)^{-1} \sum_{i=1}^n X_i u_i$  (*least squares estimator*). Then

$$\mathbb{E}_{\beta, \nu}[T] = \mathbf{c}^\top \beta \quad \text{and} \quad \text{Var}_{\beta, \nu}(T) = \nu \mathbf{c}^\top \left( \sum_{i=1}^n u_i u_i^\top \right)^{-1} \mathbf{c} \leftarrow \text{lower bound},$$

i.e.  $T$  gives the **minimum** variance among all unbiased estimators of  $\mathbf{c}^\top \beta$ ,

i.e.  $T$  is a UMVU estimator of  $\mathbf{c}^\top \beta$ .

(ii) (*Two-sample comparison*) — a special case of (i)

$$\begin{cases} X_1, \dots, X_m \text{ iid from } N(\beta_1, \nu) & \rightarrow u_1 = \dots = u_m = [1, 0]^\top, \\ X_{m+1}, \dots, X_n \text{ iid from } N(\beta_2, \nu) & \rightarrow u_{m+1} = \dots = u_n = [0, 1]^\top. \end{cases}$$



Target for estimation:  $\mathbf{c}^\top \beta = [1, -1]\beta = \beta_1 - \beta_2$ .

Take  $T = \mathbf{c}^\top \left( \sum_{i=1}^n u_i u_i^\top \right)^{-1} \sum_{i=1}^n X_i u_i = m^{-1} \sum_{i=1}^m X_i - (n-m)^{-1} \sum_{i=m+1}^n X_i$ . Then

$$\mathbb{E}_{\beta, \nu}[T] = \beta_1 - \beta_2 \quad \text{and} \quad \text{Var}_{\beta, \nu}(T) = \nu \mathbf{c}^\top \left( \sum_{i=1}^n u_i u_i^\top \right)^{-1} \mathbf{c} = \nu \left( \frac{1}{m} + \frac{1}{n-m} \right) \leftarrow \text{lower bound,}$$

i.e.  $T$  gives the **minimum** variance among all unbiased estimators of  $\beta_1 - \beta_2$ ,

i.e.  $T$  is a UMVU estimator of  $\beta_1 - \beta_2$ .

## §5.5 Maximum likelihood estimator

5.5.1 **Definition.** Suppose  $\hat{\theta}$  maximises  $\ell_{\mathbf{X}}(\theta)$ , or equivalently,  $S_{\mathbf{X}}(\theta)$ . We say  $\hat{\theta}$  is the *maximum likelihood estimator* (mle) of  $\theta$ .

5.5.2 Suppose  $\theta \in \mathbb{R}^k$ . Usually  $\hat{\theta}$  can be obtained by solving

$$(*) \text{ likelihood equations : } \mathbf{U}(\theta) = \mathbf{0}, \quad \text{i.e.} \quad \frac{\partial}{\partial \theta_j} S_{\mathbf{X}}(\theta) = 0, \quad j = 1, \dots, k.$$

- $(*)$  may have  $> 1$  solutions for  $\theta \in \Theta$ , so we must check for maximality.
- $(*)$  may be nonsense, e.g.  $\mathbf{X} = (X_1, \dots, X_n)$  iid from Binomial  $(\theta, 1/2)$  and  $\theta$  is an **integer**.
- for some  $\Theta$ , e.g.  $\Theta = [0, 1]$ ,  $S_{\mathbf{X}}(\theta)$  may be maximised at the boundary of  $\Theta$  rather than at a solution to  $(*)$ .

5.5.3 Suppose  $T$  is minimal sufficient for  $\theta$ . Then  $\ell_{\mathbf{X}}(\theta) \propto g(T, \theta)$ . Thus  $\hat{\theta}$  maximises  $g(T, \theta)$  with respect to  $\theta$  and is a function of  $T$  only.

[Note:  $\hat{\theta}$  must be a function of minimal sufficient statistic but itself needs **not** be sufficient for  $\theta$ .]

5.5.4 Maximum likelihood estimators can be biased, e.g.  $X_1, \dots, X_n$  iid from  $N(\mu, \sigma^2)$

$$\rightarrow \text{mle of } \sigma^2 \text{ is } \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \Rightarrow \mathbb{E}[\hat{\sigma}^2] = \frac{n-1}{n} \sigma^2 \neq \sigma^2.$$

5.5.5 If  $\hat{\theta}$  is mle of  $\theta$ , then  $\psi(\hat{\theta})$  is mle of  $\psi(\theta)$  for any transformation  $\psi(\cdot)$ .

5.5.6  $\mathbf{X} = (X_1, \dots, X_n) \sim$  joint probability function  $p_1(x_1|\theta) \times \dots \times p_n(x_n|\theta)$

Likelihood function:  $\ell_n(\theta) \propto \prod_{i=1}^n p_i(X_i|\theta)$ , loglikelihood function:  $S_n(\theta) = \ln \ell_n(\theta)$

Note:  $X_1, \dots, X_n$  are independent but may not be identically distributed.

Let  $\theta_0$  be true value of  $\theta \in \Theta \subset \mathbb{R}^k$ . Assume that  $S_n(\theta)$  is differentiable in an open neighbourhood of  $\theta_0$  and that the likelihood equations have a root at  $\theta = \hat{\theta}_n$  (mle).

**Theorem.** Subject to regularity conditions<sup>1</sup> on  $p_1(\cdot|\theta), \dots, p_n(\cdot|\theta)$ , we have

- (i)  $\hat{\theta}_n$  converges in probability to  $\theta_0$ ,
- (ii)  $n^{1/2}(\hat{\theta}_n - \theta_0)$  converges in distribution to  $N(\mathbf{0}, \mathcal{J}(\theta_0)^{-1})$ ,
- (iii)  $n^{-1/2}\mathbf{U}(\theta_0)$  converges in distribution to  $N(\mathbf{0}, \mathcal{J}(\theta_0))$ ,

where  $\mathcal{J}(\theta) = \lim_{n \rightarrow \infty} n^{-1}I(\theta) = -\lim_{n \rightarrow \infty} n^{-1} \sum_{\ell=1}^n \mathbb{E}_{\theta} \left[ \frac{\partial^2 \ln p_{\ell}(X_{\ell}|\theta)}{\partial \theta \partial \theta^{\top}} \right]$  has  $(i, j)$ th entry

$$\mathcal{J}_{ij}(\theta) = -\lim_{n \rightarrow \infty} n^{-1} \sum_{\ell=1}^n \mathbb{E}_{\theta} \left[ \frac{\partial^2 \ln p_{\ell}(X_{\ell}|\theta)}{\partial \theta_i \partial \theta_j} \right].$$

Note: If  $X_1, \dots, X_n$  are iid from  $p(\cdot|\theta)$ , then  $\mathcal{J}(\theta) = -\mathbb{E}_{\theta} \left[ \frac{\partial^2 \ln p(X_1|\theta)}{\partial \theta \partial \theta^{\top}} \right]$ .

.....

*Proof: (outline)*

- (i) Without loss of generality write  $S_n(\theta) = \sum_{i=1}^n \ln p_i(X_i|\theta) = \sum_{i=1}^n s_i(\theta)$ . For any  $\theta \in \Theta$ , Jensen's inequality and concavity of the logarithm imply that, for  $i = 1, \dots, n$ ,

$$\mathbb{E}_{\theta_0} [s_i(\theta) - s_i(\theta_0)] \leq \ln \mathbb{E}_{\theta_0} \left[ \frac{p_i(X_i|\theta)}{p_i(X_i|\theta_0)} \right] = 0.$$

Thus the function  $\theta \mapsto \lim_{n \rightarrow \infty} n^{-1} \mathbb{E}_{\theta_0} [S_n(\theta)]$  is maximised at  $\theta = \theta_0$ . We may assume under regularity conditions that  $\lim_{n \rightarrow \infty} n^{-1} \mathbb{E}_{\theta_0} [S_n(\theta)] < \lim_{n \rightarrow \infty} n^{-1} \mathbb{E}_{\theta_0} [S_n(\theta_0)]$  for  $0 < \|\theta - \theta_0\| < \Delta$ , some  $\Delta > 0$ .

Fix an arbitrary  $\epsilon \in (0, \Delta)$ . Then, by the (uniform) Weak Law of Large Numbers,

$$n^{-1} \sum_{i=1}^n \{s_i(\theta) - s_i(\theta_0)\} \xrightarrow{P} \lim_{n \rightarrow \infty} n^{-1} \mathbb{E}_{\theta_0} [S_n(\theta) - S_n(\theta_0)] < 0 \quad \text{uniformly over } \|\theta - \theta_0\| = \epsilon.$$

---

<sup>1</sup>For details see Bradley, R.A. and Gart, J.J. (1962). The asymptotic properties of ML estimators when sampling from associated populations. *Biometrika*, **49**, 205–214.

Put  $\delta_\epsilon = -\sup \left\{ \lim_{n \rightarrow \infty} n^{-1} \mathbb{E}_{\theta_0} [S_n(\theta) - S_n(\theta_0)] : \|\theta - \theta_0\| = \epsilon \right\} > 0$ . Assume that  $\hat{\theta}_n$  solves the likelihood equations uniquely. Then

$$\begin{aligned}
& \mathbb{P}_{\theta_0}(\|\hat{\theta}_n - \theta_0\| \leq \epsilon) \\
& \geq \mathbb{P}_{\theta_0}(S_n(\theta) < S_n(\theta_0) \forall \|\theta - \theta_0\| = \epsilon) \\
& \geq \mathbb{P}_{\theta_0}\left(n^{-1} \sum_{i=1}^n \{s_i(\theta) - s_i(\theta_0)\} \leq -\delta_\epsilon/2 \quad \forall \|\theta - \theta_0\| = \epsilon\right) \\
& \geq \mathbb{P}_{\theta_0}\left(n^{-1} \sum_{i=1}^n \{s_i(\theta) - s_i(\theta_0)\} - \lim_{n \rightarrow \infty} n^{-1} \mathbb{E}_{\theta_0} [S_n(\theta) - S_n(\theta_0)] \leq \delta_\epsilon/2 \quad \forall \|\theta - \theta_0\| = \epsilon\right) \rightarrow 1.
\end{aligned}$$

(ii, iii) Write  $W(\theta)$  for the  $k \times k$  matrix  $\frac{\partial^2}{\partial \theta \partial \theta^\top} S_n(\theta)$ , with  $(i, j)$ th entry  $W_{ij}(\theta) = \frac{\partial^2}{\partial \theta_i \partial \theta_j} S_n(\theta)$ .

By the Strong Law of Large Numbers,

$$-n^{-1}W_{ij}(\theta_0) \longrightarrow -\lim_{n \rightarrow \infty} n^{-1} \sum_{\ell=1}^n \mathbb{E}_{\theta_0} \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} s_\ell(\theta) \Big|_{\theta_0} \right] = \mathcal{J}_{ij}(\theta_0) \text{ almost surely.}$$

By Lemma §5.4.3, we have, for  $\ell, \ell' = 1, \dots, n$ ,

$$\begin{aligned}
& \mathbb{E}_{\theta_0} \left[ \frac{\partial}{\partial \theta_i} s_\ell(\theta) \Big|_{\theta_0} \right] = 0 \\
& \text{Cov}_{\theta_0} \left( \frac{\partial}{\partial \theta_i} s_\ell(\theta) \Big|_{\theta_0}, \frac{\partial}{\partial \theta_j} s_{\ell'}(\theta) \Big|_{\theta_0} \right) = -\mathbb{E}_{\theta_0} \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} s_\ell(\theta) \Big|_{\theta_0} \right] \mathbf{1}\{\ell = \ell'\}.
\end{aligned}$$

Then, by Central Limit Theorem,

$$n^{-1/2} \mathbf{U}(\theta_0) = n^{-1/2} \sum_{\ell=1}^n \left[ \frac{\partial}{\partial \theta_1} s_\ell(\theta) \Big|_{\theta_0}, \dots, \frac{\partial}{\partial \theta_k} s_\ell(\theta) \Big|_{\theta_0} \right]^\top \longrightarrow N(\mathbf{0}, \mathcal{J}(\theta_0)) \text{ in distribution,}$$

which proves (iii). Taylor expansion gives

$$\mathbf{0} = n^{-1/2} \mathbf{U}(\hat{\theta}_n) = n^{-1/2} \mathbf{U}(\theta_0) + n^{-1/2} W(\theta_0)(\hat{\theta}_n - \theta_0) + \boldsymbol{\epsilon}_n,$$

where  $\boldsymbol{\epsilon}_n \rightarrow \mathbf{0}$  in probability, so that (ii) follows by noting that

$$\begin{aligned}
n^{1/2}(\hat{\theta}_n - \theta_0) &= \{-n^{-1}W(\theta_0)\}^{-1} \left( n^{-1/2} \mathbf{U}(\theta_0) + \boldsymbol{\epsilon}_n \right) \\
&\longrightarrow N\left(\mathbf{0}, \mathcal{J}(\theta_0)^{-1} \mathcal{J}(\theta_0) [\mathcal{J}(\theta_0)^{-1}]^\top\right) \equiv N(\mathbf{0}, \mathcal{J}(\theta_0)^{-1}) \text{ in distribution.} \quad \blacksquare
\end{aligned}$$

### 5.5.7 Practical implication:

For large  $n$ ,

- $\hat{\theta}_n$  is distributed **approximately** as  $N(\theta_0, n^{-1}\mathcal{J}(\theta_0)^{-1}) \approx N(\theta_0, I(\theta_0)^{-1})$ ;
- if  $\psi(\theta)$  is differentiable near  $\theta_0$ , then  $\hat{\psi}_n = \psi(\hat{\theta}_n)$  is distributed **approximately** as  $N(\psi(\theta_0), \boldsymbol{\alpha}(\theta_0)^\top I(\theta_0)^{-1} \boldsymbol{\alpha}(\theta_0))$ , where  $\boldsymbol{\alpha}(\theta) = \left[ \frac{\partial \psi(\theta)}{\partial \theta_1}, \dots, \frac{\partial \psi(\theta)}{\partial \theta_k} \right]^\top$ .

### 5.5.8 **Example §5.4.1** (cont'd)

True values of  $(\beta, \nu)$ :  $(\beta_0, \nu_0)$

Likelihood function:  $\ell_n(\beta, \nu) \propto (2\pi\nu)^{-n/2} \exp \left[ -\frac{1}{2\nu} \sum_{i=1}^n \{X_i - \mu_\beta(u_i)\}^2 \right]$

Likelihood equations:

$$\begin{cases} \frac{\partial}{\partial \beta} S_n(\beta, \nu) = \nu^{-1} \sum_{i=1}^n \{X_i - \mu_\beta(u_i)\} \frac{\partial \mu_\beta(u_i)}{\partial \beta} = \mathbf{0} \\ \frac{\partial}{\partial \nu} S_n(\beta, \nu) = -\frac{n}{2\nu} + \frac{1}{2\nu^2} \sum_{i=1}^n \{X_i - \mu_\beta(u_i)\}^2 = 0 \end{cases}$$

$\Rightarrow$  mle of  $(\beta_0, \nu_0)$  is  $(\hat{\beta}_n, \hat{\nu}_n)$ , where

$$\sum_{i=1}^n \{X_i - \mu_{\hat{\beta}_n}(u_i)\} \frac{\partial \mu_{\hat{\beta}_n}(u_i)}{\partial \beta} \Big|_{\beta=\hat{\beta}_n} = \mathbf{0} \quad \text{and} \quad \hat{\nu}_n = n^{-1} \sum_{i=1}^n \{X_i - \mu_{\hat{\beta}_n}(u_i)\}^2.$$

Results of §5.4.7 imply that

$$\mathcal{J}(\beta, \nu) = \lim_{n \rightarrow \infty} n^{-1} I(\beta, \nu) = \begin{bmatrix} \nu^{-1} \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n \frac{\partial \mu_\beta(u_i)}{\partial \beta} \frac{\partial \mu_\beta(u_i)}{\partial \beta^\top} & \mathbf{0} \\ \mathbf{0}^\top & 1/(2\nu^2) \end{bmatrix}$$

and

$$\mathcal{J}(\beta, \nu)^{-1} = \begin{bmatrix} \nu \left\{ \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n \frac{\partial \mu_\beta(u_i)}{\partial \beta} \frac{\partial \mu_\beta(u_i)}{\partial \beta^\top} \right\}^{-1} & \mathbf{0} \\ \mathbf{0}^\top & 2\nu^2 \end{bmatrix}.$$

Large-sample properties of mle (Theorem §5.5.6(ii))  $\Rightarrow$

$$n^{1/2} \left( \begin{bmatrix} \hat{\beta}_n \\ \hat{\nu}_n \end{bmatrix} - \begin{bmatrix} \beta_0 \\ \nu_0 \end{bmatrix} \right) \longrightarrow N(\mathbf{0}, \mathcal{J}(\beta_0, \nu_0)^{-1}) \text{ in distribution.}$$

Large-sample properties of score function (Theorem §5.5.6(iii))  $\Rightarrow$

$$n^{-1/2} \begin{bmatrix} \frac{1}{\nu_0} \sum_{i=1}^n \{X_i - \mu_{\beta_0}(u_i)\} \frac{\partial \mu_{\beta}(u_i)}{\partial \beta} \Big|_{\beta=\beta_0} \\ \frac{1}{2\nu_0^2} \sum_{i=1}^n \left( \{X_i - \mu_{\beta_0}(u_i)\}^2 - \nu_0 \right) \end{bmatrix} \longrightarrow N(\mathbf{0}, \mathcal{J}(\beta_0, \nu_0)) \text{ in distribution.}$$

## §5.6 Exercise: comparison of UMVUE and MLE

- 5.6.1 Let  $\mathbf{X} \subset \{1, \dots, \theta\}$  be a random subset of positive integers, for some unknown  $\theta \in \{1, 2, \dots\}$ . Let  $p = 1 - q \in (0, 1)$  be a known probability such that the  $\theta$  events  $\{1 \in \mathbf{X}\}, \dots, \{\theta \in \mathbf{X}\}$  are independent and each occurs with probability  $p$ .

We wish to estimate  $\theta$  based on the data  $\mathbf{X}$ .

### Question 1

Show that  $\mathbf{X}$  has the mass function

$$f(\mathbf{X}|\theta) = p^{\#\mathbf{X}} q^{\theta - \#\mathbf{X}} \mathbf{1}\{\mathbf{X} \subset \{1, \dots, \theta\}\},$$

where  $\#\mathbf{X}$  denotes the number of elements in  $\mathbf{X}$ .

### Question 2

Find a complete sufficient statistic for  $\theta$ .

### Question 3

Find the maximum likelihood estimator (mle) of  $\theta$ .

### Question 4

Find a uniformly minimum variance unbiased (UMVU) estimator of  $\theta$ .

### Question 5

Compare the MSE's of the maximum likelihood and UMVU estimators.

## §5.7 Nonparametric estimation

- 5.7.1 Classical parametric (or *model-based*) inference requires specification of a parametric model  $\{f(\cdot|\theta) : \theta \in \Theta\}$ , with parameter space  $\Theta \subset \mathbb{R}^k$ .

*Nonparametric* (or *model-free*) inference attempts to yield more *robust* results by dropping such stringent model assumptions.

5.7.2 Let  $\mathbf{X} = (X_1, \dots, X_n)$  be  $n$  independent replicates of  $X \sim$  distribution function  $F \in \mathcal{F}$ , where  $\mathcal{F}$  denotes the class of **all** possible distribution functions of  $X$ . Based on the observed sample  $\mathbf{x} = (x_1, \dots, x_n)$ , the *nonparametric likelihood function* is

$$\ell_{\mathbf{x}}(F) = \prod_{i=1}^n \mathbb{P}_F(X = x_i), \quad F \in \mathcal{F}.$$

Here  $\mathbb{P}_F(X = x_i)$  should be interpreted as

*the probability of observing  $x_i$  if  $F$  were the true underlying distribution function.*

**Note:** If  $F$  is continuous, then  $\mathbb{P}_F(X = x_i) = 0$  and so  $\ell_{\mathbf{x}}(F) = 0$ .

5.7.3 **Definition.** The *empirical distribution* of  $\mathbf{X} = (X_1, \dots, X_n)$  is the distribution of a discrete random variable  $X^*$  with the mass function

$$\mathbb{P}(X^* = X_i | \mathbf{X}) = \frac{1}{n}, \quad \text{for } i = 1, 2, \dots, n, \quad \text{conditional on } \mathbf{X}.$$

Thus, the empirical distribution places on each observation  $X_i$  a probability mass of  $1/n$ , or equivalently,  $X^* = X_J$ , where  $J$  is randomly selected (with equal probabilities) from the set of integers  $\{1, 2, \dots, n\}$ .

5.7.4 The cdf of the empirical distribution is given by

$$\hat{F}_n(t) = \mathbb{P}(X^* \leq t | \mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i \leq t\}.$$

**Note:** “ $X_i \leq t$ ” should be interpreted *componentwise* if  $X_i$  and  $t$  are vectors.

5.7.5 The nonparametric likelihood function  $\ell_{\mathbf{X}}(F)$  is maximised by the empirical distribution function  $\hat{F}_n$  of the sample  $\mathbf{X} = (X_1, \dots, X_n)$ .

Thus we may regard  $\hat{F}_n$  ( $\in \mathcal{F}$  by assumption) as a *nonparametric* mle of  $F$ .

Consequently, any given functional  $\theta(F)$  of  $F$  can be estimated by the *nonparametric* mle  $\theta(\hat{F}_n)$ .

5.7.6 **Examples.**

- (i) If  $\theta(F) = \int x dF(x)$  (i.e. mean of  $F$ ), then  $\theta(\hat{F}_n) = \int x d\hat{F}_n(x) = \bar{X}$  (i.e. sample mean of  $X_1, \dots, X_n$ ).

- (ii) If  $\theta(F) = \int x^2 dF(x) - \left(\int x dF(x)\right)^2$  (i.e. variance of  $F$ ), then  $\theta(\hat{F}_n) = n^{-1} \sum_{i=1}^n X_i^2 - \bar{X}^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$  (i.e. sample variance of  $X_1, \dots, X_n$ ).
- (iii) If  $\theta(F) = F^{-1}(\xi)$  (i.e.  $\xi^{\text{th}}$  population quantile of  $F$ ), then  $\theta(\hat{F}_n) = \hat{F}_n^{-1}(\xi)$  (i.e.  $\xi^{\text{th}}$  sample quantile of  $X_1, \dots, X_n$ ).
- (iv) (*linear regression*)

Let  $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$  be  $n$  i.i.d. replicates of  $(\mathbf{X}, Y) \in \mathbb{R}^p \times \mathbb{R}$ , which follows an unknown  $(p+1)$ -variate distribution  $F$ . Define

$$\theta(F) = \psi(\beta_0, \boldsymbol{\beta}) = \psi\left(\underset{b_0, \mathbf{b}}{\operatorname{argmin}} \mathbb{E}_F[L(Y - b_0 - \mathbf{b}^\top \mathbf{X})]\right),$$

for some loss function  $L(\cdot)$  and some given real-valued function  $\psi(\cdot)$ . Then

$$\theta(\hat{F}_n) = \psi(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) = \psi\left(n^{-1} \underset{b_0, \mathbf{b}}{\operatorname{argmin}} \sum_{i=1}^n L(Y_i - b_0 - \mathbf{b}^\top \mathbf{X}_i)\right).$$

Exercise: Write down explicit expressions for  $(\beta_0, \boldsymbol{\beta})$  and  $(\hat{\beta}_0, \hat{\boldsymbol{\beta}})$  in the special case where  $L(u) = u^2$ .

- (v) Consider two independent random samples,  $\mathbf{X} = (X_1, \dots, X_n)$  and  $\mathbf{Y} = (Y_1, \dots, Y_m)$ , drawn respectively from the distribution functions  $F$  and  $G$ . Define a functional

$$\theta(F, G) = \mathbb{P}_{F, G}(X_1 > cY_1),$$

for a given constant  $c$ . Denote by  $\hat{F}_n$  and  $\hat{G}_m$  the empirical cdf's of the samples  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively. Then

$$\theta(\hat{F}_n, \hat{G}_m) = m^{-1} \sum_{j=1}^m \{1 - \hat{F}_n(cY_j)\} = (mn)^{-1} \sum_{j=1}^m \sum_{i=1}^n \mathbf{1}\{X_i > cY_j\}.$$

If it is known that  $F = G$ , then we should estimate  $F$  and  $G$  by the empirical cdf  $\hat{H}_{n+m}$  of the pooled sample  $(\mathbf{X}, \mathbf{Y}) = (W_1, \dots, W_{n+m})$ . In this case, we have

$$\theta(\hat{H}_{n+m}, \hat{H}_{n+m}) = (n+m)^{-1} \sum_{i=1}^{n+m} \{1 - \hat{H}_{n+m}(cW_i)\} = (n+m)^{-2} \sum_{i=1}^{n+m} \sum_{j=1}^{n+m} \mathbf{1}\{W_j > cW_i\}.$$

## §5.8 Bootstrap estimation

5.8.1 Suppose  $\mathbf{X} = (X_1, \dots, X_n)$  are i.i.d.  $\sim F \in \mathcal{F}$ . Very often inference about  $F$  calls for knowledge of the sampling distribution function of some real-valued statistic  $T(\mathbf{X}, F)$ :

$$F_T(y|F) \triangleq \mathbb{P}_F(T(\mathbf{X}, F) \leq y), \quad y \in \mathbb{R}.$$

Example: For a problem of estimating  $\theta(F)$ ,  $T(\mathbf{X}, F)$  may be defined as the  $L_p$  norm  $\|\hat{\theta}(\mathbf{X}) - \theta(F)\|_p$ , which measures the estimation error of the estimator  $\hat{\theta}(\mathbf{X})$ . Then knowledge of  $F_T(\cdot|F)$  enables us to investigate, for example, mean squared error or mean absolute deviation of  $\hat{\theta}(\mathbf{X})$ , among other applications.

5.8.2 If  $F$  were known, exact answers would have been available, derived by either *analytical calculations* or *Monte Carlo simulation*.

If  $F$  is unknown (as is often the case), one may sometimes resort to asymptotic approximations, usually derived from a certain kind of *central limit theorem* (for some properly chosen normalising constants  $a_n$ ):

$$a_n T(\mathbf{X}, F) \rightarrow N(0, \sigma_T^2) \text{ in distribution} \Rightarrow F_T(y|F) \approx \Phi(a_n y / \sigma_T) \text{ for large } n.$$

What if  $\sigma_T$  is **NOT** available and is very hard to estimate?

What if sample size  $n$  is **NOT** large?

5.8.3 The *bootstrap* method provides an approximate solution to the problem by applying the “*plug-in*” idea:

- plug in the empirical distribution  $\hat{F}_n$  as a substitute for  $F$  and

$$\text{estimate } F_T(y|F) \text{ by } F_T(y|\hat{F}_n) = \mathbb{P}_{\hat{F}_n}(T(\mathbf{X}^*, \hat{F}_n) \leq y), \quad y \in \mathbb{R},$$

where  $\mathbf{X}^* = (X_1^*, \dots, X_n^*)$  denotes an *i.i.d.* sample drawn from  $\hat{F}_n$ .

The estimate  $F_T(\cdot|\hat{F}_n)$  is known as the bootstrap distribution function, and  $\mathbf{X}^*$  is known as a *bootstrap sample* of size  $n$ .

Note: In some problem settings, it may be more appropriate to substitute  $F$  by an alternative estimate different from  $\hat{F}_n$  and generate bootstrap samples from that alternative estimate accordingly.

5.8.4 In practice, the *bootstrap distribution*  $F_T(\cdot|\hat{F}_n)$  can be approximated by *Monte Carlo simulation*:

$$F_T(y|\hat{F}_n) \approx B^{-1} \sum_{b=1}^B \mathbf{1}\{T(\mathbf{X}^{*b}, \hat{F}_n) \leq y\}, \quad y \in \mathbb{R},$$

where  $\mathbf{X}^{*1}, \dots, \mathbf{X}^{*B}$  are  $B$  independent bootstrap samples, each of size  $n$ , drawn from  $\hat{F}_n$ .

Noting that  $\hat{F}_n$  is the empirical cdf of  $\mathbf{X} = (X_1, \dots, X_n)$ , in practice we may generate each  $\mathbf{X}^{*b}$  by sampling  $n$  observations with replacement from  $\mathbf{X}$ . Thus the bootstrap is often referred to as a *resampling method*.



### 5.8.5 Example.

Let  $X_1, \dots, X_n \in \mathbb{R}$  be ordered as  $X_{(1)} \leq \dots \leq X_{(n)}$ .

Consider:  $\theta = \theta(F) = \mathbb{E}_F X_1$ ,  $\hat{\theta} = \hat{\theta}(\mathbf{X}) = \frac{1}{n-2m} \sum_{i=m+1}^{n-m} X_{(i)}$ ,  $T(\mathbf{X}, F) = |\hat{\theta} - \theta|$ .

Then  $\theta(\hat{F}_n) = \mathbb{E}_{\hat{F}_n} X_1^* = \bar{X}$  and

$$\hat{\theta}^* = \hat{\theta}(\mathbf{X}^*) = \frac{1}{n-2m} \sum_{i=m+1}^{n-m} X_{(i)}^*,$$

where  $X_{(1)}^* \leq \dots \leq X_{(n)}^*$  denotes the ordered sequence of the  $n$  observations in the bootstrap sample  $\mathbf{X}^* = (X_1^*, \dots, X_n^*)$ .

The bootstrap method estimates  $F_T(y|F)$  by:

$$F_T(y|\hat{F}_n) = \mathbb{P}_{\hat{F}_n}(|\hat{\theta}^* - \theta(\hat{F}_n)| \leq y) = \mathbb{P}_{\hat{F}_n} \left( \left| \frac{1}{n-2m} \sum_{i=m+1}^{n-m} X_{(i)}^* - \bar{X} \right| \leq y \right).$$

*Monte Carlo simulation* procedure for approximating  $F_T(y|\hat{F}_n)$ :

1. simulate from  $\hat{F}_n$  a large number,  $B$  say, of bootstrap samples

$$\mathbf{X}^{*1} = (X_1^{*1}, \dots, X_n^{*1}), \dots, \mathbf{X}^{*B} = (X_1^{*B}, \dots, X_n^{*B}),$$

2. for each  $b = 1, \dots, B$ , calculate  $\hat{\theta}^{*b} = \hat{\theta}(\mathbf{X}^{*b}) = \frac{1}{n-2m} \sum_{i=m+1}^{n-m} X_{(i)}^{*b}$  based on the ordered sequence  $X_{(1)}^{*b} \leq \dots \leq X_{(n)}^{*b}$ ,

3. approximate  $F_T(y|\hat{F}_n)$  by

$$B^{-1} \sum_{b=1}^B \mathbf{1} \left\{ |\hat{\theta}^{*b} - \theta(\hat{F}_n)| \leq y \right\} = B^{-1} \sum_{b=1}^B \mathbf{1} \left\{ \left| \frac{1}{n-2m} \sum_{i=m+1}^{n-m} X_{(i)}^{*b} - \bar{X} \right| \leq y \right\}.$$

### 5.8.6 Exercise: (Example §5.7.6(v))

Use the bootstrap method to estimate the distribution of the estimation error  $T = \theta(\hat{F}_n, \hat{G}_m) - \theta(F, G)$ .

Derive from the bootstrap distribution estimates of the following performance indicators of  $\theta(\hat{F}_n, \hat{G}_m)$ :

bias, mean squared error, mean absolute deviation.