MATH357 - Statistics

Based on lectures from Winter 2025 by Prof. Abbas Khalili. Notes by Louis Meunier

Contents

1 Review of Probability	2
2 Common Statistical Tools	
2.1 Definition of Statistics	6
2.2 Properties of Normal and other Common Distributions	7
2.3 Order Statistics	10
2.4 Large Sample/Asymptotic Theory	12
3 Parametric Inference	14
3.1 Uniformly Minimum Variance Unbiased Estimators (UMVUE), Cramér-Rau Lower Bound (CRLB)	17
3.2 Sufficiency	23
3.3 Completeness	
3.4 Existence of a UMVUE	32

§1 Review of Probability

⇒ Definition 1.1 (Measurable Space, Probability Space): We work with a set Ω = sample space = {outcomes}, and a σ -algebra \mathcal{F} , which is a collection of subsets of Ω containing Ω and closed under taking complements and countable unions. The tuple (Ω, \mathcal{F}) is called *measurable space*.

We call a nonnegative function $P: \mathcal{F} \to \mathbb{R}$ defined on a measurable space a *probability* function if $P(\Omega) = 1$ and if $\{E_n\} \subseteq \mathcal{F}$ a disjoint collection of subsets of Ω , then $P(\bigcup_{n \geq 1} E_n) = \sum_{n \geq 1} P(E_n)$. We call the tuple (Ω, \mathcal{F}, P) a *probability space*.

 \hookrightarrow Definition 1.2 (Random Variables): Fix a probability space (Ω, \mathcal{F}, P) . A Borel-measurable function $X : \Omega \to \mathbb{R}$ (namely, $X^{-1}(B) \in \mathcal{F}$ for every $B \in \mathfrak{B}(\mathbb{R})$) is called a *random variable* on \mathcal{F} .

- *Probability distribution*: X induces a probability distribution on $\mathfrak{B}(\mathbb{R})$ given by $P(X \in B)$
- *Cumulative distribution function (CDF)*:

$$F_X(x) := P(X \le x).$$

Note that $F(-\infty) = 0$, $F(+\infty) = 1$ and F right-continuous.

We say X discrete if there exists a countable set $S := \{x_1, x_2, ...\} \subset \mathbb{R}$, called the *support* of X, such that $P(X \in S) = 1$. Putting $p_i := P(X = x_i)$, then $\{p_i : i \ge 1\}$ is called the *probability mass function* (PMF) of X, and the CDF of X is given by

$$P(X \le x) = \sum_{i: x_i \le x} p_i.$$

We say X continuous if there is a nonnegative function f, called the *probability distribution* function (PDF) of X such that $F(x) = \int_{-\infty}^{x} f(t) dt$ for every $x \in \mathbb{R}$. Then,

- $\forall B \in \mathfrak{B}(\mathbb{R}), P(X \in B) = \int_B f(t) dt$
- F'(x) = f(x)
- $\int_{-\infty}^{\infty} f(x) \, \mathrm{d}x = 1$

If $X : \Omega \to \mathbb{R}$ a random variable and $g : \mathbb{R} \to \mathbb{R}$ a Borel-measurable function, then $Y := g(X) : \Omega \to \mathbb{R}$ also a random variable.

1 Review of Probability

Definition 1.3 (Moments): Let *X* be a discrete/random random variable with pmf/pdf *f* and support *S*. Then, if $\sum_{x \in S} |x| f(x) / \int_{S} |x| f(x) dx < \infty$, then we say the first moment/mean of *X* exists, and define

$$\mu_X = \mathbb{E}[X] = \begin{cases} \sum_{x \in S} x f(x) \\ \int_S x f(x) \, \mathrm{d}x \end{cases}.$$

Let $g : \mathbb{R} \to \mathbb{R}$ be a Borel-measurable function. Then, we have

$$\mathbb{E}[g(X)] = \begin{cases} \sum_{x \in S} g(x) f(x) \\ \int_{S} g(x) f(x) \end{cases}.$$

Taking $g(x) = |x|^k$ gives the so-called "kth absolute moments", and $g(x) = x^k$ gives the ordinary "kth moments". Notice that $\mathbb{E}[\cdot]$ linear in its argument.

For $k \ge 1$, if μ exists, define the central moments

$$\mu_k \coloneqq \mathbb{E}\Big[\left(X - \mu\right)^k\Big],$$

where they exist.

 \hookrightarrow **Definition 1.4** (Moment Generating Function (mgf)): If X a r.v., the mgf of X is given by

$$M(t) \coloneqq \mathbb{E}[e^{tX}],$$

if it exists for $t \in (-h, h)$, h > 0. Then, $M^{(n)}(0) = \mathbb{E}[X^n]$.

Definition 1.5 (Multiple Random Variable): $X = (X_1, ..., X_n) : \Omega \to \mathbb{R}^n$ a random vector if $X^{-1}(I) \in \mathcal{F}$ for every $I \in \mathfrak{B}_{\mathbb{R}^n}$. (It suffices to check for "rectangles" $I = (-\infty, a_1] \times \cdots \times (-\infty, a_n]$, as before.)

Let *F* be the CDF of *X*, and let $A \subseteq \{1, ..., n\}$, enumerating *A* by $\{i_1, ..., i_k\}$. Then, the CDF of the subvector $X_A = (X_{i_1}, ..., X_{i_k})$ is given by

$$F_{X_A}(x_{i_1},...,x_{i_k}) = \lim_{\substack{x_{i_j} \to \infty, \\ i_j \in \mathcal{I} \setminus A}} F(x_1,...,x_n).$$

In particular, the marginal distribution of X_i is given by

$$F_{X_i}(x) = \lim_{x_1,...,x_{i-1},x_{i+1},...,x_n \to +\infty} F(x_1,...,x,...,x_n).$$

Let $g: \mathbb{R}^n \to \mathbb{R}$ measurable. Then,

$$\mathbb{E}[g(X_1,...,X_n)] = \begin{cases} \sum_{(x_1,...,x_n)} g(x_1,...,x_n) f(x_1,...,x_n) \\ \int \cdots \int g(x_1,...,x_n) f(x_1,...,x_n) \, \mathrm{d} x_1 \cdots \, \mathrm{d} x_n \end{cases}.$$

We have the notion of a joint mgf,

$$M(t_1,...,t_n) = \mathbb{E}\left[e^{\sum_{i=1}^n t_i X_i}\right],$$

if it exists for $0 < \left(\sum_{i=1}^n t_i^2\right)^{\frac{1}{2}} < h$ for some h > 0. Notice that $M(0, ..., 0, t_i, 0, ..., 0)$ is equal to the mgf of X_i .

1 Review of Probability

Definition 1.6 (Conditional Probability): Let $(X_1,...,X_n)$ a random vector. Let $\mathcal{I} = \{1,...,n\}$ and A,B disjoint subsets of \mathcal{I} with k := |A|, h := |B|. Write $X_A = (X_{i_1},...,X_{i_k})^t$, similar for B. Then, the conditional probability of A given B is given by

$$f_{X_A|X_B}(x_a|x_b) := f_{X_A|X_B = x_B}(x_A) = \frac{f_{X_A,X_B}(x_a,x_b)}{f_{X_b}(x_b)},$$

provided the denominator is nonzero. Sometimes we have information about conditional probabilities but not the main probability function; we have that

$$f(x_1,...,x_n) = f(x_1)f(x_2 \mid x_1)f(x_3 \mid x_1, x_2) \cdots f(x_n \mid x_1,...,x_{n-1}),$$

which follows from expanding the previous definition and observing the cancellation.

Let $X = (X_1, ..., X_n) \sim F$. We say $X_1, ..., X_n$ (mutually) independent and write $\coprod_{i=1}^n X_i$ if

$$F(x_1,...,x_n) = \prod_{i=1}^n F_{X_i}(x_i),$$

where F_{X_i} the marginal cdf of X_i . Equivalently,

$$\prod_{i=1}^{n} X_i \Leftrightarrow f(x_1, ..., x_n) = \prod_{i=1}^{n} f_{X_i}(x_i)$$

$$\Leftrightarrow P(X_1 \in B_1, ..., X_n \in B_n) = \prod_{i=1}^{n} P(X_i \in B_i) \ \forall \ B_i \in \mathfrak{B}_{\mathbb{R}}$$

$$\Leftrightarrow M_X(t_1, ..., t_n) = \prod_{i=1}^{n} M_{X_i}(t_i).$$

If X, Y are two random variables with cdfs F_X , F_Y such that $F_X(z) = F_Y(z)$ for every z, we say X, Y identically distributed and write $X \stackrel{d}{=} Y$ (note that X need not equal Y pointwise). If $X_1, ..., X_n$ a collection of random variables that are independent and identically distributed with common cdf F, we write $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F$.

Further, define the covariance, correlation of two random variables *X*, *Y* respectively:

$$\operatorname{Cov}(X,Y) \coloneqq \sigma_{X,Y} \coloneqq \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mu_X \mu_Y, \qquad \rho_{X,Y} \coloneqq \frac{\sigma_{XY}}{\sigma_X \sigma_Y},$$

$$if \, \mathbb{E}[|X - \mathbb{E}[X]| \, |Y - \mathbb{E}[Y]|] < \infty.$$

Theorem 1.1: If $X_1, ..., X_n$ independent and $g_1, ..., g_n : \mathbb{R} \to \mathbb{R}$ borel-measurable functions, then $g_1(X_1), ..., g_n(X_n)$ also independent.

1 Review of Probability 5

Definition 1.7 (Conditional Expectation): Let *X*, *Y* be random variables and *g* : \mathbb{R} → \mathbb{R} a borel-measurable function. We define the following notions:

$$\mathbb{E}[g(X)|Y = y] = \begin{cases} \sum_{x \in S_X} g(x) f(x|y) \text{ discrete} \\ \int_{S_X} g(x) f(x|y) dx \text{ cnts} \end{cases}$$

$$\text{Var}(X|Y = y) = \mathbb{E}[X^2|Y = y] - \mathbb{E}^2[X|Y = y].$$

Theorem 1.2: If $\mathbb{E}[g(X)]$ exists, then $\mathbb{E}[g(X)] = \mathbb{E}[\mathbb{E}[g(X)|Y]]$, where the first nested \mathbb{E} is with respect to x, the second y.

Theorem 1.3: If $\mathbb{E}[X^2]$ < ∞, then $Var(X) = Var(\mathbb{E}[X|Y]) + \mathbb{E}[Var(X|Y)]$. In particular, $Var(X) \ge Var(\mathbb{E}[X|Y])$.

§2 Common Statistical Tools

§2.1 Definition of Statistics

- ⇒ Definition 2.1 (Inference): We consider some population with some characteristic we wish to study. We can model this characteristic as a random variable $X \sim F$. In general, we don't have access to F, but wish to take samples from our population to make inferences about its properties.
- (1) *Parametric inference*: in this setting, we assume we know the functional form of X up to some parameter, $\theta \in \Theta \subset \mathbb{R}^d$, where Θ our "parameter space". Namely, we know $X \sim F_\theta \in \mathcal{F} := \{F_\theta \mid \theta \in \Theta\}$.
- (2) *Non-parametric inference:* in this setting we know noting about *F* itself, except perhaps that *F* continuous, discrete, etc.

Other types exist. We'll focus on these two.

Definition 2.2 (Random Sample): Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F$. Then $X_1, ..., X_n$ called a *random sample* of the population.

We also call X_i the "pre-experimental data" (to be observed) and x_i the "post-experimental data" (been observed).

2.1 Definition of Statistics 6

 \hookrightarrow **Definition 2.3** (Statistics): Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F$ where X_i a d-dimensional random vector. Let

$$T: \underbrace{\mathbb{R}^d \times \mathbb{R}^d \times \cdots \times \mathbb{R}^d}_{n-\text{fold}} \to \mathbb{R}^k$$

be a borel-measurable function. Then, $T(X_1,...,X_n)$ is called a *statistic*, provided it does not depend on any unknown.

Example 2.1: $\overline{X_n} := \frac{1}{n} \sum_{i=1}^n X_i$ (the "sample mean") and $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n \left(X_i - \overline{X_n} \right)^2$, (the "sample variance") are both typical statistics.

§2.2 Properties of Normal and other Common Distributions

 \hookrightarrow Theorem 2.1: Let $x_1,...,x_n \in \mathbb{R}$, then

(a)
$$\operatorname{argmin}_{\alpha \in \mathbb{R}} \left\{ \sum_{i=1}^{n} (x_i - \alpha)^2 \right\} = \overline{x_n};$$

(b)
$$\sum_{i=1}^{n} (x_i - \overline{x_n})^2 = \sum_{i=1}^{n} (x_i^2) - n\overline{x_n}^2$$
;

(c)
$$\sum_{i=1}^{n} (x_i - \overline{x_n}) = 0$$
.

Theorem 2.2: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F$, and $g : \mathbb{R} \to \mathbb{R}$ borel-measurable such that $\text{Var}(g(X)) < \infty$. Then,

(a)
$$\mathbb{E}\left[\sum_{i=1}^{n} g(X_i)\right] = n\mathbb{E}[g(X_1)];$$

(b)
$$\operatorname{Var}\left(\sum_{i=1}^{n} g(X_i)\right) = n \operatorname{Var}(X_1)$$
.

Theorem 2.3: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F$ with $\sigma^2 < \infty$, then

1.
$$\mathbb{E}\left[\overline{X_n}\right] = \mu$$
, $\operatorname{Var}\left(\overline{X_n}\right) = \frac{\sigma^2}{n}$, $\mathbb{E}\left[S_n^2\right] = \sigma^2$.

2. If $M_{X_1}(t)$ exists in some neighborhood of 0, then $M_{\overline{X_n}}(t) = M_{X_1}(\frac{t}{n})^n$, where it exists.

→Theorem 2.4: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$. Then

1.
$$\overline{X_n} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n});$$

2. $\overline{X_n}$, S_n^2 are independent;

3.
$$\frac{(n-1)S_n^2}{\sigma^2} = \frac{\sum_{i=1}^n (X_i - \overline{X_n})^2}{\sigma^2} \sim \chi_{(n-1)}^2.$$

Remark 2.1:

2. actually holds iff the underlying distribution is normal.

PROOF. We prove 2. We first write S_n^2 as a function of $(X_2 - \overline{X}_n, ..., X_n - \overline{X}_n)$:

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X}_n)^2 = \frac{1}{n-1} \left\{ \sum_{i=2}^n (X_i - \overline{X}_n)^2 + (X_1 - \overline{X}_n)^2 \right\}$$
$$= \frac{1}{n-1} \left\{ \sum_{i=2}^n (X_i - \overline{X}_n)^2 + \left(\sum_{i=2}^n (X_i - \overline{X}_n) \right)^2 \right\}.$$

Then, it suffices to show that \overline{X}_n and $(X_2 - \overline{X}_n, ..., X_n - \overline{X}_n)$ are independent.

Consider now the transformation

$$\begin{cases} y_1 = \overline{x}_n \\ y_2 = x_2 - \overline{x}_n \\ \vdots \\ y_n = x_n - \overline{x}_n \end{cases} \Rightarrow \begin{cases} x_1 = y_1 - \sum_{i=2}^n y_i \\ x_2 = y_2 + y_1 \\ \vdots \\ x_n = y_n + y_1 \end{cases},$$

which gives Jacobian

$$|J| = \begin{vmatrix} \begin{pmatrix} 1 & -1 & \cdots & -1 \\ 1 & 1 & 0 & \cdots \\ \vdots & \ddots & \ddots & \vdots \\ 1 & 0 & \cdots & 1 \end{vmatrix} = n,$$

and so

$$\begin{split} f_{Y_{1},...,Y_{n}}(y_{1},...,y_{n}) &= |J| \cdot f_{X_{1},...,X_{n}}(x_{1}(y_{1},...,y_{n}),...,x_{n}(y_{1},...,y_{n})) \\ &= n \cdot \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^{2}}} e^{-\frac{1}{2\sigma^{2}}(x_{i}(y_{1},...,y_{n})-\mu)^{2}} \\ &\approx \underbrace{e^{-n\frac{(y_{1}-\mu)^{2}}{2\sigma^{2}}} \cdot \underbrace{e^{-\frac{1}{2\sigma^{2}}\left\{\left(\sum_{i=2}^{n}y_{i}\right)^{2} + \sum_{i=2}^{n}y_{i}^{2}\right\}}_{\text{no } y_{1} \text{ dependence}}, \end{split}$$

and hence as the PDFs split, we conclude Y_1 independent of $Y_2, ..., Y_n$ and so \overline{X}_n independent of $(X_2 - \overline{X}_n, ..., X_n - \overline{X}_n)$ and so in particular of any Borel-measurable function of this vector such as S_n^2 , completing the proof.

For 3, note that

$$\begin{split} V \coloneqq \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma}\right)^2 &= \frac{1}{\sigma^2} \sum_{i=1}^n \left(\left(X_i - \overline{X}_n\right) - \left(\mu - \overline{X}_n\right)\right)^2 \\ &= \frac{\sum_{i=1}^n \left(X_i - \overline{X}_n\right)^2}{\sigma^2} + \frac{n\left(\overline{X}_n - \mu\right)^2}{\sigma^2} =: W_1 + W_2. \end{split}$$

The first part, W_1 , of this summation is just $(n-1)\frac{S_n^2}{\sigma^2}$, a function of S_n^2 , and the second, W_2 , is a function of \overline{X}_n . By what we've just shown in the previous part, these two are independent. In addition, $V \sim \chi^2_{(n)}$ and

$$W_2 = \frac{n \left(\overline{X}_n - \mu\right)^2}{\sigma^2} = \left(\frac{\overline{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}}\right)^2 \sim \chi_{(1)}^2,$$

since the inner random variable is a standard normal. Then, since W_1,W_2 independent, $M_V(t)=M_{W_1}(t)M_{W_2}(t)$, so for $t<\frac{1}{2}$,

$$M_{W_1}(t) = \frac{M_V(t)}{M_{W_2}(t)} = \frac{(1-2t)^{-\frac{n}{2}}}{(1-2t)^{-\frac{1}{2}}} = (1-2t)^{-\frac{(n-1)}{2}},$$

hence $W_1 \sim \chi^2_{(n-1)}$.

 \hookrightarrow **Proposition 2.1**: Let $X \sim t(\nu)$, the Student *t*-distribution i.e

$$f(x) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\pi\nu} \cdot \Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}},$$

then

- $Var(X) = \frac{\nu}{\nu 2}$ for $\nu > 2$
- If $Z \sim \mathcal{N}(0,1)$ and $V \sim \chi^2_{(\nu)}$ are independent random variables, then $T = \frac{Z}{\sqrt{V/\nu}} \sim t(\nu)$.

→Theorem 2.5: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$. Then,

$$T = \frac{\overline{X}_n - \mu}{\sqrt{S_n^2/n}} = \frac{\sqrt{n}(\overline{X}_n - \mu)}{S_n} \sim t(n-1).$$

Remark 2.2: By combing CLT and Slutsky's Theorem, T asymptotes to $\mathcal{N}(0,1)$, but this gives a general distribution. Note that for large n, t(n-1) approximately normal too.

PROOF. Notice that

$$W_1 \coloneqq \sqrt{n} \frac{\overline{X}_n - \mu}{\sigma} \sim \mathcal{N}(0,1), \qquad W_2 \coloneqq \frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{(n-1)}^2$$

are independent, and

$$T = \frac{W_1}{\sqrt{W_2/(n-1)}}$$

so by the previous proposition $T \sim t(n-1)$.

Proposition 2.2: Given $U \sim \chi^2_{(m)}$, $V \sim \chi^2_{(n)}$ independent, then $F = \frac{U/m}{V/n} \sim F(m,n)$. If $T \sim t(v)$, $T^2 \sim F(1,v)$.

Theorem 2.6: Let $X_1, ..., X_m \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_1, \sigma_1^2)$ and $Y_1, ..., Y_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_2, \sigma_2^2)$ be mutually independent random samples. Then,

$$F = \frac{S_m^2/\sigma_1^2}{S_n^2/\sigma_2^2} \sim F(m-1, n-1).$$

PROOF. We have that $U=\frac{(m-1)S_m^2}{\sigma_1^2}\sim \chi_{(m-1)}^2$ and $V=\frac{(n-1)S_n^2}{\sigma_2^2}$ are independent so by the previous proposition

$$F = \frac{U/(m-1)}{V/(n-1)} \sim F(m-1, n-1).$$

§2.3 Order Statistics

Definition 2.4: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F$. Then, the *order statistics* are

$$X_{(1)} \leq X_{(2)} \leq \cdots \leq X_{(n)},$$

where $X_{(i)}$ the *i*th largest of $X_1, ..., X_n$.

→ Definition 2.5 (Related Functions of Order Statistcs): The sample range is defined

$$R_n := X_{(n)} - X_{(1)}.$$

The sample median is defined

$$M(X_1,...,X_n) := \begin{cases} X_{\left(\frac{n+1}{2}\right)} & \text{if } n \text{ odd} \\ X_{\left(\frac{n}{2}\right)} + X_{\left(\frac{n+1}{2}\right)} & \text{if } n \text{ even.} \end{cases}$$

→Theorem 2.7 (Distribution of Max, Min): Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F, f$.

(Discrete)

(a)
$$P(X_{(1)} = x) = [1 - F(x^{-})]^{n} - [1 - F(x)]^{n}$$

(b)
$$P(X_{(n)} = y) = [F(y)]^n - [F(y^-)]^n$$

(Continuous)

(c)
$$F_{X_{(1)}}(x) = P(X_{(1)} \le x) = 1 - [1 - F(x)]^n$$
, $f_{X_{(1)}}(x) = n \cdot f(x)[1 - F(x)]^{n-1}$

(d)
$$F_{X_{(n)}}(y) = [F(y)]^n$$
, $f_{X_{(n)}}(y) = n \cdot f(y)[F(y)]^{n-1}$

Proof. (a) Notice

$$P(X_{(1)} = x) = P(X_{(1)} \le x) - P(X_{(1)} < x).$$

2.3 Order Statistics

We have

$$\begin{split} P\big(X_{(1)} \leq x\big) &= 1 - P\big(X_{(1)} > x\big) \\ &= 1 - P\big(X_1 > x, X_2 > x, ..., X_n > x\big) \\ &= 1 - P\big(X_1 > x\big) P\big(x_2 > x\big) \cdots P\big(X_n > x\big) \\ &= 1 - \big[1 - F(x)\big]^n, \end{split}$$

and similarly

$$P(X_{(1)} < x) = 1 - P(X_{(1)} \ge x) = 1 - [1 - F(x^{-})]^{n}$$

where $F(x^-) = \lim_{z \to x^-} F(z)$. So in all,

$$P(X_{(1)} = x) = [1 - F(x^{-})]^{n} - [1 - F(x)]^{n}.$$

(b) is very similar. For (c), we have

$$P(X_{(1)} \le x) = 1 - P(X_{(1)} > x)$$

$$= 1 - P(X_1 > x, ..., X_n > x)$$

$$= 1 - [1 - F(x)]^n.$$

(d) is similar.

Theorem 2.8 (Distribution of *j*th Order Statistics): Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F, f$.

(*Discrete*) Suppose the X_i 's take values in $S_x = \{x_1, x_2, ...\}$ and put $p_i = P(X_i)$. Then,

$$F_{X_{(j)}}(x_i) = P(X_{(j)}(x_i) \le x_i) = \sum_{k=j}^n \binom{n}{k} P_i^k (1 - P_i)^{n-k},$$

where $P_i = P(X_i \le x_i) = \sum_{\ell=1}^i p_\ell$.

(Continuous)

$$F_{X_{(j)}}(x) = \sum_{k=j}^{n} {n \choose k} F^k(x) [1 - F(x)]^{n-k},$$

so

$$f_{X_{(j)}}(x) = \frac{n!}{(j-1)!(n-j)!} f(x) [F(x)]^{j-1} [1 - F(x)]^{n-j}.$$

Proof. For discrete, we have

$$P(X_{(j)}(x_i) \le x_i) = P(\text{at least } j \text{ out of } X_1, ..., X_n \le x_i).$$

Then,

$$P(\text{at least } j \text{ out of } X_1, ..., X_n \le x_i) = \sum_{k=i}^n \binom{n}{k} P_i^k (1 - P_i)^{n-k}.$$

Continuous is similar.

§2.4 Large Sample/Asymptotic Theory

 \hookrightarrow **Definition 2.6** (Convergence in Probability): We say $T_n = T(X_1, ..., X_n)$ converges *in* probability to θ $T_n \stackrel{P}{\to} \theta$ as $n \to \infty$ if for any $\varepsilon > 0$,

$$\lim_{n \to \infty} P(|T_n - \theta| > \varepsilon) = 0.$$

 \hookrightarrow **Definition 2.7** (Convergence in Distribution): Find a positive sequence $\{r_n\}$ with $r_n \to \infty$ such that

$$r_n(T_n-\theta)\stackrel{d}{\to} T$$
,

where *T* a random variable.

Theorem 2.9 (Slutsky's): Suppose $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{P} a$ for some $a \in \mathbb{R}$. Then,

$$X_n + Y_n \stackrel{d}{\to} X + a$$

$$X_n Y_n \stackrel{d}{\to} aX$$
,

and if $a \neq 0$,

$$\frac{X_n}{Y_n} \xrightarrow{d} \frac{X}{a}$$
.

→Theorem 2.10 (Continuous Mapping Theorem (CMT)): Suppose $X_n \stackrel{P}{\to} X$ and g is continuous on the set C such that $P(X \in C) = 1$. Then,

$$g(X_n) \stackrel{P}{\to} g(X).$$

Example 2.2: Let $X_1,...,X_n \stackrel{\text{iid}}{\sim} F$ with $\mu = \mathbb{E}[X_i], \sigma^2 = \text{Var}(X_i) < \infty$. Then,

$$\frac{\sqrt{n}(\overline{X}_n - \mu)}{S_n} \stackrel{d}{\to} \mathcal{N}(0, 1),$$

since we may rewrite

$$\frac{\sqrt{n}(\overline{X}_n - \mu)/\sigma}{S_n/\sigma}.$$

The numerator $\stackrel{d}{\to} \mathcal{N}(0,1)$ by CLT. $S_n^2 \stackrel{P}{\to} \sigma^2$, so the denominator goes to 1 in probability.

- \hookrightarrow **Definition 2.8** (Big O, Little o Notation): Let $\{a_n\}$, $\{b_n\}$ ⊆ \mathbb{R} real sequences.
- We say $a_n = O(b_n)$ if $\exists 0 < c \in \mathbb{R}$ and $N \in \mathbb{N}$ such that $|\frac{a_n}{b_n}| \le c$ for every $n \ge N$.
- We say $a_n = o(b_n)$ if $\lim_{n\to\infty} \frac{a_n}{b_n} = 0$.
- \hookrightarrow **Definition 2.9** (Big O_p , Little o_p Notation): Let $\{X_n\}$, $\{Y_n\}$ sequences of random variables.
- We say $X_n = O_p(1)$ if $\forall \ \varepsilon > 0$ there is a $N_\varepsilon \in \mathbb{N}$ and $C_\varepsilon \in \mathbb{R}$ such that

$$P(|X_n| > C_{\varepsilon}) < \varepsilon$$

for every $n > N_{\varepsilon}$.

- We say $X_n = O_p(Y_n)$ if $X_n/Y_n = O_p(1)$.
- We say $X_n = o_p(1)$ if $X_n \stackrel{P}{\to} 0$.
 - We say $X_n = o_p(Y_n)$ if $X_n/Y_n = o_p(1)$.
- \hookrightarrow **Proposition 2.3**: If $X_n \stackrel{d}{\to} X$, then $X_n = O_p(1)$.
- **Proposition 2.4** (The Delta Method (First Order)): Let $\sqrt{n}(X_n \mu) \stackrel{d}{\rightarrow} V$ and *g* a real-valued function such that *g'* exists at *x* = *μ* and *g'*(*μ*) ≠ 0. Then,

$$\sqrt{n}(g(X_n) - g(\mu)) \stackrel{d}{\to} g'(\mu)V.$$

In particular, if $V \sim \mathcal{N}(0, \sigma^2)$ then

$$\sqrt{n}(g(X_n) - g(\mu)) \stackrel{d}{\to} \mathcal{N}(0, g'(\mu)^2 \sigma^2).$$

PROOF. Taylor expanding the LHS,

$$\sqrt{n}\{g(X_n)-g(\mu)\}=g'(\mu)\sqrt{n}(X_n-\mu)+o_p(1)\to g'(\mu)V.$$

Proposition 2.5 (The Delta Method (Second Order)): Suppose $\sqrt{n}(X_n - \mu) \stackrel{d}{\to} \mathcal{N}(0, \sigma^2)$ and $g'(\mu) = 0$ but $g''(\mu) \neq 0$. Then,

$$n\{g(X_n) - g(\mu)\} \stackrel{d}{\to} \sigma^2 \frac{g''(\mu)}{2} \cdot \chi^2_{(1)}.$$

Proof.

$$g(X_n) - g(\mu) = \frac{g''(\mu)}{2} (X_n - \mu)^2 + o_p(1),$$

so

$$n(g(X_n)-g(\mu))=\sigma^2\frac{g''(\mu)}{2}\left\lceil\frac{\sqrt{n}(X_n-\mu)}{\sigma}\right\rceil^2+o_p(1).$$

The bracketed term converges in distribution to $\mathcal{N}(0,1)$ and the $o_p(1)$ term converges in probability to zero, so the proposition follows by applying Slutsky's Theorem.

Example 2.3:
$$\sqrt{n}(\overline{X}_n - \mu) \stackrel{d}{\to} \mathcal{N}(0, \sigma^2)$$
 by CLT. Letting $g(x) = x^2$, and assuming $\mu \neq 0$, then $\sqrt{n}(\overline{X}_n^2 - \mu^2) \to \mathcal{N}(0, 4\mu^2\sigma^2)$,

by the first-order delta method.

- **Proposition 2.6**: Let $X_1,...,X_n \stackrel{\text{iid}}{\sim} F$, and denote the ECDF $F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq x)$. Then,
- 1. $\mathbb{E}[F_n(x)] = F(x)$;
- 2. Var $(F_n(x)) = \frac{1}{n}F(x)(1 F(x));$
- 3. $nF_n(x) = \sum_{i=1}^n \mathbb{1}(X_i \le x) \sim \text{Bin}(n, F(x));$ 4. $\frac{\sqrt{n}(F_n(x) F(x))}{\sqrt{F(x)(1 F(x))}} \stackrel{d}{\to} \mathcal{N}(0, 1).$ 5. $F_n(x) \stackrel{P}{\to} F(x).$

- 6. $P(|F_n(x) F(x)| \ge \varepsilon) \le 2e^{-2n\varepsilon^2}$, by Hoeffing's Inequality.
- 7. $\sup_{x \in \mathbb{R}} |F_n(x) F(x)| = ||F_n F||_{\infty} \stackrel{\text{a.s.}}{\to} 0$, by the Glivenko-Cantelli Theorem.
- 8. $P(\|F_n F\|_{\infty} > \varepsilon) \le C\varepsilon e^{-2n\varepsilon^2}$ for some constant C (Dvoretzky-Kiefer-Wolfowitz Theorem).

Remark 2.3: The constant in 8. was shown to be 2 by Massart.

§3 PARAMETRIC INFERENCE

 \hookrightarrow **Definition 3.1** (Point Estimator): Let $X_1, ..., X_n$ a random sample. A *point estimator* $\hat{\theta} :=$ $\hat{\theta}(X_1,...,X_n)$ is an estimator of a parameter θ if it is a statistic.

Example 3.1: Let X be a random variable denoting whether a randomly selected electronic chip is operational or not, i.e. $X = \begin{cases} 1 \text{ operational} \\ 0 \text{ else} \end{cases}$, supposing $X \sim \text{Ber}(\theta)$, then $0 < \theta < 1$ is the probability a randomly selected chip is operational. Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \text{Ber}(\theta)$. Then,

$$\mathcal{F} = \{ \operatorname{Ber}(\theta) : 0 < \theta < 1 \}, \qquad \Theta = (0,1).$$

Then, possible estimators are \overline{X}_n , $\frac{X_1+X_2}{2}$, just X_2 , etc.

 \hookrightarrow **Definition 3.2** (Bias): An estimator $\hat{\theta}_n$ is an *unbiased* estimator of θ if

$$\mathbb{E}_{\theta} \Big[\hat{\theta}_n \Big] = \theta, \qquad \forall \, \theta \in \Theta,$$

where the expected value is taken with respect to the distribution of $\hat{\theta}_n$ (and thus depends on the distribution F_{θ}).

Generally, the *bias* of an estimator $\hat{\theta}_n$ is defined

$$\operatorname{Bias}(\hat{\theta}_n) \coloneqq \mathbb{E}_{\theta}[\hat{\theta}_n] - \theta, \quad \theta \in \Theta.$$

If $\hat{\theta}_n$ unbiased, then $\text{Bias}(\hat{\theta}_n) = 0$.

★ Example 3.2: For instance, recalling the previous example,

$$\mathbb{E}_{\theta}\left[\overline{X}_{n}\right] = \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}_{\theta}[X_{i}] = \frac{1}{n} n\theta = \theta,$$

so \overline{X}_n unbiased. Also,

$$\mathbb{E}_{\theta}[X_1] = \theta,$$

so just X_1 also unbiased, as is $\frac{X_1+X_2}{2}$.

⊗ Example 3.3: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F_\theta$, $\theta = \binom{\mu}{\sigma^2}$, $\mu = \mathbb{E}[X_i]$, $\sigma^2 \operatorname{Var}(X_i)$. Then, $\hat{\mu}_n = \overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ an unbiased estimator of μ . Let $\hat{\sigma}_n^2 = S_n^2 = \frac{1}{n-1} \sum_{i=1}^n \left(X_i - \overline{X}_n \right)^2$, then recalling $\mathbb{E}[\hat{\sigma}_n^2] = \sigma^2$, this is an unbiased estimator of σ^2 . However, changing the constant term, to get

$$\hat{\sigma}_n^{*2} = \frac{1}{n} \sum_{i=1}^n \left(X_i - \overline{X}_n \right)^2,$$

is biased, with

$$\mathbb{E}_{\theta}[\hat{\sigma}_n^{*2}] = \frac{n-1}{n}\sigma^2,$$

so

$$\operatorname{Bias}(\hat{\sigma}_n^{*2}) = -\frac{\sigma^2}{n} < 0,$$

i.e. $\hat{\sigma}_n^{*2}$ underestimates the true parameter on average. Of course, in the limit it becomes 0.

Example 3.4: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \mathcal{U}(0, \theta), \theta > 0, \Theta = (0, \infty)$. Recall $\mathbb{E}_{\theta}[X_i] = \frac{\theta}{2}$. Consider $\hat{\theta}_{n,1} \coloneqq 2X_3, \qquad \hat{\theta}_{n,2} \coloneqq 2\overline{X}_n, \qquad \hat{\theta}_{n,3} \coloneqq X_{(n)}$.

Then, $\mathbb{E}\left[\hat{\theta}_{n,i}\right] = \theta$ for i = 1, 2 and $\frac{n}{n+1}\theta$ for i = 3. Hence, we can scale the last one, $\hat{\theta}_{n,4} := \frac{n+1}{n}\hat{\theta}_{n,3}$, to get an unbiased estimator.

→ Definition 3.3 (Mean-Squared Error): The *Mean-Squared Error* (MSE) of an estimator is defined

$$MSE_{\theta}(\hat{\theta}_{n}) := \mathbb{E}_{\theta} \Big[(\hat{\theta}_{n} - \theta)^{2} \Big]$$

$$= \mathbb{E}_{\theta} \Big[((\hat{\theta}_{n} - \mathbb{E}_{\theta} [\hat{\theta}_{n}]) + (\mathbb{E}_{\theta} [\hat{\theta}_{n}] - \theta))^{2} \Big]$$

$$= Var_{\theta}(\hat{\theta}_{n}) + [Bias(\hat{\theta}_{n})]^{2}.$$

Remark that if $\mathbb{E}_{\theta} [\hat{\theta}_n] = \theta$, i.e. $\hat{\theta}_n$ unbiased, then $MSE_{\theta} (\hat{\theta}_n) = Var_{\theta} (\hat{\theta}_n)$.

Definition 3.4 (Consistency): We say an estimator $\hat{\theta}_n$ of θ is *consistent* if $\hat{\theta}_n \stackrel{P}{\to} \theta$ as $n \to \infty$.

Remark 3.1: There are many ways of establishing consistency; by direct definition of convergence in probability, the WLLN (maybe continuous mapping theorem), or checking if $\mathbb{E}_{\theta}[\hat{\theta}_n] \to \theta$ (if this happens we say $\hat{\theta}_n$ "asymptotically unbiased") and $\mathrm{Var}_{\theta}(\hat{\theta}_n) \to 0$ as $n \to \infty$, for in this case by Chebyshev's Inequality we have consistency.

3 Parametric Inference

- \otimes **Example 3.5**: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} F_{\theta}$.
- 1. $\hat{\mu}_n := \overline{X}_n \xrightarrow{P} \mu$ by WLLN, and $S_n^2 \xrightarrow{P} \sigma^2$ similarly.
- 2. If $X_1,...,X_n \stackrel{\text{iid}}{\sim} \mathcal{U}(0,\theta)$, then $\mathbb{E}[X_i] = \frac{\theta}{2}$. Note that $\hat{\theta}_{n,1} = 2\overline{X}_n$ and $\hat{\theta}_{n,2} = \frac{n+1}{n}X_{(n)}$ are both unbiased estimators of θ , and both are consistent. To see the second one, we have that for any $\varepsilon > 0$,

$$\begin{split} P\big(|X_{(n)} - \theta| > \varepsilon\big) &= P\big(\theta - X_{(n)} > \varepsilon\big) \\ &= P\big(X_{(n)} < \theta - \varepsilon\big) \\ &= \left(\frac{\theta - \varepsilon}{\theta}\right)^n \to 0 \text{ as } n \to \infty. \end{split}$$

We have too that

$$MSE_{\theta}(\hat{\theta}_{n,1}) = Var_{\theta}(\hat{\theta}_{n,1}) = 4Var_{\theta}(\overline{X}_n) = \frac{4}{n} Var(X_i) = \frac{4}{n} \frac{\theta^2}{12} = \frac{\theta^2}{3n}.$$

Also

$$\begin{split} \mathrm{MSE}_{\theta}\Big(\hat{\theta}_{n,2}\Big) &= \mathrm{Var}_{\theta}\Big(\hat{\theta}_{n,2}\Big) = \left(\frac{n+1}{n}\right)^2 \mathrm{Var}\big(X_{(n)}\big) \\ &= \cdots = \frac{\theta^2}{n(n+2)} = \frac{\theta^2}{3n} \cdot \frac{3}{n+2} \leq \mathrm{MSE}_{\theta}\Big(\hat{\theta}_{n,1}\Big) \ \forall \ n \geq 1. \end{split}$$

We will focus on the class of unbiased estimators of a real-valued parameter, $\tau(\theta)$, $\tau:\Theta\to\mathbb{R}$.

§3.1 Uniformly Minimum Variance Unbiased Estimators (UMVUE), Cramér-Rau Lower Bound (CRLB)

Definition 3.5 (UMVUE): Let $X = (X_1, ..., X_n)^t$ be a random variable with a joint pdf/pmf given by

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_1, ..., x_n),$$

where θ some parameter in $\Theta \subseteq \mathbb{R}^d$. An estimator T(X) of a real valued parameter $\tau(\theta)$: $\Theta \to \mathbb{R}$ is said to be a UMVUE of $\tau(\theta)$ if

- 1. $\mathbb{E}_{\theta}[T(X)] = \tau(\theta)$ for every $\theta \in \Theta$;
- 2. for any other unbiased estimator $T^*(X)$ of $\tau(\theta)$, we have

$$\operatorname{Var}_{\theta}(\mathsf{T}(X)) \leq \operatorname{Var}_{\theta}(\mathsf{T}^*(X)), \forall \ \theta \in \Theta.$$

- \hookrightarrow Proposition 3.1 (Cramér-Rau Lower Bound): We define in the case d=1 ($\Theta\subseteq\mathbb{R}$) for convenience. Assume that
- (1) the family $\{p_{\theta}: \theta \in \Theta\}$ has a common support $S = \{x \in \mathbb{R}^n: p_{\theta}(x) > 0\}$ that does not depend on θ ;
 - (2) for $x \in S$, $\theta \in \Theta$, $\frac{d}{d\theta} \log p_{\theta}(x) < \infty$;
 - (3) for any statistic h(x) with $\mathbb{E}_{\theta}[|h(x)|] < \infty$ for every $\theta \in \Theta$, we have

$$\frac{\mathrm{d}}{\mathrm{d}\theta} \int_{S} h(x) p_{\theta}(x) \, \mathrm{d}x = \int_{S} h(x) \frac{\mathrm{d}}{\mathrm{d}\theta} p_{\theta}(x) \, \mathrm{d}x,$$

whenever the right-hand side is finite.

Let T(X) be such that $Var_{\theta}(T(X)) < \infty$ and $\mathbb{E}_{\theta}[T(X)] = \tau(\theta)$ for every every $\theta \in \Theta$. Then if $0 < \mathbb{E}_{\theta}\left[\left(\frac{d}{d\theta}\log(p_{\theta}(x))\right)^2\right] < \infty$ for every $\theta \in \Theta$, then the Cramér-Rao Lower Bound (CRLB) holds:

$$\mathrm{Var}_{\theta}(\mathrm{T}(X)) \geq \frac{\left[\tau'(\theta)\right]^2}{\mathbb{E}_{\theta}\left[\left(\frac{\mathrm{d}}{\mathrm{d}\theta}\log p_{\theta}(x)\right)^2\right]}, \qquad \forall \, \theta \in \Theta.$$

Remark 3.2: The quantity

$$I(\theta) \coloneqq \mathbb{E}_{\theta} \left[\left(\frac{\mathrm{d}}{\mathrm{d}\theta} \log(p_{\theta}(x)) \right)^2 \right]$$

is called the *Fisher information* contained in X about θ .

PROOF. Note that $\tau(\theta) = \mathbb{E}_{\theta}[T(X)]$ implies

$$\tau'(\theta) = \frac{\mathrm{d}}{\mathrm{d}\theta} \mathbb{E}[\mathsf{T}(X)]$$

$$= \frac{\mathrm{d}}{\mathrm{d}\theta} \left[\int_{S} \mathsf{T}(x) p_{\theta}(x) \, \mathrm{d}x \right]$$
by ass. 2, 3
$$= \int_{S} \mathsf{T}(x) \frac{\mathrm{d}}{\mathrm{d}\theta} p_{\theta}(x) \, \mathrm{d}x$$

$$= \int_{S} \mathsf{T}(x) \frac{\mathrm{d}}{\mathrm{d}\theta} [\log p_{\theta}(x)] p_{\theta}(x) \, \mathrm{d}x$$

$$= \mathbb{E}_{\theta} \left[\mathsf{T}(X) \frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(X) \right], \quad \forall \, \theta \in \Theta. \quad (I)$$

On the other hand, by (3) with $h \equiv 1$, then

$$0 = \int_{S} \frac{\mathrm{d}}{\mathrm{d}\theta} p_{\theta}(\mathbf{x}) \, \mathrm{d}\mathbf{x} = \int_{S} \left[\frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(\mathbf{x}) \right] p_{\theta}(\mathbf{x}) \, \mathrm{d}\mathbf{x} \qquad \forall \, \theta \in \Theta$$

$$\Rightarrow \mathbb{E}_{\theta} \left[\frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(\mathbf{X}) \right] = 0. \quad \text{(II)}$$

Combining (I) and (II),

$$\tau'(\theta) = \operatorname{Cov}_{\theta}\left(\mathrm{T}(X), \frac{\mathrm{d}}{\mathrm{d}\theta}\log p_{\theta}(X)\right),$$

since $Cov(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$, but the second of these terms vanishes by (II). Thus,

$$\left[\tau'(\theta)^2\right] = \mathrm{Cov}_{\theta}^2\left(\mathrm{T}(\boldsymbol{X}), \frac{\mathrm{d}}{\mathrm{d}\theta}\log p_{\theta}(\boldsymbol{X})\right).$$

By Cauchy-Schwarz, we find

$$\begin{split} \left[\tau'(\theta)\right]^2 &\leq \mathrm{Var}_{\theta}(\mathrm{T}(\boldsymbol{X})) \mathrm{Var}_{\theta}\left(\frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(\boldsymbol{X})\right) \\ &\leq \mathrm{Var}_{\theta}(\mathrm{T}(\boldsymbol{X})) \mathbb{E}_{\theta}\left\{\left[\frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(\boldsymbol{X})\right]^2\right\}, \end{split}$$

the last line following by the Bartlett Identity.

Remark 3.3: If $X_1, ..., X_n \stackrel{\text{iid}}{\sim} f_{\theta}$, then $p_{\theta}(x) = \prod_{i=1}^n f(x_i; \theta)$, and

$$I(\theta) = \mathbb{E}_{\theta} \left\{ \left[\frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(X) \right]^{2} \right\} = \mathbb{E}_{\theta} \left\{ \left[\sum_{i=1}^{n} \frac{\mathrm{d}}{\mathrm{d}\theta} \log f(X_{i}; \theta) \right]^{2} \right\}$$
$$= n \mathbb{E}_{\theta} \left\{ \left(\frac{\mathrm{d}}{\mathrm{d}\theta} \log f(X_{1}; \theta) \right)^{2} \right\},$$
$$= I_{1}(\theta)$$

so the CRLB in this case reads

$$\operatorname{Var}_{\theta}(\mathrm{T}(X)) \ge \frac{\left[\tau'(\theta)\right]^2}{nI_1(\theta)},$$

and moreover if $\tau(\theta) = \theta$ itself,

$$\operatorname{Var}_{\theta}(\mathbf{T}(X)) \ge \frac{1}{nI_1(\theta)}.$$

Example 3.6: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \text{Ber}(\theta)$, so $f(x; \theta) = \theta^x (1 - \theta)^{1 - x}$ for x = 0, 1. Then, $\log(f(x; \theta)) = x \log(\theta) + (1 - x) \log(1 - \theta)$

so

$$\frac{\mathrm{d}}{\mathrm{d}\theta}\log(f(x;\theta)) = \frac{x}{\theta} - \frac{1-x}{1-\theta'}$$

so the Fisher information in one X_1 is given

$$I_1(\theta) = \mathbb{E}_{\theta} \left\{ \left(\frac{X}{\theta} - \frac{1-X}{1-\theta} \right)^2 \right\} = \frac{1}{\theta(1-\theta)}.$$

For any unbiased estimator of $\tau(\theta) = \theta$, the CRLB gives

$$\operatorname{Var}_{\theta}(\mathrm{T}(X)) \ge \frac{1}{nI_1(\theta)} = \frac{\theta(1-\theta)}{n}.$$

Recall our estimator $\hat{\theta}_n = \overline{X}_n$. We have that $\operatorname{Var}_{\theta}(\overline{X}_n) = \frac{1}{n}\operatorname{Var}_{\theta}(X_1) = \frac{\theta(1-\theta)}{n}$.

Remark 3.4: If p_{θ} additionally twice differentiable in θ and $\mathbb{E}_{\theta} \left\{ \frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(X) \right\}$ is also differentiable under the \mathbb{E}_{θ} ,

$$\frac{\mathrm{d}}{\mathrm{d}\theta}\log p_{\theta}(X) = \int \frac{\mathrm{d}}{\mathrm{d}\theta} \left\{ \left[\frac{\mathrm{d}}{\mathrm{d}\theta}\log p_{\theta}(x) \right] p_{\theta}(x) \right\} \mathrm{d}x.$$

In particular, this implies $\int p''_{\theta}(x) dx = 0$. Then,

$$I(\theta) = \mathbb{E}_{\theta} \left\{ \left[\frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(X) \right]^{2} \right\} = -\mathbb{E}_{\theta} \left\{ \frac{\mathrm{d}^{2}}{\mathrm{d}\theta^{2}} p_{\theta}(X) \right\},$$

making it easier to compute $I(\theta)$. This follows from the fact that

$$\frac{\mathrm{d}^2}{\mathrm{d}\theta^2}\log p_\theta(x) = \frac{{p_\theta}''(x)}{p_\theta(x)} - \left[\frac{\mathrm{d}}{\mathrm{d}\theta}\log p_\theta(x)\right]^2,$$

and so taking the expected value of both sides cancels the inner-most term by the differentiability condition of p_{θ} ;

$$\mathbb{E}\left[\frac{\mathrm{d}^2}{\mathrm{d}\theta^2}\log p_{\theta}(x)\right] = \mathbb{E}\left[\frac{p_{\theta}''(x)}{p_{\theta}(x)}\right] - \mathbb{E}\left[\left[\frac{\mathrm{d}}{\mathrm{d}\theta}\log p_{\theta}(x)\right]^2\right]$$
$$= \int p_{\theta}''(x)\,\mathrm{d}x - I(\theta).$$

★ Example 3.7: Returning to the previous example, remark that

$$\frac{\mathrm{d}^2}{\mathrm{d}\theta^2}\log(f(x;\theta)) = -\frac{x}{\theta^2} - \frac{x-1}{(1-\theta)^2},$$

and so

$$\mathbb{E}\left[\frac{\mathrm{d}^2}{\mathrm{d}\theta^2}\log f(x;\theta)\right] = \frac{1}{\theta} + \frac{1}{1-\theta}$$

so $I_1(\theta) = \frac{1}{\theta(1-\theta)}$ as we found before.

Remark 3.5: The CRLB is *not* a sharp bound, in the sense that the UMVUE for a particular parameter may be strictly larger than the CRLB.

Example 3.8: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \theta^2)$. Then, $\hat{\mu}_n$ the UMVUE for μ . If μ known, then $\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ is the UMVUE for σ^2 . If μ is unknown, then $\frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X}_n)^2$ would be the UMVUE for σ^2 .

However, if $X_i \stackrel{\text{iid}}{\sim} \exp(\beta)$, with $f(x; \beta) = \frac{1}{\beta} e^{-\frac{x}{\beta}}$ for x > 0, S_n^2 is not the UMVUE for $\operatorname{Var}_{\beta}(X_i) = \beta^2$.

Theorem 3.1 (Attaining the CRLB): Suppose $X = (X_1, ..., X_n) \sim p_\theta$. Let T(X) be unbiased for $\tau(\theta)$. Then, T(X) attains the CRLB if and only if

$$a(\theta)\{T(x) - \tau(\theta)\} = \frac{\mathrm{d}}{\mathrm{d}\theta} \log p(x;\theta),$$

for some function $a(\theta)$, for every $\theta \in \Theta$ and x in the support of p.

PROOF. In the proof of the CRLB, the only inequality arose from using Cauchy-Schwarz with bounding the covariance of T(X) and $\frac{d}{d\theta} \log p_{\theta}(X)$. Equality in this inequality holds if and only if the terms are linearly dependent, namely if there is some function $a(\theta)$ and $b(\theta)$ such that $a(\theta)T(x) + b(\theta) = \frac{d}{d\theta} \log p_{\theta}(x)$.

On the other hand,

$$\mathbb{E}_{\theta}\{a(\theta)T(\boldsymbol{X}) + b(\theta)\} = \mathbb{E}_{\theta}\left\{\frac{\mathrm{d}}{\mathrm{d}\theta}\log p_{\theta}(x)\right\} = 0 \Rightarrow b(\theta) = -\mathbb{E}_{\theta}\{a(\theta)T(\boldsymbol{X})\} = -a(\theta)\tau(\theta),$$

so combining these two gives the desired linear relation.

Example 3.9 (Exponential family): $X_i \stackrel{\text{iid}}{\sim} f(x;\theta) = h(x)c(\theta) \exp\{\omega(\theta)T_1(x)\}$, where h a nonnegative function of only x and x a nonnegative function of only x, with the support of x being independent of x. Then

$$p_{\theta}(x) = \prod_{i=1}^{n} f(x_i; \theta) = \left[\prod_{i=1}^{n} h(x_i) \right] (c(\theta))^n \exp\left(\omega(\theta) \sum_{i=1}^{n} T_1(x_i)\right).$$

Taking the log:

$$\begin{split} \frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(x) &= n \frac{c'(\theta)}{c(\theta)} + \omega'(\theta) \sum_{i=1}^{n} T_{1}(x_{i}) \\ &= \omega'(\theta) \left\{ \sum_{i=1}^{n} T_{1}(x_{i}) - \frac{-nc'(\theta)}{c(\theta)\omega'(\theta)} \right\}. \end{split}$$

Let

$$\tau(\theta) = -\frac{c'(\theta)}{c(\theta)\omega'(\theta)}.$$

Then, since

$$\mathbb{E}_{\theta} \left[\frac{\mathrm{d}}{\mathrm{d}\theta} \log p_{\theta}(x) \right] = 0,$$

then

$$\mathbb{E}_{\theta} \left[\sum_{i=1}^{n} T_1(X_i) \right] = n\tau(\theta),$$

so

$$T(X) = \frac{1}{n} \sum_{i=1}^{n} T_1(X_i)$$

is a UMVUE for $\tau(\theta)$ by the previous theorem.

Example 3.10: Let $X_i \stackrel{\text{iid}}{\sim} \text{Poisson}(\theta)$ so

$$f(x;\theta) = \frac{e^{-\theta}}{x!} \theta^x = \frac{e^{-\theta}}{x!} e^{x \log(\theta)},$$

with support $x \in \{0, 1, ...\}$. Then, we notice that with

$$h(x) = \frac{1}{x!}, c(\theta) = e^{-\theta}, \omega(\theta) = \log(\theta), T_1(x) = x,$$

that X_i in the exponential family. Then, according to the previous example,

$$\tau(\theta) = -\frac{-e^{-\theta}}{e^{-\theta}\frac{1}{\theta}} = \theta,$$

has UMVUE

$$T(\boldsymbol{X}) = \frac{1}{n} \sum_{i=1}^{n} X_i = \overline{X}_n.$$

Example 3.11: Recall we found, for $X_i \stackrel{\text{iid}}{\sim} \mathcal{U}(0,\theta)$, that $\hat{\theta}_n := \frac{n+1}{n} X_{(n)}$ was an unbiased estimator but cannot obtain the CRLB since the regularity conditions are not satisfied (namely, the support of the pdfs depends on the parameter). Moreover, we found

$$\mathbb{E}_{\theta}\left\{\frac{n+1}{n}X_{(n)}\right\} = \theta, \operatorname{Var}_{\theta}\left\{\frac{n+1}{n}X_{(n)}\right\} = \frac{\theta^2}{n(n+2)}.$$

If we temporarily ignore that we cannot apply CRLB, we would find

$$CRLB = \frac{1}{nI_1(\theta)} = \frac{\theta^2}{n},$$

so our estimator actually has a "better" variance. We'll see later that this estimator actually the UMVUE.

§3.2 Sufficiency

We can't always find unbiased estimators; here we look for other ways for comparing different estimators.

3.2 Sufficiency 23

Example 3.12: Let $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, and consider the following estimators of σ^2 :

$$S_1^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X}_n)^2,$$

$$S_2^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{X}_n)^2,$$

$$S_3^2 = \frac{1}{n+1} \sum_{i=1}^n (X_i - \overline{X}_n)^2.$$

One verifies these have respective means, variances

$$\begin{array}{c|ccccc} & S_1^2 & S_2^2 & S_3^2 \\ \hline \mathbb{E} & \frac{n-1}{n}\sigma^2 & \sigma^2 & \frac{n-1}{n+1}\sigma^2 \\ \text{Var} & \frac{2(n-1)\sigma^4}{n^2} & \frac{2\sigma^4}{n-1} & \frac{2(n-1)}{(n+1)^2}\sigma^4 \end{array}$$

. We notice then that

$$MSE(S_3^2) < MSE(S_2^2) < MSE(S_1^2),$$

so despite the fact that S_2^2 is unbiased, it does not minimize the MSE.

Definition 3.6 (Sufficiency): Suppose $X = (X_1, ..., X_n)$ has joint pdf (pmf) $p(x; \theta)$ for $\theta \in \Theta$. A statistic $T(X) : \mathbb{R}^n \supseteq X \to S_T \subseteq \mathbb{R}^k$, $k \le n$, is *sufficient* for θ or the parametric family $\{p_\theta : \theta \in \Theta\}$ if the conditional distribution of $(X_1, ..., X_n)$ given T(X) = t for any $\theta \in \Theta$ and $t \in S_T$ in the support such that $P_\theta(t \in S_T) = 1$, does not depend on θ . Namely,

$$f_{X|T(X)=t}(x_1,...,x_n),$$

does *not* depend on θ .

Example 3.13: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \text{Ber}(\theta)$. Let $T(X) = \sum_{i=1}^n X_i$. We know that then $T(X) \sim \text{Bin}(n, \theta)$. We claim T sufficient; we have

$$f_{\theta}(x_1,...,x_n \mid T(X) = t) = \begin{cases} \frac{1}{\binom{n}{t}} & \text{if } \sum_{i=1}^n x_i = t\\ 0 & \text{else} \end{cases}$$

which is independent of θ so indeed sufficient.

Remark 3.6: A sufficient statistic induces a partitioning of the sample space $X \subseteq \mathbb{R}^n$; namely,

$$X = \bigcup_{t \in S_T} \Pi_t,$$

such that

$$\Pi_t = \{ x = (x_1, ..., x_n) \in X \mid T(x) = t \},$$

and S_T the support of T.

Example 3.14: Return to the Bernoulli example from before, and consider specifically the case when n = 2, so $T(X) = X_1 + X_2$ is a sufficient statistic as we showed. Then, the sample space is given by

$$X = \{(0,0), (0,1), (1,0), (1,1)\},$$

and T has support

$$T(x) = x_1 + x_2 \in \{0, 1, 2\} =: S_T.$$

This induces the partitioning

$$X = \Pi_0 \sqcup \Pi_1 \sqcup \Pi_2 = \{(0,0)\} \sqcup \{(0,1),(1,0)\} \sqcup \{(1,1)\}.$$

Theorem 3.2 (Neyman-Fisher Factorization Theorem): Let $X = (X_1, ..., X_n)^t$ be a random vector with a joint pdf/pmf $p_\theta(x) = p(x; \theta)$. A statistic T(X) is sufficient for θ if and only if there exist functions $g(\cdot; \theta)$ and $h(\cdot)$ such that

$$p_{\theta}(\mathbf{x}) = h(\mathbf{x}) \cdot g(\theta, T(\mathbf{x})),$$

for every $\theta \in \Theta$ and $x \in X$.

Note that g depends on x only through T(x), and h does not depend on θ .

Proof. We prove in the discrete case.

Note that

$$f_{X|T(X)=t_x}(x) = \frac{P_{\theta}(X_1 = x_1, ..., X_n = x_n, T(X) = t_x)}{P_{\theta}(T(X) = t_x)},$$

for every x such that $T(x) = t_x$, and 0 otherwise;

$$= \frac{P_{\theta}(X_1 = x_1, ..., X_n = x_n)}{\sum_{y = (y_1, ..., y_n): T(y) = t_x} P(X_1 = y_1, ..., X_n = y_n)}.$$

If T(X) a sufficient statistic for θ , then the above ratio, by definition, does not depend on θ ; hence, putting h(x) to be the ratio above, it is independent of θ (is only a function of the data), and if we take g to be the denominator of the ratio above, then g depends on the data only through T. Hence, we can write $p_{\theta}(x) = h(x) \cdot g(t_x; \theta)$.

3.2 Sufficiency 25

Conversely, suppose $p_{\theta}(x) = g(T(x); \theta)h(x)$. Then,

$$f_{X|T(X)=t_x}(x;\theta) = \frac{g(t_x;\theta)h(x)}{\sum_{y:T(y)=t_x}g(T(y);\theta)h(y)} = \frac{h(x)}{\sum_{y:T(y)=t_x}h(y)},$$

which depends only on x and hence T(X) a sufficient statistic.

Example 3.15: Let again $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \text{Ber}(\theta)$ so

$$p_{\theta}(x_1, ..., x_n) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1 - x_i} = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i} \prod_{i=1}^n \mathbb{1}\{x_i \in \{0, 1\}\}.$$

for $x_i = 0, 1$.

One notices that the LHS (not the product) can be written as a function of θ and $\sum_{i=1}^{n} x_i$ only, and the remaining term is independent of θ . Hence by the previous theorem $T(X) = \sum_{i=1}^{n} X_i$ a sufficient statistic for θ .

Example 3.16: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \mathcal{U}(0, \theta)$, so $f(x; \theta) = \begin{cases} \frac{1}{\theta} & \text{if } 0 < x < \theta \\ 0 & \text{else} \end{cases}$. Then

$$\begin{split} p_{\theta}(x) &= \prod_{i=1}^{n} \frac{1}{\theta} \mathbb{1}(0 < x_i < \theta) \\ &= \underbrace{\frac{1}{\theta^n} \mathbb{1}\big(0 < x_{(n)} < \theta\big)}_{=:g(T(x;\theta))} \underbrace{\mathbb{1}\big(0 < x_{(1)} < \theta\big)}_{=:h(x)}, \end{split}$$

so $X_{(n)}$ is a sufficient statistic for θ .

Remark 3.7: If T is a sufficient statistic for θ and $T(X) = \Phi(T^*(X))$ where Φ is a measurable function and T^* another statistic, then T^* is also a sufficient statistic.

- **Example 3.17**: In the exponential family, we claim $T(X_1,...,X_n) = \sum_{i=1}^n T_1(X_i)$.
- **Example 3.18**: Let $X_1, ..., X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$ and $\theta = (\mu, \sigma^2)$ both unknown. Using the factorization theorem, we can see that

$$T(\mathbf{X}) = \left(\sum_{i=1}^{n} X_i, \sum_{i=1}^{n} X_i^2\right)$$

is a sufficient statistic for θ , as is (\overline{X}_n, S_n^2) .

Remark 3.8: This does *not* imply that say $\sum_{i=1}^{n} X_i$ sufficient for μ ! Namely T is a sufficient statistic for the 2-dimensional parameter θ . We cannot simply separate the dependence.

⊛ Example 3.19: Recall the Bernoulli example once again. We claim that

$$T_m^*(X) = \left(\sum_{i=1}^m X_i, \sum_{i=m+1}^n X_i\right), \qquad 1 \le m \le n-1)$$

is also sufficient for $0 < \theta < 1$. Clearly this is no different then just using the one-dimensional statistic $\sum_{i=1}^{n} X_i$; we'd like to formalize how to differentiate such statistics. Namely, $\sum_{i=1}^{n} X_i$ is called a *minimal* sufficient statistic for θ .

 \hookrightarrow **Definition 3.7** (Minimal Sufficient Statistic): A statistic T(X) is a *minimal sufficient statistic* for θ iff

- T(X) is sufficient;
- For any other sufficient statistic $T^*(X)$ of θ , T(X) is a function of $T^*(X)$, i.e.

$$T(\boldsymbol{X}) = \varphi(T^*(\boldsymbol{X})),$$

where $\varphi(\cdot)$ some measurable function, or equivalently, $\forall x, y \in X \subseteq \mathbb{R}^n$, if $T^*(x) = T^*(y)$ then T(x) = T(y).

Remark 3.9: If T(X) minimally sufficient and induces a partitioning

$$X = \bigcup_{t \in S_T} \Pi_t, \qquad \Pi_t \coloneqq \{x \in X : T(x) = t\}$$

and $T^*(X)$ any sufficient statistic that induces a partitioning

$$X = \bigcup_{t^* \in S^*_{T^*}} \Pi^*_{t^*}, \qquad \Pi^*_{t^*} \coloneqq \{x \in X : T^*(x) = t^*\},$$

then we find that $\forall t^* \in S_{T^*}^*$, there is some $t \in S_T$ such that $\Pi_{t^*}^* \subseteq \Pi_t$; namely, the partition induced by T(X) is the *coarsest* possible partition of X.

Theorem 3.3 (Lehmann-Scheffé): For a parametric family $p_{\theta}(\cdot)$ (the joint pdf/pmf of X), suppose a statistic $T(X) = T(X_1, ..., X_n)$ is such that for every $x, y \in X \subseteq \mathbb{R}^n$ $T(x) = T(y) \Leftrightarrow \frac{p_{\theta}(x)}{p_{\theta}(y)}$ does not depend on θ . Then, T(X) is a minimal sufficient statistic for θ .

3.2 Sufficiency 27

Example 3.20: Suppose $X_i \stackrel{\text{iid}}{\sim} \mathcal{U}(0,\theta)$, then $p_{\theta}(x) = \frac{1}{\theta^n} \mathbb{1}\{x_{(n)} < \theta\} \mathbb{1}\{x_{(1)} > 0\}$; then $T(X) = X_{(n)}$ is a sufficient statistic for θ . For any $x, y \in X$, we find

$$\frac{p_{\theta}(x)}{p_{\theta}(y)} = \frac{\mathbb{1}\{x_{(n)} < \theta\} \mathbb{1}\{x_{(1)} > 0\}}{\mathbb{1}\{y_{(n)} < \theta\} \mathbb{1}\{y_{(1)} > 0\}},$$

which does not depend on θ iff $x_{(n)} = y_{(n)}$ iff T(x) = T(y) and therefore by the previous theorem T(X) is a minimally sufficient statistic.

Example 3.21: If $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$ and $\theta = (\mu, \sigma^2)$, it can be shown that

$$T(X) = \left(\sum_{i=1}^{n} X_i, \sum_{i=1}^{n} X_i^2\right)$$

is a minimal sufficient statistic for θ . Any one-to-one function of a minimally sufficient statistic also minimally sufficient, hence this implies (\overline{X}_n, S_n^2) is also minimally sufficient for θ .

§3.3 Completeness

Definition 3.8 (Completeness): Let *X* be a random variable with a pmf/pdf belonging to a parametric family $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$. This family is said to be *complete* if for any measurable function *g* with $\mathbb{E}_{\theta}[g(X)] < \infty$, then $\mathbb{E}_{\theta}[g(X)] = 0$ for all $\theta \in \Theta$ implies $P_{\theta}(g(X) = 0) = 1$.

A statistic $T(X) = T(X_1, ..., X_n)$ is said to be *complete* if the family of its distributions is complete.

Remark 3.10: Complete and sufficient ⇒ minimal, but minimally sufficient may not be complete, as we'll see.

Example 3.22: Let $X_i \stackrel{\text{iid}}{\sim} \text{Ber}(\theta)$, then note $T(X) = \sum_{i=1}^n X_i \sim \text{Bin}(n, \theta)$. Let g a measurable function. Then,

$$\begin{split} 0 &= \mathbb{E}_{\theta}[g(X)] \Rightarrow 0 = \sum_{t=0}^{n} g(t) \binom{n}{t} \theta^{t} (1-\theta)^{n-t} \\ &= (1-\theta)^{n} \sum_{t=0}^{n} g(t) \binom{n}{t} \left(\frac{\frac{-:\eta}{\theta}}{1-\theta}\right)^{t} \\ &= \sum_{t=0}^{n} g(t) \binom{n}{t} \eta^{t}. \end{split}$$

Then, this is just a polynomial in η , which, being equal to zero implies all the coefficients $g(t)\binom{n}{t}=0$ for every t and hence g(t)=0. Hence, T(X) is a complete statistic.

Example 3.23: If $X \sim \mathcal{N}(0, \theta)$, the family is not complete. For instance with g(x) := x, $\mathbb{E}_{\theta}(X) = 0$ but g(x) is not identically zero. On the other hand, $T(X) = X^2$ is a complete statistic. To see this, we know $\frac{X^2}{\theta} \sim \chi^2_{(1)}$, so

$$\begin{split} \mathbb{E}_{\theta}\big(g(T)\big) &= 0 \Rightarrow 0 = \int_0^\infty g(t) f_T(t;\theta) \, \mathrm{d}t \\ &= \int_0^\infty \frac{1}{\sqrt{2\pi\theta}} g(t) t^{-\frac{1}{2}} e^{-\frac{t}{2\theta}} \, \mathrm{d}t \\ &= \mathcal{L}\bigg\{g(t) t^{-\frac{1}{2}} \frac{1}{\sqrt{2\pi\theta}}\bigg\}. \end{split}$$

By uniqueness of the Laplace transform, it must be that $g(t)t^{-\frac{1}{2}} \equiv 0$ hence g(t) = 0 and thus $T(X) = X^2$ is a complete statistic.

Example 3.24: In the exponential family, $\sum_{i=1}^{n} T_1(X_i)$ is a complete statistic.

Note that an unbiased estimator of a parameter of interest may not even exist. For instance,

Example 3.25: If $X \sim \text{Bin}(n, \theta)$, let $\tau(\theta) = \frac{1}{\theta}$. If $\delta(X)$ is an unbiased estimator of $\tau(\theta)$, we must have $\mathbb{E}_{\theta}[\delta(X)] = \frac{1}{\theta}$ i.e.

$$\sum_{x=0}^{n} \delta(x) \binom{n}{x} \theta^{x} (1-\theta)^{n-x} = \frac{1}{\theta}.$$

As $\theta \to 0$, the left-hand side will just be $\delta(0)$, while the right-hand side will diverge to ∞ , so no such estimator exists.

Theorem 3.4 (Rao-Blackwell): Let U(X) be an unbiased estimator of $\tau(\theta)$ and let T(X) be a sufficient statistic for the parametric family. Set

$$\delta(t) = \mathbb{E}_{\theta}[U(X) \mid T(X) = t], \quad t \in S_T.$$

Then,

- $\delta(T(X))$ is a statistic, i.e. only depends on X;
- $\mathbb{E}_{\theta}[\delta(T(X))] = \tau(\theta);$
- $\operatorname{Var}_{\theta}(\delta(T(X))) \leq \operatorname{Var}_{\theta}[U(X)].$

Proof.

- $\delta(T(X)) = \mathbb{E}_{\theta}[U(X)|T(X)]$ is a random variable in its own right, and is a statistic because T(X) is sufficient, hence conditioning on T(X) will result in no reliance on θ .
- $\mathbb{E}_{\theta}[\delta(T(X))] = \mathbb{E}_{\theta}[\mathbb{E}_{\theta}[U(X)|T(X)]] = \mathbb{E}_{\theta}[U(X)] = \tau(\theta)$ (using the law of total expectation), since U(X) is an unbiased estimator of $\tau(\theta)$.
- Using the law of total variance, we find

$$\operatorname{Var}_{\theta}(U(X)) = \operatorname{Var}_{\theta}(\underbrace{\mathbb{E}_{\theta}[U(X)|T(X)]}_{=\delta(T(X))}) + \mathbb{E}_{\theta}[\operatorname{Var}_{\theta}(U(X)|T(X))]$$

$$= \operatorname{Var}_{\theta}[\delta(T(X))] + \mathbb{E}_{\theta}[\underbrace{\operatorname{Var}_{\theta}(U(X)|T(X))}_{\geq 0}]$$

$$\geq \operatorname{Var}_{\theta}[\delta(T(X))].$$

Remark 3.11: This theorem gives a systematic manner of improving unbiased estimators, by taking an unbiased estimator and a sufficient statistic, and "Rao-Blackwell-izing", leading to a uniform improvement in variance.

Theorem 3.5 (Lehmann-Scheffé: Uniqueness): Let T(X) be a complete sufficient statistic. Let U(X) = h(T(X)), for a measurable function h, an unbiased estimator of $\tau(\theta)$ such that $\mathbb{E}_{\theta} \left[U(X)^2 \right] < \infty$. Then, U(X) is the unique unbiased estimator of $\tau(\theta)$ with the smallest variance in the class of unbiased estimators of $\tau(\theta)$.

PROOF. By the Rao-Blackwell Theorem, it suffices to restrict attention to unbiased estimators that are only functions of T(X); for any other such unbiased statistic, applying Rao-Blackwell to it results in a new statistic with smaller variance.

Now, let $V(X) = h^*(T(X))$ be any other unbiased estimator of $\tau(\theta)$. Then,

$$\mathbb{E}_{\theta}[V(X)] = \mathbb{E}_{\theta}[U(X)] = \tau(\theta)$$

hence

$$\mathbb{E}_{\theta}[V(\boldsymbol{X}) - U(\boldsymbol{X})] = \mathbb{E}_{\theta}\big[h^*(T(\boldsymbol{X})) - h(T(\boldsymbol{X}))\big] = 0.$$

Let $g(T(X)) = h^*(T(X)) - h(T(X))$; then, since T(X) complete, it must be that P(g = 0) = 1 i.e.

$$P(h(T(X)) = h^*(T(X))) = 1,$$

so U(X), V(X) are almost surely identical, hence we indeed have uniqueness.

Remark 3.12: This, combined with the Rao-Blackwell theorem, provides a method for obtaining the UMVUE for $\tau(\theta)$ starting with a complete sufficient statistic and an unbiased statistic.

Example 3.26: Let $X_i \stackrel{\text{iid}}{\sim} \text{Ber}(\theta)$, i = 1, ..., n and $\hat{\theta}_n = \overline{X}_n$. This is unbiased, and $\sum_{i=1}^n X_i$ is a complete and sufficient statistic. Hence, $\hat{\theta}_n$ is a unbiased estimator that is a function of a complete and sufficient statistic and thus is the UMVUE for θ by the Lehmann-Scheffé Theorem.

Example 3.27: Let $X_i \stackrel{\text{iid}}{\sim} \operatorname{Pos}(\theta)$, i = 1, ..., n and $\hat{\theta}_n = \overline{X}_n$. This is unbiased, and again $\sum_{i=1}^n X_i$ is a complete sufficient statistic hence $\hat{\theta}_n$ is the UMVUE of θ .

Suppose now $\tau(\theta) = P_{\theta}(X = 0) = e^{-\theta}$; can we obtain a UMVUE for this (function of) a parameter? Define

$$U(X_1) = \mathbb{1}\{X_1 = 0\},\$$

which will be unbiased for $\tau(\theta)$. We already have a complete and sufficient statistic. Applying now the Rao-Blackwell theorem, we obtain

$$\delta(t) = \mathbb{E}_{\theta} \left[U(X_1) \mid \sum_{j=1}^{n} X_j = t \right].$$

One verifies that

$$\left(X_i \mid \sum_{j=1}^n X_j = t\right) \sim \text{Bin}\left(t, \frac{1}{n}\right),\,$$

therefore

$$\delta(t) = P_{\theta}(X_1 = 0 \mid T(X) = t) = \left(1 - \frac{1}{n}\right)^t.$$

So, $\delta(T(X)) = \left(1 - \frac{1}{n}\right)^{\sum_{i=1}^{n} X_i}$ is the UMVUE of $e^{-\theta}$. Remark that

$$\delta(T(X)) = \left(1 - \frac{1}{n}\right)^{nX_n} \approx e^{-\overline{X}_n} \text{ for large } n.$$

Example 3.28: Let $X_i \stackrel{\text{iid}}{\sim} \text{Ber}(\theta)$, i = 1, ..., n, and suppose $\tau(\theta) = \text{Var}(X_i) = \theta(1 - \theta)$. Recall the UMVUE for θ is $\hat{\theta}_n$. Note that

$$T(\boldsymbol{X}) = \sum_{i=1}^{n} X_i \sim \text{Bin}(n, \theta),$$

is complete and sufficient. We know $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n \left(X_i - \overline{X}_n \right)^2 = U(X)$ is unbiased for $\tau(\theta)$. We may write

$$U(X) = \frac{1}{n-1} \left[\sum_{i=1}^{n} X_i^2 - n \overline{X}_n^2 \right]$$
since $X_i \in \{0, 1\}$

$$= \frac{1}{n-1} \left[\sum_{i=1}^{n} X_i - n \overline{X}_n^2 \right]$$

$$= \frac{1}{n-1} \left(T(X) - \frac{T^2(X)}{n} \right)$$

$$= \frac{n}{n-1} \overline{X}_n \left(1 - \overline{X}_n \right)$$

Hence, U(X) a function of T(X), a complete sufficient statistic, and U(X) is unbiased, so we conclude U(X) the UMVUE for $\tau(\theta)$.

§3.4 Existence of a UMVUE

Definition 3.9 (Unbiased Estimators of Zero): An estimator $\delta(X)$ satisfying $\mathbb{E}_{\theta}[\delta(X)] = 0$ is called an *unbiased estimator of zero*.

3.4 Existence of a UMVUE