

# MATH324 (Statistics) – Lecture Notes

McGill University

Prof. Masoud Asgharian

Sam K.H.Targhi-Dunn

sam.targhi@mail.mcgill.ca

Winter 2019

## Contents

<b>1</b>	<b>Lecture 1</b>	<b>4</b>
1.1	Overview . . . . .	5
1.2	Parametric Models . . . . .	5
1.3	Nonparametric Models . . . . .	5
1.4	Point Estimator . . . . .	6
1.5	Estimation Error . . . . .	6
<b>2</b>	<b>Lecture 2</b>	<b>7</b>
2.1	Markov's Inequality . . . . .	7
2.2	Tchebyshev's Inequality . . . . .	7
2.3	Application to Voting . . . . .	10
<b>3</b>	<b>Lecture 3</b>	<b>12</b>
3.1	MSE . . . . .	12
3.2	Unbiased Estimators . . . . .	13

3.3	Stein's Paradox . . . . .	14
3.4	Admissibility . . . . .	15
<b>4</b>	<b>Lecture 4</b>	<b>16</b>
4.1	Estimating Variance . . . . .	18
<b>5</b>	<b>Lecture 5 : Confidence Intervals</b>	<b>22</b>
5.1	Confidence Intervals . . . . .	22
5.2	Large Sample Confidence Interval . . . . .	27
5.3	Small Sample Confidence Intervals . . . . .	30
5.4	Pivotal Quantity and Probability Integral Transform . . . . .	34
<b>6</b>	<b>Lecture 6</b>	<b>36</b>
6.1	Small Sample Confidence Interval(general case): . . . . .	36
6.2	Probability Integral Transform(PIT) . . . . .	36
6.3	Pivotal Quantity . . . . .	38
6.4	Small Size Determination . . . . .	40
6.5	Sample Size Determination For Other Parameters . . . . .	41
6.5.1	Sample Size Determination (Small Sample) . . . . .	42
6.5.2	Sample Size Determination(Two Sample Case) . . . . .	43
<b>7</b>	<b>Lecture 7</b>	<b>45</b>
7.1	Chapter 9 - Relative Efficiency . . . . .	45
7.2	Consistency . . . . .	48
<b>8</b>	<b>Lecture 8</b>	<b>51</b>
8.1	Consistency . . . . .	51
8.2	Markov's Inequality . . . . .	51
8.3	Kolmogorov's Law of Large Numbers(LLN) . . . . .	54
8.4	Sufficiency . . . . .	57
<b>9</b>	<b>Lecture 9</b>	<b>59</b>
9.1	Sufficiency . . . . .	59

9.2 Likelihood . . . . .	63
<b>10 Lecture 10</b>	<b>66</b>
10.1 The Rao-Blackwell Theorem . . . . .	66
<b>11 Lecture 11</b>	<b>72</b>
11.1 Method of Maximum Likelihood (ML) . . . . .	72
11.2 Method of Moments . . . . .	76
<b>12 Lecture 12</b>	<b>78</b>
12.1 Section 9.8 - Large Sample Property of the MLEs . . . . .	78

# 1 Lecture 1

## Overview of statistics

### Point estimation:

Statistic and estimator + examples

### Bias and Mean Square Error

Unbiasedness, Bias,  $MSE(\hat{\theta})$ , Decomposition of MSE (#8.8 ( $\hat{\theta}_1, \hat{\theta}_5$ ), page 3294, #8.6, , page 394)

### Common Unbiased Estimators

$\mu, p, \mu_1 - \mu_2, p_1 - p_2$ , &  $\sigma^2 (S_{n-1}^2 \rightarrow S^2$  in the textbook) **Confidence Interval**

### Pivotal Quantities

#MISSING Graph - lecture 1 - p0

#### Pivotal

$$X_i \sim F_{\theta}(x) \implies Y_i = F_{\theta}^{-1}(X_i) \sim \text{Unif}(0, 1) \implies -\log Y_i \sim \text{Exp}(1)$$

$$\implies \sum_{i=1}^n -\log F_{\theta}^{-1}(X_i) = \sum_{i=1}^n Y_i \sim G(n, 1)$$

Example :  $X_i \sim \text{Exp}(\lambda)$

### Large n:

$$\frac{\hat{\theta}_n - \theta}{\sigma_{\hat{\theta}}} \sim N(0, 1) \text{ for large } n \quad \text{Example : } X_i \sim N(\mu, 1)$$

### Sample Size Determination:

Use the notes for 203.

## 1.1 Overview

#MISSING GRAPH LECTURE 1 PAGE 1

In MATH324 we cover statistical Inference (Theory of Point Estimator (TPE) , Testing Statistical Hypothesis (TSH) , Confidence Interval (C.I) and Data Generating Process (DGP).

## 1.2 Parametric Models

Model is known up to finitely many unknown parameters.

E.g.  $X_i \stackrel{iid}{\sim}_{i=1,2,\dots,n} N(\mu, \sigma^2)$  where  $\mu$  and  $\sigma^2$  are unknown.

**Note:** "iid" means Independent identically distributed

" $\sim$ " means Distributed according to

## 1.3 Nonparametric Models

$X_i \stackrel{iid}{\sim} F_x(x)$  where the cdf  $F$  is completely unknown, but we may assume that  $F$  is smooth, for instance continuous or differentiable.

In the non-parametric setting  $F_x(x)$  should be estimated for every  $x$ . Thus for a random variable  $X$  that can assume infinitely many values, we need to estimate  $F(x)$  at infinitely many values of  $x$ . This is, particularly , the case when  $X$  is a continuous random variable. Recall that  $F_x(x) = P(X \leq x)$ . Then the sample counterpart of  $F_x(x)$  is  $\frac{\#X_i \leq x}{n}$  for a sample  $X_1, \dots, X_n$  . Define:

$$\mathcal{E}(t) = \begin{cases} 1 & \text{if } t \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Then  $\frac{\#X_i \leq x}{n} = \frac{1}{n} \sum_i 1^n \mathcal{E}(x - X_i)$  and  $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n (x - X_i)$  is the Empirical Cumulative Distribution Function (ECDF) .

## 1.4 Point Estimator

Suppose  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, 1)$  where:

$$N(\mu, 1) : f_x(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}}, \quad x \in \mathbb{R}, \mu \in \mathbb{R}$$

We want to have an estimate of  $\mu$ ; i.e. a scientific guess, based on the observations,  $X_1, \dots, X_n$ . Recall that  $\mathbb{E}(X_i) = \mu$ ,  $\mu = 1, 2, \dots, n$ ,  $\mu$  is the population mean.

**What is an "estimate"?**

**Statistic:** A function of observations that does not depend on any unknown parameter.

**Estimator:** An estimator is a statistic that aims at estimating a function of the population unknown parameters.

**Example.**  $X_i \stackrel{iid}{\sim} N(\mu, 1)$

$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is a statistic and as an estimator of  $\mu$

$(\bar{X}_n - \mu)$  is NOT a statistic since it depends on  $\mu$ , which is an unknown parameter.

$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  is a statistic, but not an estimator of  $\mu$ . Note that  $\dim(S^2) = (\dim \mu)^2$ . For instance, if  $X_i$ s are #MISSING-lec1-p3 of a fund and measured in dollars (\$), then  $\dim$  of  $\mu$  is \$ while the  $\dim$  of  $S^2$  is  $\$^2$ . Besides,  $\mu$  can be negative while  $S^2$  is always positive.

## 1.5 Estimation Error

Going back to the example above ( $X_i \stackrel{iid}{\sim} N(\mu, 1)$ ,  $i = 1, 2, \dots, n$ ) and choosing  $\bar{X}_n$  as the estimator of  $\mu$ . We often want to study  $\mathcal{E} = |\bar{X}_n - \mu|$  or a function of  $\mathcal{E}$ . Starting with  $\mathcal{E}$  itself, the first thing that comes to mind is  $P(\mathcal{E} \geq \delta)$  for a prespecified  $\delta$  or perhaps  $\mathbb{E}(\mathcal{E})$ . A well known tool for studying the latter is *Tchbyshev's Inequality*.

## 2 Lecture 2

### 2.1 Markov's Inequality

Let  $X$  be a random variable and  $h$  be a **non-negative** function; ie:

$$h : \mathbb{R} \rightarrow \mathbb{R}^+ \cup \{0\} = [0, \infty)$$

Suppose  $E(h(X)) < \infty$ , then for some  $\lambda > 0$ , we have:

$$P(h(X) \geq \lambda) \leq \frac{E[h(X)]}{\lambda} \quad (1)$$

*Proof.* Suppose  $X$  is a continuous random variable:

$$\begin{aligned} E[h(x)] &= \int_{\mathbb{R}} h(x) f_X(x) dx \\ &= \left( \int_{x:h(x) \geq \lambda} h(x) f_X(x) dx + \int_{x:h(x) < \lambda} h(x) f_X(x) dx \right) \\ &\geq \int_{x:h(x) \geq \lambda} h(x) f_X(x) dx && \text{since } h \geq 0 \\ &\geq \lambda \int_{x:h(x) \geq \lambda} f_X(x) dx = \lambda P(h(X) \geq \lambda) \\ \implies P(h(X) \geq \lambda) &\leq \frac{E(h(X))}{\lambda} \end{aligned}$$

The proof for the discrete case is similar. □

### 2.2 Tchebyshev's Inequality

Tchebyshev's Inequality is a special case of Markov's Inequality. Consider  $h(x) = (x - \mu)^2$ , then:

$$\begin{aligned} P(|X - \mu| \geq \lambda) &= P((X - \mu)^2 \geq \lambda^2) \\ &\leq \frac{E[(X - \mu)^2]}{\lambda^2} && \text{if } E[(X - \mu)^2] < \infty \end{aligned}$$

Let  $\mu = E(X)$ , then  $E[(X - \mu)^2] = \text{Var}(X)$  denoted by  $\sigma_x^2$ . We therefore have:

$$P(|X - \mu_x| \geq \lambda) \leq \frac{\sigma_x^2}{\lambda^2} \quad \text{where } \mu_x = E(X) \quad (2)$$

Now consider  $\lambda = K\sigma_x$  where  $K$  is a known number. Then:

$$P(|X - \mu_x| \geq K\sigma_x) \geq \frac{\sigma_x^2}{K^2\sigma_x^2} = \frac{1}{K^2} \quad (3)$$

This is called **Tchbyshev's Inequality**.

**Example 2.1.** Suppose  $K = 3$ .

$$P(|X - \mu_x| \geq 3\sigma_x) \leq \frac{1}{9}$$

*In other words, at least 88% of the observations are within 3 standard deviation from the population mean.*

Going back to the our example:

$$X_i \sim (\mu, 1) \quad , \quad \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

We want to study  $P(\epsilon \geq \delta) = P(|\bar{X}_n - \mu| \geq \delta)$ , first we note that:

$$E(X_i) = \mu \quad , \quad i = 1, 2, \dots, n$$

Then:

$$\begin{aligned} E(\bar{X}_n) &= E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E(X_i) \\ &= \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n} (n\mu) = \frac{1}{n} \cdot (n\mu) \\ &= \mu \end{aligned}$$

(\*)



Thus, using (2) we have:

$$P(|\bar{X}_n - \mu| \geq \delta) \leq \frac{Var(\bar{X}_n)}{\delta^2}$$

Now:

$$\begin{aligned} Var(\bar{X}_n) &= Var\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} Var\left(\sum_{i=1}^n X_i\right) \\ &= \frac{1}{n^2} \left[ \sum_{i=1}^n Var(X_i) + \sum_{1 \leq i < j \leq n} \sum_{1 \leq i < j \leq n} Cov(X_i, X_j) \right] \quad \text{using Thm 5.12(b) - page 271} \\ &= \frac{1}{n^2} \sum_{i=1}^n Var(X_i) \quad \text{since } \prod_{i=1}^n X_i \\ &= \frac{1}{n^2} n Var(X) = \frac{Var(X)}{n} \quad \text{since } x_i \text{'s are identically distributed} \\ &= \frac{\delta_X^2}{n} \end{aligned} \quad (**)$$

In our case  $X \sim N(\mu, 1)$  so  $Var(X) = \delta_X^2 = 1$ . Thus  $Var(\bar{X}_n) = \frac{1}{n}$

**Remark.**  $X \perp\!\!\!\perp Y \implies Cov(X, Y) = 0$ . Note that:

$$X \perp\!\!\!\perp Y \implies E[g_1(X)g_2(Y)] = E[g_1(X)].E[g_2(Y)]$$

in particular:

$$X \perp\!\!\!\perp Y \implies E[XY] = E[X].E[Y]$$

on the other hand:

$$Cov(X, Y) = E[XY] - E(X)E(Y)$$

thus:

$$X \perp\!\!\!\perp Y \implies Cov(X, Y) = 0.$$

recall that  $X \perp\!\!\!\perp Y$  means  $X$  and  $Y$  are independent, i.e.  $f_{X,Y}(x, y) = f_X(x)f_Y(y)$

where  $f_{X,Y}, f_X$  and  $f_Y$  represent respectively the

We therefore have:

$$P(|\bar{X}_n - \mu| \geq \delta) \leq \frac{1}{n\delta^2} \quad (4)$$

Using (4) and the sample size,  $n$ , we can find an upper bound for the proportion of deviations which are greater than a given threshold  $\delta$ .

We can also use (4) for Sample Size Deterministic:

Suppose  $\delta$  is given and we want  $P(|\bar{X}_n - \mu| \geq \delta) \leq \beta$  where  $\beta$  is also given. Then setting  $\frac{1}{n\delta^2} = \beta$ , we can estimate  $n \approx \frac{1}{\beta\delta^2}$ .

## 2.3 Application to Voting

Define  $X_i = \begin{cases} 1 & \text{NDP} \\ 0 & \text{otherwise} \end{cases}$ . Associated to each eligible voter in Canada we

have a binary variable  $X$ . Let  $p = P(X = 1)$ . So  $p$  represents the proportion of eligible voters who favor *NDP*. Of interest is often estimation of  $p$ . Suppose we have a sample of size  $n$ ,  $X_1, X_2, \dots, X_n$ .

$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is the sample proportion; The counterpart of  $p$  which can be denoted by  $\hat{p}$ . Note that:

$$\mu_x = E(X) = 1 \times P(X = 1) + 0 \times P(X = 0) = 1 - p + 0 \times (1 - p) = p$$

and:

$$E(X^2) = 1^2 \times P(X = 1) + 0^2 \times P(X = 0) = 1 - p + 0 \times (1 - p) = p$$

From (\*) and (\*\*) we find that :

$$E(\hat{p}_n) = E(\bar{X}_n)\mu_x = p$$

and:

$$\text{Var}(\hat{p}_n) = E(\bar{X}_n) = \frac{\text{Var}(X)}{n} = \frac{\sigma_X^2}{n} = \frac{p(1-p)}{n}$$

Thus using (2), we have:

$$P(|\hat{p}_n - p| \geq \delta) \leq \frac{\text{Var}(\hat{p}_n)}{\delta^2} = \frac{p(1-p)}{n\delta^2}$$

Note that the above bound on the probability of derivation depends on  $p$  which is *unknown*. We however notice that  $p(1-p) \leq \frac{1}{4}$ .

Define  $\mathcal{C}(x) = x(1-x)$  for  $0 < x < 1$ . Then:

$$\begin{aligned} \mathcal{C}'(x) = 1 - 2x &\implies \mathcal{C}'(x) = 0 \implies x = \frac{1}{2} \\ \mathcal{C}''\left(\frac{1}{2}\right) = -2 &\implies x = \frac{1}{2} \quad \text{which is a **maximizer**} \\ \mathcal{C}\left(\frac{1}{2}\right) = \frac{1}{2}\left(1 - \frac{1}{2}\right) &= \frac{1}{4} \end{aligned}$$

(Note that  $\mathcal{C}''(x) = -2$  for all  $0 < x < 1$ )

We therefore find:

$$P(|\hat{p}_n - p| \geq \delta) \leq \frac{1}{4n\delta^2} \tag{5}$$

Using (5) and a given sample size  $n$  we can find an upper bound for the probability of derivation by  $\delta$  and the amount for any given  $\delta$ .

We can also use (5) for sample size deterministic for a size bound  $\beta$  and derivative  $\delta$  as follows:

$$\frac{1}{4n\delta^2} = \beta \implies n \geq \frac{1}{4\beta\delta^2}$$

This is of course conservative since  $p(1-p) \leq \frac{1}{4}$ .

### 3 Lecture 3

#### 3.1 MSE

**MSE:** To study estimation error we started by studying  $P(|\hat{\Theta}_n - \Theta| > \delta)$ , deviation above a given threshold  $\delta$ , by bounding this probability. One may take a different approach by studying average Euclidean distance, i.e.  $E[|\hat{\Theta}_n - \Theta|^2]$ , which denoted by  $\text{MSE}(\hat{\Theta}_n)$ .

We note that if  $\Theta = E(\hat{\Theta}_n)$ , i.e.  $\hat{\Theta}_n$  is an unbiased estimation of  $\Theta$ , then:

$$\text{MSE}(\hat{\Theta}_n) = E[|\hat{\Theta}_n - \Theta|^2] = E[(\hat{\Theta}_n - \mu_{n_{\Theta_n}})^2] = \text{Var}(\hat{\Theta}_n)$$

Now recall that  $\text{Var}(X) = 0 \implies P(X = \text{constant}) = 1$  which essentially means random variable  $X$  is a constant.

The same comment applies to  $\text{MSE}(\hat{\Theta}_n)$ . We want to find the closest estimator  $\hat{\Theta}_n$  to  $\Theta$  which means that we want to minimize  $E[(\hat{\Theta}_n - \Theta)^2]$  over all possible estimators, ideally at least the above comment tells us that in real applications we cannot expect to find an estimator whose MSE is equal to zero. Let's try to understand the MSE a bit more:

$$\begin{aligned} \text{MSE}(\hat{\Theta}_n) &= E[(\hat{\Theta}_n - \Theta)^2] \\ &= E\left[\left((\hat{\Theta}_n - E(\hat{\Theta}_n)) + (E(\hat{\Theta}_n) - \Theta)\right)^2\right] \\ &= E\left[(\hat{\Theta}_n - E(\hat{\Theta}_n))^2\right] + \underbrace{E\left[(E(\hat{\Theta}_n) - \Theta)^2\right]}_{\text{not a r.v.}} + 2 \cdot \underbrace{E\left[(\hat{\Theta}_n - E(\hat{\Theta}_n)) \cdot (E(\hat{\Theta}_n) - \Theta)\right]}_{\text{not a r.v.}} \\ &= E\left[(\hat{\Theta}_n - E(\hat{\Theta}_n))^2\right] + \underbrace{E\left[(E(\hat{\Theta}_n) - \Theta)^2\right]}_{\text{Bias}(\hat{\Theta}_n)} + 2 \cdot \underbrace{E\left[(\hat{\Theta}_n - E(\hat{\Theta}_n)) \cdot (E(\hat{\Theta}_n) - \Theta)\right]}_{E(\hat{\Theta}_n) - E(\hat{\Theta}_n) = 0} \\ &= \text{Var}(\hat{\Theta}_n) + \underbrace{\left[E(\hat{\Theta}_n) - \Theta\right]^2}_{\text{Bias}(\hat{\Theta}_n)} + 2 \cdot \underbrace{\text{Bias}(\hat{\Theta}_n) \cdot E[(\hat{\Theta}_n - E(\hat{\Theta}_n))]}_{E(\hat{\Theta}_n) - E(\hat{\Theta}_n) = 0} \\ &= \text{Var}(\hat{\Theta}_n) + \text{Bias}^2(\hat{\Theta}_n) \end{aligned}$$

Roughly speaking, **bias** measures how far off the target we hit on the average while **variance** measures how much fluctuation our estimator may show from one sample to another.

### 3.2 Unbiased Estimators

In almost all real applications, the class of possible estimators for an **ESTIMANAL** is huge and the best estimator, i.e. the one that minimizes MSE no matter what the value of the **ESTIMANAL** is, almost never exists. Thus we try to reduce the class of potential estimators by improving a plausible restriction, for example  $\text{Bias}(\hat{\Theta}_n) = 0$ .

**Definition.** An estimator  $\hat{\Theta}_n$  of an **ESTIMANAL**  $\Theta$  is said to be **unbiased** if  $E(\hat{\Theta}_n) = \Theta$ , for all possible values of  $\Theta$ .

**Example 3.1.**  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2) \quad i = 1, 2, \dots, n$

Suppose both  $\mu$  and  $\sigma^2$  are unknown. Consider  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ .

$$E(\bar{X}_n) = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \overbrace{E(X_i)}^{\mu} = \frac{1}{n} \cdot n\mu = \mu$$

Thus  $\bar{X}_n$  is an unbiased estimator of  $\mu$ . As for the  $\text{MSE}(\bar{X}_n)$ , we need to find  $\text{Var}(\bar{X}_n)$ .

$$\begin{aligned} \text{Var}(\bar{X}_n) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) \\ &= \frac{1}{n^2} \left[ \sum_{i=1}^n \text{Var}(X_i) + 2 \cdot \sum_{1 \leq i < j \leq n} \overbrace{\text{Cov}(X_i, X_j)}^0 \right] && \text{Theorem 5.12(b) - page 271} \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) && \prod_{i=1}^n X_i \\ &= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{1}{n^2} \cdot n\sigma^2 = \frac{\sigma^2}{n} && \text{identically distributed} \\ \implies \text{MSE}(\bar{X}_n) &= \text{Var}(\bar{X}_n) + \overbrace{\text{Bias}^2(\bar{X}_n)}^0 = \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n} \end{aligned}$$

An inspection of the above calculation shows that for unbiased  $\mu$  we only require a common mean  $\mu$  while for calculating the variance we would only

require a common variance  $\sigma^2$  and orthogonality, i.e:

$$\text{Cov}(X_i, X_j) = 0 \quad \text{where } i \neq j$$

Suppose  $X_1, \dots, X_n$  have the same mean value  $\mu$ . Then:

$$E(\bar{X}_n) = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} n \mu = \mu$$

Suppose further that  $X_1, \dots, X_n$  have the same variance  $\sigma^2$  and  $\text{Cov}(X_i, X_j) = 0, i \neq j$ .

Then:

$$\begin{aligned} \text{Var}(\bar{X}_n) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) \\ &= \frac{1}{n^2} \left[ \sum_{i=1}^n \text{Var}(X_i) + 2 \sum_{1 \leq i < j \leq n} \text{Cov}(X_i, X_j) \right] && \text{Theorem 5.12(b) - Page 271} \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) && \text{Orthogonality: i.e. } \text{Cov}(X_i, X_j) = 0 \text{ if } i \neq j \\ &= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{1}{n^2} n \sigma^2 = \frac{\sigma^2}{n} && \text{having the same variance} \\ &\implies \text{MSE}(\bar{X}_n) = \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n} \end{aligned}$$

If  $X_1, \dots, X_n$  have the same mean value and variance and they are orthogonal.

### 3.3 Stein's Paradox

We will learn later that if  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$  then  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  has many optimal properties. A paradox due to Charles Stein, however, shows that such a nice optimal properties are not preserved in higher dimensions. In fact if:

$$X_i \stackrel{iid}{\sim} N(\mu_x, 1), \quad Y_i \stackrel{iid}{\sim} N(\mu_y, 1) \text{ and } Z_i \stackrel{iid}{\sim} (\mu_z, 1)$$

then, we can find the biased estimators of  $\begin{pmatrix} \mu_x \\ \mu_y \\ \mu_z \end{pmatrix}$  which are closer to  $\begin{pmatrix} \mu_x \\ \mu_y \\ \mu_z \end{pmatrix}$  than  $\begin{pmatrix} \bar{X}_n \\ \bar{Y}_n \\ \bar{Z}_n \end{pmatrix}$  for any  $\begin{pmatrix} \mu_x \\ \mu_y \\ \mu_z \end{pmatrix}$ . We may then say that  $\begin{pmatrix} \bar{X}_n \\ \bar{Y}_n \\ \bar{Z}_n \end{pmatrix}$  is an **inadmissible estimator** of  $\begin{pmatrix} \mu_x \\ \mu_y \\ \mu_z \end{pmatrix}$ .

### 3.4 Admissibility

An estimator  $\hat{\Theta}$  is called admissible if there is no estimator  $\tilde{\Theta}$  such that:

$$MSE(\tilde{\Theta}) \leq MSE(\hat{\Theta}) \quad \text{for all possible values of } \Theta$$

and this inequality is strict for some values of  $\Theta$ .

What this example tells us is that by allowing a bit of bias we may be able to reduce variance considerably and hence find an estimator which is closer to the target than the most natural unbiased estimator. Note that this phenomena happens only when the dimension is at least 3.

## 4 Lecture 4

We now want to restrict the class of estimators even further. Suppose  $X_1, \dots, X_n$  have the same mean  $\mu$  and variance  $\sigma^2$  and they are orthogonal; i.e.  $\text{Cov}(X_i, X_j) = 0$ ,  $i \neq j$ . Consider  $\tilde{X}_{n,\tilde{C}} = \sum_{i=1}^n C_i X_i$  and

$$\mathcal{C} = \left\{ \tilde{X}_{n,\tilde{C}} : \tilde{C} = (C_1, \dots, C_n) \in \mathbf{R}^n, \sum_{i=1}^n C_i = 1 \right\}$$

Note that

$$\begin{aligned} E(\tilde{X}_{n,\tilde{C}}) &= E\left(\sum_{i=1}^n C_i X_i\right) = \sum_{i=1}^n C_i E(X_i) \\ &= \sum_{i=1}^n C_i \mu = \mu \underbrace{\sum_{i=1}^n C_i}_1 = 1 \cdot \mu \\ &= \mu \end{aligned}$$

Thus  $\tilde{X}_{n,\tilde{C}}$  is an unbiased estimator of  $\mu$  for any  $\tilde{C} \in \mathbf{R}^n$  as long as  $\sum_{i=1}^n C_i = 1$ . Then  $\mathcal{C}$  is the class of all unbiased linear estimators of  $\mu$ . We want to find the best estimator with  $\mathcal{C}$ ; i.e.:

$$\underset{\tilde{C} \in \mathbf{R}^n}{\text{Min}} \text{MSE}(\tilde{X}_{n,\tilde{C}}) \quad \text{s.t.} \quad \sum_{i=1}^n C_i = 1 \quad (*)$$



First we note that  $MSE(\tilde{X}_{n,\tilde{C}}) = Var(\tilde{X}_{n,\tilde{C}})$  since  $\tilde{X}_{n,\tilde{C}}$  is an unbiased estimator of  $\mu$  when  $\sum_{i=1}^n C_i = 1$ . On the other hand:

$$\begin{aligned}
 Var(\tilde{X}_{n,\tilde{C}}) &= Var\left(\sum_{i=1}^n C_i X_i\right) \\
 &= \sum_{i=1}^n C_i^2 Var(X_i) + 2 \sum_{1 \leq i < j \leq n} Cov(C_i X_i, C_j X_j) \quad \text{Theorem 5.12 page 271} \\
 &= \sum_{i=1}^n C_i^2 \sigma^2 + 2 \sum_{1 \leq i < j \leq n} \overbrace{C_i C_j Cov(X_i, X_j)}^0 \\
 &= \sigma^2 \sum_{i=1}^n C_i^2
 \end{aligned}$$

Thus (\*) is equivalent to :

$$\underset{\tilde{C} \in \mathbb{R}^n}{\text{Min}} \sigma^2 \sum_{i=1}^n C_i^2 \quad (**)$$

Using the *Lagrange Theorem*, (\*\*) is equivalent to:

$$\tilde{C} = (C_1, \dots, C_n) \in \mathbb{R}^n \quad \overbrace{\left\{ \sigma^2 \sum_{i=1}^n C_i + \lambda \left( \sum_{i=1}^n C_i - 1 \right) \right\}}^{\mathcal{C}_\lambda(\tilde{C})} .$$

Note that:  $\frac{\partial \mathcal{C}_\lambda(\tilde{C})}{\partial C_i} = 2 \sigma^2 C_i + \lambda \quad , \quad i = 1, 2, 3, \dots$

$$\frac{\partial}{\partial \lambda} \mathcal{C}_\lambda(\tilde{C}) = \sum_{i=1}^n C_i - 1$$

$$\begin{cases} \frac{\partial}{\partial C_i} \mathcal{C}_\lambda(\tilde{C}) = 2 \sigma^2 C_i + \lambda = 0 \quad , \quad i = 1, 2, 3, \dots \\ \frac{\partial}{\partial \lambda} \mathcal{C}_\lambda = 0 \implies \sum_{i=1}^n C_i = 1 \end{cases}$$

Thus  $C_i = -\frac{\lambda}{2 \sigma^2} \quad , \quad i = 1, 2, 3, \dots, n$  and using the last equation:

$$\sum_{i=1}^n -\frac{\lambda}{2 \sigma^2} = 1 \implies \lambda = -\frac{2 \sigma^2}{n}$$

and therefore:

$$C_i = -\frac{\lambda}{2\sigma^2} = -\frac{-\frac{2\sigma^2}{n}}{2\sigma^2} = \frac{1}{n}, \quad i = 1, 2, 3, \dots, n$$

We can further find:

$$\mathcal{H} = [\frac{\partial^2}{\partial C_i \partial C_j} \mathcal{C}_\lambda(\tilde{C})] \quad , \quad i, j = 1, 2, \dots, n$$

and show that:

$$\begin{aligned} \tilde{x}^T \mathcal{H} \tilde{x} &\geq 0 \quad \forall \tilde{x} \in \mathbb{R}^n \\ &= 0 \quad \text{if and only if } \tilde{x} = 0 \end{aligned}$$

This then guarantees that  $\tilde{C}^* = (\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$  is indeed a minimizer; in fact, the *unique minimizer*. To summarize:

$$\tilde{X}_{n, \tilde{C}^*} = \sum_i i = 1^n \frac{1}{n} X_i = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n$$

Thus  $\bar{X}_n$  is the best unbiased linear estimator.

## 4.1 Estimating Variance

So far we confirmed ourselves to estimation of the population mean.

Now suppose we are interested in estimating variance from  $X_1, \dots, X_n$  where  $X_i$ s have the same mean value  $\mu$ , the same variance  $\sigma^2$  and they are orthogonal, i.e.  $Cov(X_i, X_j) = 0$ ,  $i \neq j$ , then a *natural estimator* of:

$$\sigma^2 = Var(X) = \mathbb{E}[(x - \mu)^2]$$

is its sample counterpart, i.e.

$$S_{n,*}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

Now the first question is if  $S_{n,*}^2$  is an unbiased estimator of  $\sigma^2$ , i.e.  $\mathbb{E}(S_{n,*}^2) = \sigma^2$

$$\begin{aligned}
(X_i - \mu)^2 &= [(X_i - \bar{X}_n) + (\bar{X}_n - \mu)]^2 \\
&= (X_i - \bar{X}_n)^2 + (\bar{X}_n - \mu)^2 + 2 \cdot (X_i - \bar{X}_n)(\bar{X}_n - \mu) \\
\sum_{i=1}^n (X_i - \mu)^2 &= \sum_{i=1}^n (X_i - \bar{X}_n)^2 + n(\bar{X}_n - \mu)^2 + 2 \cdot (\bar{X}_n - \mu) \overbrace{\sum_{i=1}^n (X_i - \bar{X}_n)}^0 \\
&= \sum_{i=1}^n (X_i - \bar{X}_n)^2 + n(\bar{X}_n - \mu)^2
\end{aligned} \tag{I}$$

Taking estimation we find:

$$\begin{aligned}
\mathbb{E}\left[\sum_{i=1}^n (X_i - \mu)^2\right] &= \mathbb{E}[n \cdot S_{n,*}^2] + \mathbb{E}[n(\bar{X}_n - \mu)^2] \\
RHS &= \sum_{i=1}^n \overbrace{\mathbb{E}(X_i - \mu)^2}^{\sigma^2} = n \cdot \sigma^2
\end{aligned} \tag{II}$$

Note that  $\mathbb{E}(\bar{X}_n - \mu) = 0$ , i.e.  $\mathbb{E}(\bar{X}_n) = \mu$ . Thus:

$$\mathbb{E}[n(\bar{X}_n - \mu)^2] = n \mathbb{E}[(\bar{X}_n - \mu)^2] = n \text{Var}(\bar{X}_n).$$

On the other hand  $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$ . We therefore have:

$$\mathbb{E}[n(\bar{X}_n - \mu)^2] = n \cdot \text{Var}(\bar{X}_n) = n \cdot \frac{\sigma^2}{n} = \sigma^2$$

and hence from (II):

$$n\sigma^2 = \mathbb{E}(n S_{n,*}^2) + \sigma^2$$

which implies:

$$\implies \mathbb{E}(S_{n,*}^2) = \left(\frac{n-1}{n}\right)\sigma^2 = \left(1 - \frac{1}{n}\right)\sigma^2$$

meaning that  $S_{n,*}^2$  is **NOT** an unbiased estimator of  $\sigma^2$ .

Multiplying both sides of the last equation by the reciprocal of  $(1 - \frac{1}{n})$  we find

$\mathbb{E}(\frac{n}{n-1} S_{n,*}^2) = \sigma^2$ . Note however that:

$$\frac{n}{n-1} S_{n,*}^2 = \frac{n}{n-1} \cdot \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

Thus  $\boxed{S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2}$  is an unbiased estimator.

**Question: Why  $(n-1)$ ?**

" $n-1$ " is the dimension of  $\overbrace{\text{span}\{X_i - \bar{X}_n : i = 1, 2, \dots, n\}}^V$ .

$n-1 = \dim(\text{span } V)$ . Note however  $\dim(\text{span } W) = n$  where  $W = X_i - \mu$ ,  $i = 1, 2, \dots, n$ .

We discuss these issues further in Chapter 11 where we learn about the regression.

So far we only considered sampling from one population. We may have samples from two or more populations and may want to make inference about differences between the populations.

#### Example 4.1.

Suppose we want to study the differences between the average salaries of men and women:

Men	Women
$X_1$	$Y_1$
$\vdots$	$\vdots$
$X_m$	$Y_n$

where  $X_i$ s have the common mean  $\mu_x$  and  $Y_j$ s have the common mean  $\mu_y$ .

We want to estimate  $\mu_x - \mu_y$ . The natural estimator is  $\bar{X}_m - \bar{Y}_n$ . Show that:

$$\mathbb{E}[\bar{X}_m - \bar{Y}_n] = \mu_x - \mu_y$$

Hence  $\bar{X}_m - \bar{Y}_n$  is an unbiased estimator of  $\mu_x - \mu_y$ .

Assume further that  $X$ s and  $Y$ s are independent and  $X$ s have common vari-

ance  $\sigma_x^2$  and  $Y$ s have common variance  $\sigma_y^2$  and  $Cov(X_i, X_j) = 0$ ,  $i \neq j$  and  $Cov(Y_i, Y_j) = 0$ ,  $i \neq j$ .

Find  $Var(\bar{X}_m - \bar{Y}_n)$ . Hint: use *Thm* 5.12.

The difference between two proportions can be treated similarly. Note that proportions are essentially means of binary variables.

## 5 Lecture 5 : Confidence Intervals

**Definition.** *Random Interval* An interval whose endpoint(s) are random variables is called a **Random Interval**.

### 5.1 Confidence Intervals

A  $100(1 - \alpha)\%$  confidence interval for a parameter  $\theta$  is a *random interval*  $(\hat{\Theta}_L(X), \hat{\Theta}_V(X))$  such that:

$$P(\hat{\Theta}_L(X) < \Theta < \hat{\Theta}_V(X))$$

**Pivotal Quantity** : A function of the observation  $X_1, \dots, X_n$  and some unknown parameters, ideally just the parameter(s) of interest, whose distribution DOES NOT depend on any unknown parameter is called **Pivotal Quantity**.

Pivotal Quantities play a central role in theory confidence intervals.

**Example.** Let  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $i = 1, 2, \dots, n$  where  $\sigma^2$  is known, but  $\mu$  is unknown. We show that

$$\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right) \quad \text{where} \quad \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Recall that there are three methods for finding the distribution of a function of random variables:

1. **Method of Transformation:** This is essentially theorem of change of variables in calculus.
2. **Method of Distribution:** On this method we connect the *cdf* of the new variable to the *cdf* of the original variables.

**Example.** Suppose  $X_i \stackrel{iid}{\sim} f$ ,  $i = 1, 2, \dots, n$  are continuous random variables with pdf  $f$  and cdf  $F$ . Define  $X_{(n)} = \max_{1 \leq i \leq n}$ .

$$\begin{aligned}
F_{X(n)}(t) &= P(X_{(n)} \leq t) = P(X_1 \leq t, X_2 \leq t, \dots, X_n \leq t) \\
&= \prod_{i=1}^n P(X_i \leq t) && \text{(by } \prod_{i=1}^n X_i \text{)} \\
&= \prod_{i=1}^n F_{X_i}(t) \\
&= \prod_{i=1}^n F(t) = F^n(t) && \text{identically distributed}
\end{aligned}$$

Thus

$$\begin{aligned}
f_{X(n)}(t) &= \frac{d}{dt} F_{X(n)}(t) = \frac{d}{dt} F^n(t) \\
&= n f(t) F^{n-1}(t)
\end{aligned}$$

**3. Method of Moment Generating Function(mgf):** This method is essentially based on the *mgf* of the new variable of the *mgf* if the original variables.

**Example.** Suppose  $X_i \sim N(\mu_i, \sigma_i^2)$ ,  $i = 1, 2, \dots, n$  and  $X_i$ s are independent. Define  $S = \sum_{i=1}^n X_i$ .

$$\begin{aligned}
m_s(t) &= \mathbb{E}[e^{tS}] = \mathbb{E}[e^{t \sum_{i=1}^n X_i}] \\
&= \mathbb{E}\left[\prod_{i=1}^n e^{tX_i}\right] && \text{using independence: } (\prod) \\
&= \prod_{i=1}^n m_{X_i}(t) \\
&= \prod_{i=1}^n e^{\mu_i t + \frac{\sigma_i^2 t^2}{2}} \\
&= \exp\left\{t \sum_{i=1}^n \mu_i + \frac{t^2}{2} \sum_{i=1}^n \sigma_i^2\right\} \\
\Rightarrow S &\sim N\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)
\end{aligned}$$

If we further assume that  $X_i$ 's are identically distributed, then:

$$\mu_i = \mu \quad \text{and} \quad \sigma_i^2 = \sigma^2 \quad \forall i = 1, 2, \dots, n$$

Therefore we have:

$$m_S(t) = \exp\left\{n\mu t + \frac{n\sigma^2 t^2}{2}\right\}$$

and hence:  $\boxed{S \sim N(n\mu, n\sigma^2)}$  Then:

$$\begin{aligned} m_{\bar{X}_n}(t) &= \mathbb{E}[e^{t\bar{X}_n}] = \mathbb{E}[e^{t\frac{1}{n}\sum_{i=1}^n X_i}] \\ \text{by } t^* = \frac{t}{n} &\implies = E[e^{t^* S}] \\ &= m_S(t^*) = e^{n\mu t^* + \frac{n\sigma^2 t^{*2}}{2}} \\ &= \exp\left\{n\mu t^* + \frac{n\sigma^2 t^{*2}}{2}\right\} \\ \implies \boxed{\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right)} &\quad (1) \end{aligned}$$

Note further that if  $X \sim N(\mu, \sigma^2)$ , then  $Z = \frac{X-\mu}{\sigma} \sim N(0, 1)$ . We prove a general form of this. Let  $X \sim N(\mu, \sigma^2)$ ; then:

$$aX + b \sim N(a\mu + b, a^2\sigma^2) \quad \text{for any constant } a, b$$



Let  $V = ax + b$ , then:

$$\begin{aligned}
 m_v(t) &= \mathbb{E}[e^{tV}] = \mathbb{E}[e^{t(ax+b)}] \\
 &= \mathbb{E}[e^{taX+tb}] = \mathbb{E}\left[\underbrace{e^{tb}}_{\text{constant}} \cdot \overbrace{e^{taX}}^{t^*}\right] \\
 &= e^{tb} \cdot \mathbb{E}[e^{t^*X}] \\
 &= e^{tb} \cdot e^{t\mu + \frac{\sigma^2 t^2}{2}} \\
 &= \exp\left\{t(a\mu + b) + \frac{t^2(a^2\sigma^2)}{2}\right\}
 \end{aligned}$$

Thus :

$$(ax + b) \sim N(a\mu + b, a^2\sigma^2)$$

Now :

$$\begin{aligned}
 Z &= \frac{X - \mu}{\sigma} = \frac{X}{\sigma} - \frac{\mu}{\sigma} \\
 &= \frac{1}{\sigma}X - \frac{\mu}{\sigma} \\
 &= aX + b \quad \text{where } a = \frac{1}{\sigma} \text{ and } b = -\frac{\mu}{\sigma}
 \end{aligned}$$

Hence :

$$Z \sim N\left(\overbrace{\frac{1}{\sigma}\mu + (-\frac{\mu}{\sigma})}^0, \overbrace{(\frac{1}{\sigma})^2\sigma^2}^1\right)$$

Thus :

$$\boxed{Z \sim N(0, 1)} \quad (2)$$

Using (1) and (2) :

$$\frac{\bar{X}_n - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1).$$

This means that :

$$\boxed{\frac{\bar{X}_n - \mu}{\sqrt{h}} \text{ is a Pivotal Quantity}}$$

To summarize:

$$X_i \stackrel{iid}{\sim} N(\mu, \sigma^2) \quad \Rightarrow \quad \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \quad \text{is a \textbf{pivotal quantity}}$$

Notice that using the table for the normal distribution:

$$P\left(\left|\frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}}\right| \leq 1.96\right) = 0.95$$

Equivalently:

$$P\left(\bar{X}_n - 1.96 \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + 1.96 \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

This means that:

$$\left(\bar{X}_n - 1.96 \frac{\sigma}{\sqrt{n}}, \bar{X}_n + 1.96 \frac{\sigma}{\sqrt{n}}\right)$$

covers the true  $\mu$  with 95% probability.

Thus a  $100(1-\alpha)\%$  confidence interval for  $\mu$  where  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$  and  $\sigma^2$  is known:

$$\bar{X}_n \pm \zeta_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

where:

$$P(Z > \zeta_{\frac{\alpha}{2}}) = \frac{\alpha}{2}, \quad Z \sim N(0, 1)$$

#MISSING GRAPH - (LECTURE 5 - P5)

**Remark.** In real applications we compute  $\bar{X}_n$  and obtain an interval, say (125, 135). Now either this interval covers the true  $\mu$  or it does not. Then the question is what do we mean by a 95% #MISSING ?

Note that the  $100(1-\alpha)\%$  confidence is the property of the procedure. It means that out of the all possible intervals of the form  $(\bar{X}_n \pm 1.96 \frac{\sigma}{\sqrt{n}})$  that we can make by taking samples of size  $n$  from  $N(\mu, \sigma^2)$ , 95% of them cover the true  $\mu$ . Now this is a real application when we make one of the such intervals by taking a random sample of size  $n$  from  $N(\mu, \sigma^2)$ , it is like taking one of those intervals

randomly. Since that 95% of them cover  $\mu$ , my chance of selecting an interval that covers  $\mu$  is 95%. Thus I can take a bet 19 to 1 that the interval I select covers  $\mu$ .

## 5.2 Large Sample Confidence Interval

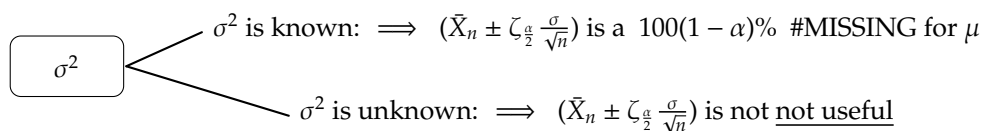
The derivation of the pivotal quantity in the above example totally hinges over the normality assumption, i.e.  $X_i \sim N(\mu, \sigma^2)$ .

What happens if we do not know the parametric for the population distribution?

**Theorem (General Limit Theorem - GLT (baby version)).** Suppose  $X_1, \dots, X_n$  are independent random variables with common  $\mu$  and variance  $\sigma^2$ . Then:

$$\frac{\bar{X}_n - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)} \stackrel{app}{\sim} N(0, 1) \quad \text{when } n \text{ is large enough}$$

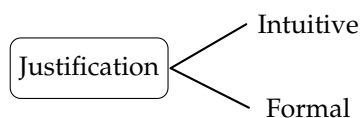
This is a powerful theorem that implies that  $\frac{\bar{X}_n - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)}$  is approximately a pivotal quantity distributed according to  $N(0, 1)$  for large enough  $n$  regardless of population distribution provided that the condition of the GLT are met.



Note: if  $\sigma^2$  is unknown,  $(\bar{X}_n \pm z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}})$  is still a 100(1 -  $\alpha$ )% #MISSING for  $\mu$  but not useful.

We need to somehow get rid of the #MISSING parameter  $\sigma$ . We can replace  $\sigma$  by  $S_n$  where:

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$



**Intuitive** :  $S_n^2$  is the sample counterpart, almost, for  $\sigma^2$  . Thus as  $n$  increases, greater portion of the population and hence our sample sets closer to the population.

**Formal** : The formal proof comprises three steps:

1. **GLT** of

$$\frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}}$$

2. Consistency of  $S_n^2$  for  $\sigma^2$ , i.e.  $S_n^2 \xrightarrow{P} \sigma^2$  which we learn in *Chapter 9*. We then use a theorem called **Continuous Mapping Theorem** which says that if  $S_n^2 \xrightarrow{P} \sigma^2$ , then:

$$g(S_n^2) \xrightarrow{P} g(\sigma^2) \quad \text{for any continuous function}$$

Considering  $g(x) = \sqrt{x}$ , we obtain  $S_n \xrightarrow{P} \sigma$  and hence  $\frac{\sigma}{S_n} \xrightarrow{P} 1$ .

3. **Cramer's Theorem** :

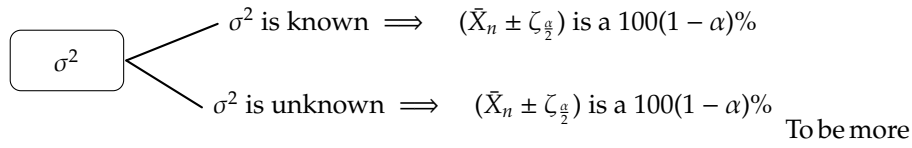
This result says that if  $V_n \xrightarrow{D} X$  and  $Y_n \xrightarrow{P} 1$ , then  $Y_n \cdot V_n \xrightarrow{D} X$  :

$$\frac{\bar{X}_n - \mu}{\frac{S_n}{\sqrt{n}}} = \underbrace{\frac{\sigma}{S_n}}_{Y_n} \cdot \underbrace{\frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}}}_{V_n}$$

Note that GLT implies that  $V_n \xrightarrow{D} 2$ , i.e:

$$\overbrace{F_{V_n}(t)}^{\text{cdf of } V_n} \rightarrow F_z(t) = \overbrace{\int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx}^{\Phi(t) \text{ cdf of } N(0,1)}$$

Using step (2),  $Y_n = \frac{\sigma}{S_n} \xrightarrow{P} 1$  and application of Cramer's Theorem computes the proof. To summarize:



precise, these confidence intervals are approximate  $100(1 - \alpha)\%$  confidence intervals for  $\mu$  when  $n$  is large enough.

So far we focused on C.I for population mean. How can we make C.I for other estimates?

A common, perhaps the most common, method of estimation that we will learn about in Chapter 9 is the method of maximum likelihood. Suppose  $\Theta$  is a parameter of interest. Suppose  $\hat{\Theta}_n = \hat{\Theta}(X_1, \dots, X_n)$  is the maximum likelihood estimate (MLE) of  $\Theta$  based on  $X_1, \dots, X_n$ . Then relatively several condition we have:

$$\frac{\hat{\Theta}_n - \Theta}{\sqrt{\text{Var}(\hat{\Theta}_n)}} \stackrel{app}{\sim} N(0, 1) \text{ when } n \text{ is large enough} \quad (*)$$

We therefore have a several recipe for confidence interval when the sample size  $n$  is large enough, namely:

$$\hat{\Theta}_n \pm \zeta_{\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\Theta}_n)} \quad (\dagger)$$

that is a  $100(1 - \alpha)\%$  C.I for  $Q$ .

**Example 5.1.**  $X_i \stackrel{iid}{\sim} N(\mu, \overbrace{\sigma^2}^{\text{known}})$ ,  $i = 1, 2, \dots, n$ .

We show in chapter 9 that  $\bar{X}_n$  is the MLE of  $\mu$ . Note that  $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$ . Then using  $(\dagger)$ :

$$\bar{X}_n \pm \zeta_{\frac{\alpha}{2}} \sqrt{\frac{\sigma^2}{n}} \text{ is a } 100(1 - \alpha)\% \text{ C.I for } \mu$$

**Example 5.2.**  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$   $\forall i = 1, 2, \dots, n$ , i.e:

$$X_i = \begin{cases} 1 & p \\ 0 & 1 - p \end{cases}$$

Then  $\hat{p}_n = \frac{1}{n} \sum_{i=1}^n x_i$  is the MLE of  $p$ . Thus using (†) :

$$\hat{p}_n \pm \zeta_{\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{p}_n)} \text{ is a } 100(1 - \alpha)\% \text{ C.I for } p .$$

Note that  $\text{Var}(\hat{p}_n) = \frac{p(1-p)}{n}$  . We have two choices:

1. replace  $p$  by  $\hat{p}_n$  in  $\text{Var}(\hat{p}_n)$  :

$$\hat{p}_n \pm \zeta_{\frac{\alpha}{2}} \frac{\sqrt{\hat{p}_n(1 - \hat{p}_n)}}{\sqrt{n}}$$

2. replace  $p(1 - p)$  in  $\text{Var}(\hat{p}_n)$  by  $\frac{1}{4}$  to find a conservatively large C.I for  $p$  :

$$\hat{p}_n \pm \zeta_{\frac{\alpha}{2}} \frac{1}{2\sqrt{n}}$$

**Example 5.3.** Suppose  $X_i \stackrel{iid}{\sim} \text{Ber}(p)$  ,  $i = 1, 2, \dots, n$  and we are interested in  $\Theta = p(1-p)$  , the variance.

An interesting property of MLE is the invariance , i.e. if  $\hat{\Theta}_n$  is the MLE of  $\Theta$  , then  $h(\hat{\Theta}_n)$  is the MLE of  $h(\Theta)$ . The invariance property then implies that:  $\hat{\Theta}_n = \hat{p}_n(1 - \hat{p}_n)$  is the MLE of  $p(1 - p) = \Theta$  .

The  $100(1 - \alpha)\%$  C.I for  $\Theta = p(1 - p)$  is  $\hat{\Theta}_n \pm \zeta_{\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\Theta}_n)}$

### 5.3 Small Sample Confidence Intervals

Unlike the large sample case, there is no general recipe like (\*) using which we can find an approximate pivotal quantity. In fact, there is on the paper, but only gives #MISSING in special cases.

To summarize , small sample probabilities are solved mostly case by case. A case of particular importance is the *normal case* . We will learn about the importance of this case when we discuss regression and ANOVA (Analysis of Variance) .

**Normal Case:**

Suppose  $X_i \stackrel{N}{\sim} (\overbrace{\mu}^{\text{of interest}}, \underbrace{\sigma^2}_{\text{nuisance}})$ ,  $i = 1, 2, \dots, n$  where  $n$ , the sample size is *NOT* large.

We learned that when  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $i = 1, 2, \dots, n$ :

$$\frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \stackrel{Exact}{\sim} N(0, 1) \quad (\ddagger)$$

This by itself is not useful since  $\sigma$  is not known. We discussed in previous section at length why we can replace  $\sigma$  by  $S$  when  $n$  is large enough. The formal justification is **not** applicable now since  $n$  is small, the intuitive justification still stands though.

Replacing  $\sigma$  with  $(\ddagger)$  changes the picture a bit. Given that  $S$  has the same spirit as  $\sigma$ , though in a small #MISSING the distribution of  $T = \frac{\bar{X}_n - \mu}{\frac{s}{\sqrt{n}}}$  still has a bell curve shape. The tails of the distribution, however, die out much more slowly than those of normal distribution. Heavier tails mean much more variability and this should perhaps be expected since by replacing  $\sigma$  by  $S$  which can be crude estimate when  $n$  is small, can add quite a bit to the variability. This is, of course, a intuitive argument. Following we present the sketch of a formal argument:

$$\begin{aligned} \text{step 1) } X_i &\stackrel{iid}{\sim} N(\mu, \sigma^2) \implies \bar{X}_n \sim N(\mu, \frac{\sigma^2}{n}) \\ &\implies \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1) \end{aligned}$$

$$\text{step 2) } X_i \stackrel{iid}{\sim} N(\mu, \sigma^2) \implies \frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{(n-1)}$$

**proof**

$$\begin{aligned}
 \sum_{i=1}^n (X_i - \mu)^2 &= \sum_{i=1}^n \left[ (X_i - \bar{X}_n) + (\bar{X}_n - \mu) \right]^2 \\
 &= \sum_{i=1}^n (X_i - \bar{X}_n)^2 + n(\bar{X}_n - \mu)^2 + 2(\bar{X}_n - \mu) \underbrace{\sum_{i=1}^n (X_i - \bar{X}_n)}_0 \\
 &= (n-1)S^2 + n(\bar{X}_n - \mu)^2
 \end{aligned}$$

by dividing both sides by  $\sigma^2$  we obtain:

$$\underbrace{\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2}_W = \underbrace{\frac{(n-1)S^2}{\sigma^2}}_U + \underbrace{\left( \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \right)^2}_V$$

Now note that :

$$\begin{aligned}
 X_i &\stackrel{iid}{\sim} N(\mu, \sigma^2) \implies \frac{X_i - \mu}{\sigma} \sim N(0, 1) \\
 &\implies \left( \frac{X_i - \mu}{\sigma} \right)^2 \sim \chi_1^2 \\
 &\implies \sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 \sim \chi_n^2
 \end{aligned}$$

(Exercise: Theorem 7.2 , page 356)

Thus  $W \sim \chi_n^2$  . On the other hand, using step 1:

$$\left( \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \right)^2 \sim \chi_1^2$$

step 3)    If     $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$     then     $\bar{X}_n \perp\!\!\!\perp S^2$



step 4)

$$\begin{aligned}
 m_w(t) &= \mathbb{E}e^{tW} \\
 &= \mathbb{E}[e^{t(U+V)}] \\
 &= \mathbb{E}[e^{tU} \cdot e^{tV}] \\
 &= \mathbb{E}[e^{tU}] \cdot \mathbb{E}[e^{tV}] \\
 &= m_u(t) + m_v(t)
 \end{aligned}$$

$U \amalg V$  using step 3

Thus

$$\begin{aligned}
 m_u(t) &= \frac{m_w(t)}{m_v(t)} \\
 &= \frac{(1-2t)^{-\frac{n}{2}}}{(1-2t)^{-\frac{1}{2}}} \\
 &= (1-2t)^{-\frac{n-1}{2}}
 \end{aligned}$$

which implies that  $U \sim \chi^2_{(n-1)}$

step 5) If  $Z \sim N(0,1)$ ,  $U \sim \chi^2_V$  and  $Z \amalg U$  then:

$$\frac{Z}{\sqrt{\frac{U}{V}}} \sim T_{n-1} \quad (\text{Exercise 7.30, page 367})$$

step 6)

$$T_{n-1} = \frac{\bar{X}_n - \mu}{\frac{s}{\sqrt{n}}} = \frac{\left(\frac{\bar{X}_n - \mu}{\frac{s}{\sqrt{n}}}\right)}{\sqrt{\frac{(n-1)s^2}{\sigma^2}} \frac{1}{\sqrt{n-1}}} = \frac{Z}{\sqrt{\frac{U}{V}}}$$

The pdf of  $T_v$  is :

$$f_{T_v}(t) = \frac{\Gamma(\frac{v+1}{2})}{\Gamma(\frac{v}{2})\sqrt{v\pi}} \left(1 + \frac{t^2}{v}\right)^{-\frac{v+1}{2}} \quad -\infty < t < +\infty$$

#MISSING GRAPH Lecture 5 - page 13

$$\mathbb{E}[T_v^r] = \begin{cases} 0 & \text{if } r < v \text{ and } r \text{ is odd} \\ v^{\frac{r}{2}} \cdot \frac{\Gamma(\frac{1+r}{2})\Gamma(\frac{v-r}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{v}{2})} & \text{if } r < v \text{ and } r \text{ is even} \end{cases}$$

Thus, if  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $i = 1, 2, \dots, n$  and  $\mu$  and  $\sigma^2$  are both unknown :

$$\bar{X}_n \pm t_{(n-1), \frac{\alpha}{2}} \frac{S}{\sqrt{n}}$$

provides a  $100(1 - \alpha)\%$  C.I for  $\mu$  where  $P(T_{(n-1)} > t_{(n-1), \frac{\alpha}{2}}) = \frac{\alpha}{2}$

## 5.4 Pivotal Quantity and Probability Integral Transform

Suppose  $X$  is a continuous random variable with *p.d.f*  $f$  and *cdf*  $F$ . Then  $F(X) \sim \text{Uniform}(0, 1)$  (Exercise)

This result is referred to as the **Probability Integral Transform**. Now suppose  $X_i \stackrel{iid}{\sim} F$ . Then :

$$\begin{aligned} F(X_i) \sim \text{Unif}(0, 1) &\implies -2 \ln F(X_i) \sim \chi_2^2 \\ &\implies -2 \sum_{i=1}^n \ln F(X_i) \sim \chi_{2n}^2 \\ &\implies -2 \sum_{i=1}^n \ln [1 - F(X_i)] \sim \chi_{2n}^2 \end{aligned}$$

There is hence a general recipe for finding a pivotal quantity when we have samples from continuous random variables. The usefulness of this pivotal quantity *depends* on the form of  $F$ , the *cdf* of  $X$ .

Suppose  $X_i \stackrel{iid}{\sim} \exp(\lambda)$ ,  $i = 1, 2, \dots, n$ , i.e:

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x > 0 \\ 0 & o/w \end{cases}$$

Then:

$$F(x) = \int_0^x f(t) dt = 1 - e^{-\lambda x}, \quad x > 0$$

and:

$$F(x) = \begin{cases} 1 - e^{-\lambda x} & x > 0 \\ 0 & o/w \end{cases}$$

Using the above discussion:

$$-2 \sum_{i=1}^n \ln F(X_i) \sim \chi^2_{2n} \quad \text{and} \quad -2 \sum_{i=1}^n \ln [1 - F(X_i)] \sim \chi^2_{2n}$$

for this example it is easier to work with the latter, i.e. :

$$\begin{aligned} 2 \sum_{i=1}^n \ln [1 - F(X_i)] &= -2 \sum_{i=1}^n \ln (e^{-\lambda X_i}) \\ &= 2\lambda \sum_{i=1}^n X_i = 2n\lambda \bar{X}_n \\ \text{so } \implies 2n\lambda \bar{X}_n &\sim \chi^2_{2n} \end{aligned}$$

Using the  $\chi^2$  table (Application 3, page 850-851) , we can find  $\chi^2_{(2n),0.025}$  and  $\chi^2_{(2n),0.975}$  such that:

$$P(\chi^2_{(2n),0.975} < 2n\lambda \bar{X}_n < \chi^2_{(2n),0.025}) = 0.95$$

Thus:

$$\left( \frac{\chi^2_{(2n),0.975}}{2n\bar{X}_n}, \frac{\chi^2_{(2n),0.025}}{2n\bar{X}_n} \right)$$

provides that a 95% C.I for  $\lambda$  . Note that  $\chi^2_{(2n),\alpha}$  is such that  $P(\chi^2_{(2n)} > \chi^2_{(2n),\alpha}) = \alpha$

#MISSING GRAPH LECTURE 5 - PAGE 15

## 6 Lecture 6

### 6.1 Small Sample Confidence Interval(general case):

We learned in the last lecture how to find C.I. for the population mean when the population distribution is normal. The two main pivotal quantities are:

$$(a) \quad \frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{(n-1)}^2 \quad \& \quad (b) \quad \frac{\bar{X}_n - \mu}{\frac{S}{\sqrt{n}}} \sim T_{(n-1)}$$

when  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  and  $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ .

The first result can be used to make a C.I. for  $\sigma^2$  and  $\sigma$  which the latter is used for making a C.I. for  $\mu$ .

We now consider the general case.

### 6.2 Probability Integral Transform(PIT)

Suppose  $X_i \stackrel{iid}{\sim} F_X$  and  $f$  is the pdf of  $X_i$ s:

$$X \sim F_X, Y = F_X(X)$$

$$\begin{aligned} F_Y(t) &= P(Y \leq t) = P(F_X(X) \leq t) \\ &= P(X \leq F_X^{-1}(t)) \\ &= F_X(F_X^{-1}(t)) = t \quad \text{for } 0 \leq t \leq 1 \end{aligned}$$

Thus  $F_Y(t) = \begin{cases} 0 & \text{if } t < 0 \\ t & \text{if } 0 \leq t < 1 \\ 1 & \text{if } 1 \leq t \end{cases}$

and hence  $Y \sim Unif(0,1)$ . This is called **Probability Integral Transform(PIT)**.

**Example 6.1.**  $X \sim Exp(\lambda)$ ,  $f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x > 0 \\ 0 & x \leq 0 \end{cases}, \lambda > 0.$

$$\begin{aligned}
F_x(x) &= P(X \leq x) = \int_{-\infty}^x f_x(t)dt = \int_0^x \lambda e^{-\lambda t} dt \\
&= -e^{-\lambda t} \Big|_0^x \\
&= 1 - e^{-\lambda x}
\end{aligned} \tag{1}$$

Now consider  $Y = F_x(x) = 1 - e^{-\lambda x}$  :

$$\begin{aligned}
F_Y(t) &= P(Y \leq t) = P(1 - e^{-\lambda x} \leq t) \\
&= P(e^{-\lambda x} \geq 1 - t) = P(X \leq -\frac{\ln(1-t)}{\lambda}) \\
&= F_x(-\frac{\ln(1-t)}{\lambda}) \\
&= 1 - e^{-\lambda(-\frac{\ln(1-t)}{\lambda})} \quad \text{using (1)} \\
&= 1 - e^{\ln(1-t)} = 1 - (1 - t) \\
&= t
\end{aligned}$$

Thus  $Y \sim \text{Unif}(0, 1)$ .

**Remark.** Using *PIT* we can essentially generate random numbers from any continuous distributions. In fact, suppose we want samples from cdf  $F$ . Then:

**Step 1:** Generate  $U_i \stackrel{iid}{\sim} \text{Unif}(0, 1)$  ,  $i = 1, 2, \dots, n$ .

**Step 2:**  $X_i = F^{-1}(U_i) \stackrel{iid}{\sim} F$  ,  $i = 1, 2, \dots, n$

This algorithm then works as long as we can generate uniform random numbers and  $F^{-1}$  can be explicitly found or well approximated.

**Example.**  $f_x(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$  ,  $-\infty < x < +\infty$  ( $X \sim N(0, 1)$ )

$$\text{Then: } F_x(x) = \int_{-\infty}^x f_x(t)dt = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt .$$

In this case,  $F^{-1}$  does not have an explicit nice form, but it can be well approximated.

**Remark.** A simple and useful transformation:

$$X \sim F \implies \overbrace{Y = F(X) \sim Unif(0,1)}^{P.I.T} \implies -2 \cdot \log(Y) \sim \chi_2^2$$

### 6.3 Pivotal Quantity

Suppose  $X_i \stackrel{iid}{\sim} F$ ,  $i = 1, 2, \dots, n$ . Define:

$$Y_i = F(X_i) \stackrel{iid}{\sim} Unif(0,1), \quad i = 1, 2, \dots, n$$

Now consider:

$$V_i = -2 \cdot \log(Y_i) \stackrel{iid}{\sim} \chi_{2n}^2, \quad i = 1, 2, \dots, n$$

Then:

$$\sum_{i=1}^n V_i \sim \chi_{2n}^2.$$

Having established the first two results, i.e. :

**Step 1:**  $X_i \sim F \implies Y_i = F(X_i) \sim Unif(0,1)$  (PIT)

**Step 2:**  $V_i = -2 \cdot \log(Y_i) \sim \chi_2^2$  (method of transformation)

The last result can be established using the method of moments:

$$\begin{aligned} m_{\sum_{i=1}^n V_i}(t) &= \mathbb{E}[e^{-t \sum_{i=1}^n V_i}] = \mathbb{E}\left[\prod_{i=1}^n e^{-t V_i}\right] \\ \prod_{i=1}^n V_i &\implies = \prod_{i=1}^n \mathbb{E}[e^{-t V_i}] = \prod_{i=1}^n m_{V_i}(t) \\ \text{identically distributed} &\implies = [m_V(t)]^n = [(1-2t)^{-\frac{2}{2}}]^n \\ &= (1-2t)^{-\frac{2n}{2}} \implies \sum_{i=1}^n V_i \sim \chi_{2n}^2 \end{aligned}$$

Then a pivotal quantity based on  $X_i \stackrel{iid}{\sim} F_\theta$ ,  $i = 1, 2, \dots, n$  is:

$$-2 \sum_{i=1}^n \log(F_\theta(X_i)) \sim \chi_{2n}^2 \quad (1)$$

**Example.**  $X_i \stackrel{iid}{\sim} \text{Exp}(\lambda)$  ,  $f_x(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$

$$F_x(x) = \int_{-\infty}^x f_x(t)dt = \int_0^x \lambda e^{-\lambda t} dt = 1 - e^{-\lambda x}$$

Now we notice that  $-2 \cdot \log(F) = -2 \cdot \log(1 - e^{-\lambda x})$  does not provide an useful form for the purpose of making a C.I for  $\lambda$ . There is a dual to (1) that is useful in this case, however:

$$\sum_{i=1}^n W_i = \sum_{i=1}^n -2\log(1 - F(X_i)) \sim \chi_{2n}^2 \quad (2)$$

This quickly follows from the fact that:

$$U \sim \text{Unif}(0,1) \implies 1 - U \sim \text{Unif}(0,1) \quad .$$

Using (2) we have:

$$\begin{aligned} \sum_{i=1}^n -2\log(1 - F(X_i)) &= \sum_{i=1}^n -2\log(e^{-\lambda X_i}) \\ &= 2\lambda \sum_{i=1}^n X_i = 2\lambda n \bar{X}_n \sim \chi_{2n}^2 \end{aligned}$$

Using the  $\chi^2$ -table (App.3 m page 850-851) , we find  $\chi_{2n,0.025}^2$  and  $\chi_{2n,0.975}^2$  such that:

$$P(\chi_{2n,0.975}^2 < 2\lambda n \bar{X}_n < \chi_{2n,0.025}^2) = 0.95$$

and hence:

$$P\left(\frac{\chi_{2n,0.975}^2}{2n\bar{X}_n} < \lambda < \frac{\chi_{2n,0.025}^2}{2n\bar{X}_n}\right) = 0.95$$

Thus:

$$\left(\frac{\chi_{2n,0.975}^2}{2n\bar{X}_n}, \frac{\chi_{2n,0.025}^2}{2n\bar{X}_n}\right) \quad \text{is a } \underline{95\%} \text{ confidence interval for } \lambda$$

#MISSING Graph Lecture 6 - page 36

## 6.4 Small Size Determination

Suppose we want to estimate the proportion of Canadian voters who are in favor of NDP and want our estimate to be one-percentage point from the actual population with 95% confidence. Define:

$$X = \begin{cases} 1 & \text{NDP} \\ 0 & \text{other parties} \end{cases} \quad \text{associated to each potential voter.}$$

We learned that to estimate the proportion of interest  $p = P(X = 1)$ , we can use  $\hat{p}_n = \frac{1}{n} \sum_{i=1}^n X_i$  from a random sample of size  $n$ . We further learned that if the sample size  $n$  is large enough, then:

$$\hat{p}_n \pm 1.96 \sqrt{\frac{p(1-p)}{n}}$$

is a 95% confidence interval for  $p$ . Thus the margin of error is  $\beta = 1.96 \sqrt{\frac{p(1-p)}{n}}$  which is controlled by  $n$ , the sample size. We should therefore choose  $n$  such that:

$$0.01 = 1.96 \sqrt{\frac{p(1-p)}{n}}$$

Given that  $p$  is unknown, we can either replace  $p$  by  $\hat{p}_n$  or take a conservative approach and replace  $p$  by  $\frac{1}{2}$  which maximizes  $p(1-p)$ . Thus we find:

$$n = \frac{p(1-p)\zeta_{\frac{\alpha}{2}}^2}{\beta^2} = \begin{cases} \frac{\hat{p}_n(1-\hat{p}_n)\zeta_{\frac{\alpha}{2}}^2}{\beta^2} & \text{replacing } p \text{ by } \hat{p}_n \\ \frac{\zeta_{\frac{\alpha}{2}}^2}{4\beta^2} & \text{replacing } p \text{ by } \frac{1}{2} \end{cases}$$

Taking the conservative approach, we have:

$$n = \frac{\zeta_{\frac{\alpha}{2}}^2}{4\beta^2} = \frac{(1.96)^2}{4(0.01)^2} = 9604$$

Likewise we can find the sample size formula for estimating the population mean with a given confidence  $1 - \alpha$  and margin of error  $\beta$ , we should in fact



solve the following equation for  $n$ :

$$\beta = \zeta_{\frac{\alpha}{2}}^2 \frac{\sigma}{\sqrt{n}} \quad \text{where } \sigma^2 \text{ is the population variance.}$$

We then find  $n = \frac{\zeta_{\frac{\alpha}{2}}^2 \sigma^2}{\beta^2}$  where  $\sigma^2$  should be estimated from a prior sample.

## 6.5 Sample Size Determination For Other Parameters

So far we only considered the population mean. Now consider a parameter  $\theta$ . In chapter 9 we learn about different methods of estimation, among them there is a method called the method of maximum likelihood (ML). Suppose  $\hat{\theta}_n$  is the maximum likelihood estimate (MLE) of  $\theta$ . Then under some reasonable conditions for a considerably large class of parametric distributions, we have:

$$\frac{\hat{\theta}_n - \theta}{\sqrt{\text{Var}(\hat{\theta}_n)}} \stackrel{app}{\sim} N(0, 1) \quad , \text{ for large } n$$

Thus:

$$\hat{\theta} \pm \underbrace{\zeta_{\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\theta}_n)}}_{\beta} \quad \text{is a } 100(1 - \alpha)\% \text{ C.I for } \theta$$

Let  $\beta = \zeta_{\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\theta}_n)}$ . In many interesting cases,  $\text{Var}(\hat{\theta}_n)$  is an explicit function of  $n$  and  $\sigma^2$ , the variance in the target population, say  $h(\sigma^2, n)$ . Then the sample size can be determined by the solution of the following equation:

$$h(\sigma^2, n) = \frac{\beta^2}{\zeta_{\frac{\alpha}{2}}^2}$$

Recall that for Bernoulli case, i.e.  $X_i = \begin{cases} 1 \\ 0 \end{cases}$  ,  $i = 1, 2, \dots, n$  :

$$h(\sigma^2, n) = \text{Var}(\hat{\theta}_n) = \text{Var}(\hat{p}_n) = \frac{\overbrace{p(1-p)}^{\sigma^2}}{n}$$

while for estimating the population mean:

$$h(\sigma^2, n) = \text{Var}(\hat{\theta}_n) = \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n} .$$

### 6.5.1 Sample Size Determination (Small Sample)

- Normal Case: We learned that if  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$  ,  $i = 1, 2, \dots, n$  :

$$\frac{\bar{X}_n - \mu}{\frac{S}{\sqrt{n}}} \sim T_{n-1} \implies \bar{X}_n \pm t_{\frac{\alpha}{2}, (n-1)} \frac{S}{\sqrt{n}} \quad \text{is a } 100(1 - \alpha)\% \text{ C.I for } \mu$$

We can therefore find sample size from the following equation:

$$\beta = t_{\frac{\alpha}{2}, (n-1)} \frac{S}{\sqrt{n}} \implies \boxed{n = \frac{S^2 t_{\frac{\alpha}{2}, (n-1)}^2}{\beta^2}}$$

Now note that the sample size determination based on large sample and small sample in the normal case had to  $N(0, 1)$  and  $T_{(n-1)}$  respectively. These distributions are both symmetric. As such in the sample size determination we only deal with the half length of the confidence intervals when the sample size is small and the population from which the samples are taken is not normal, the pivotal quantities do not necessarily have asymmetric distribution and hence the confidence interval do not have the form of  $\boxed{\hat{\theta}_n \pm \beta}$ .

In such cases we try to control the total length of the confidence interval.

**Example.**  $X_i \stackrel{iid}{\sim} \text{Exp}(\lambda)$  ,  $i = 1, 2, \dots, n$

We found that  $(\frac{\chi_{2n,0.975}^2}{2n\bar{X}_n}, \frac{\chi_{2n,0.025}^2}{2n\bar{X}_n})$ ; let  $C$  be the desired length for C.I for  $\lambda$ .

Then:

$$C = \frac{\chi_{2n, \frac{\alpha}{2}}^2 - \chi_{2n, 1-\frac{\alpha}{2}}^2}{2n\bar{X}_n}$$

represents the length of a C.I for  $\lambda$  based on a sample of size  $n$  with  $100(1 - \alpha)\%$  confidence.

### 6.5.2 Sample Size Determination(Two Sample Case)

So far we just confirmed ourselves to one population. We might, however, have two samples,  $X_1, X_2, \dots, X_m$  from the population of men with population mean  $\mu_M$ , and  $Y_1, Y_2, \dots, Y_n$  from the population of women with population mean  $\mu_W$ . Suppose the parameter of interest is  $\theta = \mu_M - \mu_W$ . Then the natural estimate of  $\theta$  is  $\hat{\theta} = \bar{X}_m - \bar{Y}_n$  and using the central limit theorem:

$$\frac{(\bar{X}_m - \bar{Y}_n) - (\mu_M - \mu_W)}{\sqrt{Var(\bar{X}_m - \bar{Y}_n)}} \stackrel{app}{\sim} N(0, 1) \quad \text{for large } m \text{ \& } n$$

Now:

$$Var(\bar{X}_m - \bar{Y}_n) = Var(\bar{X}_m) + Var(\bar{Y}_n) - 2Cov(\bar{X}_m, \bar{Y}_n)$$

Assuming that  $X$ s and  $Y$ s are independent (i.e  $Cov(\bar{X}_m, \bar{Y}_n) = 0$ ):

$$Var(\bar{X}_m - \bar{Y}_n) = Var(\bar{X}_m) + Var(\bar{Y}_n) = \frac{\sigma_M^2}{m} + \frac{\sigma_W^2}{n}$$

Therefore:

$$(\bar{X}_m - \bar{Y}_n) \pm \zeta_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_M^2}{m} + \frac{\sigma_W^2}{n}} \quad \text{is a } 100(1 - \alpha)\% \text{ C.I for } \mu_M - \mu_W.$$

To find the sample size we should solve:

$$\beta = \zeta_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_M^2}{m} + \frac{\sigma_W^2}{n}}$$

We should assume that  $\sigma_M^2$  &  $\sigma_W^2$  are known or estimated from prior samples, we will have one equation with two unknowns,  $m$  &  $n$ . In order to have a unique solution we need another equation. We often consider  $n = K \cdot m$ , where  $K$  is a known value as the second equation. Suppose  $C_M$  &  $C_W$  represent respectively, the cost of taking a sample from population of men and women. Then  $K \propto (\frac{C_W}{C_M})^{-1}$ . In case that  $C_W = C_M$ , we choose  $K = 1$ . Now:

$$\begin{cases} \beta = \zeta_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_M^2}{m} + \frac{\sigma_W^2}{n}} \\ n = Km \end{cases}$$

Solving the above system we find:

$$m = \left(\frac{\zeta_{\frac{\alpha}{2}}}{\beta}\right)^2 \cdot \left(\sigma_M^2 + \frac{\sigma_W^2}{K}\right)$$

For proportions:  $\sigma_M^2 = p_M(1-p_M)$  &  $\sigma_W^2 = p_W(1-p_W)$  Taking the conservative approach and replacing both  $p_M$  &  $p_W$  by  $\frac{1}{2}$  we find:

$$m = \left(\frac{\zeta_{\frac{\alpha}{2}}}{2\beta}\right)^2 \left(1 + \frac{1}{K}\right)$$

## 7 Lecture 7

### 7.1 Chapter 9 - Relative Efficiency

**Definition.** The relative efficiency of two unbiased estimators,  $\hat{\theta}_1$  and  $\hat{\theta}_2$ , is defined to be:

$$eff(\hat{\theta}_1, \hat{\theta}_2) = \frac{Var(\hat{\theta}_2)}{Var(\hat{\theta}_1)}$$

We learned that between the two unbiased estimators the one with smaller variance is closer to the target on the average, i.e. has smaller MSE.

We also learned that the length of confidence intervals for large sample size is controlled by the variance of the estimator; so, the smaller the variance, the shorter the confidence interval using that estimator is.

We now want to quantify the gain in using the estimator with smaller variance.

**Example.** Suppose  $X_i \stackrel{iid}{\sim} f$ ,  $i = 1, 2, \dots, n$  where  $f$  is a symmetric pdf.

The mean and median of  $f$  are the same, say  $\mu$ . Given that  $f$  is symmetric, we can use:  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  (the sample average), or:

$$M_n = \begin{cases} X_{\frac{n+1}{2}} & \text{if } n \text{ is odd} \\ \frac{1}{2}[X_{\frac{n}{2}} + X_{\frac{n}{2}+1}] & \text{if } n \text{ is even} \end{cases} \quad \text{where } X_{(1)} < X_{(2)} < \dots < X_{(n)} \text{ are the order statistics.}$$

We learned that  $Var(\hat{X}_n) = \frac{\sigma^2}{n}$  where  $\sigma^2$  is the population variance, i.e. :

$$\sigma^2 = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx$$

It can be shown (beyond the scope of this course), that:

$$Var(M_n) \approx \frac{1}{4 \cdot [f(\mu)]^2 n} \quad \text{for large } n$$

For instance if  $f$  is Normal, i.e.  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , then  $f(\mu) = \frac{1}{\sigma\sqrt{2\pi}}$  and hence:

$$Var(M_n) = \frac{2\pi}{4} \cdot \frac{\sigma^2}{n}$$

Thus:

$$eff(\bar{X}_n, M_n) = \frac{Var(M_n)}{Var(\bar{X}_n)} = \frac{\frac{2\pi}{4} \cdot \frac{\sigma^2}{n}}{\frac{\sigma^2}{n}} = \frac{2\pi}{4} = 1.57$$

This then essentially means that if you can make a confidence interval of a given length using  $\bar{X}_n$  with 100 observations, to make a confidence interval of the same length for  $\mu$  using  $M_n$ , you need  $100 \times 1.57 = 157$  observations.

**Example** (9.1, page 446).  $Y_i \stackrel{iid}{\sim} Unif(0, \theta)$ ,  $i = 1, 2, \dots, n$ ,  $\theta > 0$  and  $\theta$  is unknown.

Consider  $\hat{\theta}_1 = 2\bar{Y}_n$  and  $\hat{\theta}_2 = (\frac{n+1}{n})Y(n)$  where  $Y(n) = \max_{1 \leq i \leq n} Y_i$

For  $\hat{\theta}_1$ :

$$\begin{aligned} \mathbb{E}(\hat{\theta}_1) &= \mathbb{E}(2\bar{Y}_n) = 2\mathbb{E}(\bar{Y}_n) \\ &= 2\mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \frac{2}{n} \sum_{i=1}^n \mathbb{E}(Y_i) \\ &= \frac{2}{n} \sum_{i=1}^n \frac{\theta}{2} = \frac{2}{n} \cdot n \cdot \frac{\theta}{2} \\ &= \theta \end{aligned}$$

$$\begin{aligned} Var(\hat{\theta}_1) &= Var(2\bar{Y}_n) = 4Var(\bar{Y}_n) \\ &= 4 \cdot \frac{Var(Y)}{n} = 4 \cdot \frac{\sigma_Y^2}{n} \\ \sigma_Y^2 &= Var(Y) = \mathbb{E}(Y^2) - [\mathbb{E}(Y)]^2 \\ \mathbb{E}(Y^2) &= \int_{-\infty}^{+\infty} y^2 f_Y(y) dy = \int_0^\theta y^2 \cdot \frac{dy}{\theta} \\ &= \frac{1}{3\theta} y^3 \Big|_0^\theta = \frac{\theta^3}{3\theta} = \frac{\theta^2}{3} \\ \sigma_Y^2 &= \mathbb{E}(Y^2) - [\mathbb{E}(Y)]^2 = \frac{\theta^2}{3} - \left[\frac{\theta}{2}\right]^2 \\ &= \frac{\theta^2}{3} - \frac{\theta^2}{4} = \frac{\theta^2}{12}, \quad \text{then:} \\ Var(\hat{\theta}_1) &= \frac{4\sigma_Y^2}{n} = \frac{4\frac{\theta^2}{12}}{n} = \frac{\theta^2}{3n} \end{aligned}$$

For  $\hat{\theta}_2$ :

$$F_{Y(n)}(t) = P(Y(n) \leq t) = P(Y_1 \leq t, Y_2 \leq t, \dots, Y_n \leq t) \quad \text{by } \prod_{i=1}^n Y_i :$$

$$\begin{aligned}
&\implies \prod_{i=1}^n P(Y_i \leq t) = \prod_{i=1}^n F_{Y_i}(t) \\
&= [F_Y(t)]^n \quad \therefore \text{ identically distributed.} \\
&\text{Thus } d_{Y(n)}(t) = \frac{d}{dt} F_{Y(n)}(t) = \frac{d}{dt} F_Y^n(t) = n f_Y(t) F_Y^{n-1}(t) \\
&f_{Y(n)}(t) = \begin{cases} n \frac{1}{\theta} \left(\frac{t}{\theta}\right)^{n-1} & 0 < t < \theta \\ 0 & \text{otherwise} \end{cases}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}(\hat{\theta}_2) &= \mathbb{E}\left[\left(\frac{n+1}{n}\right)Y(n)\right] = \left(\frac{n+1}{n}\right)\mathbb{E}(Y(n)) \\
&= \left(\frac{n+1}{n}\right) \int_0^\theta y \cdot n \cdot \frac{1}{\theta} \left(\frac{y}{\theta}\right)^{n-1} dy \\
&= \left(\frac{n+1}{n}\right) \cdot \frac{n}{\theta^n} \int_0^\theta y^n dy \\
&= \left(\frac{n+1}{n}\right) \frac{n}{\theta^n} \left[ \frac{1}{n+1} y^{n+1} \right]_0^\theta \\
&= \left(\frac{n+1}{n}\right) \frac{n}{\theta^n} \frac{\theta^{n+1}}{n+1} \\
&= \theta
\end{aligned}$$

$$\begin{aligned}
\text{Var}(\hat{\theta}_2) &= \text{Var}\left(\frac{n+1}{n}Y(n)\right) = \left(\frac{n+1}{n}\right)^2 \text{Var}(Y(n)) \\
\text{Var}(Y(n)) &= \mathbb{E}(Y(n)^2) - [\mathbb{E}(Y(n))]^2 \\
\mathbb{E}(Y(n)^2) &= \int_0^\theta y^2 \cdot n \cdot \frac{1}{\theta} \left(\frac{y}{\theta}\right)^{n-1} dy \\
&= \frac{n}{\theta^n} \int_0^\theta y^{n+1} dy \\
&= \frac{n}{\theta^n} \cdot \frac{1}{n+2} \cdot \theta^{n+2} = \frac{n\theta^2}{n+2} \\
\text{Var}(Y(n)) &= \frac{n\theta^2}{n+2} - \left[\frac{n}{n+1}\theta\right]^2 \\
&= n\theta^2 \left( \frac{1}{n+2} - \frac{n}{(n+1)^2} \right) \\
&= \frac{n\theta^2}{(n+2)(n+1)^2} \quad \text{thus:} \\
\text{Var}(\hat{\theta}_2) &= \left(\frac{n+1}{n}\right)^2 \cdot \frac{n\theta^2}{(n+2)(n+1)^2} = \frac{\theta^2}{n(n+2)}
\end{aligned}$$

Thus:

$$\text{eff}(\hat{\theta}_1, \hat{\theta}_2) = \frac{\text{Var}(\hat{\theta}_2)}{\text{Var}(\hat{\theta}_1)} = \frac{\frac{\theta^2}{n(n+2)}}{\frac{\theta^2}{3n}} = \frac{3}{n+2} \rightarrow 0 \text{ as } n \rightarrow \infty$$

Note that  $\text{eff}(\hat{\theta}_1, \hat{\theta}_2) < 1$  for  $n \geq 2$ . This means that  $\hat{\theta}_2$  is more efficient than  $\hat{\theta}_1$  for  $n \geq 2$ .

We also notice that the efficiency gap increases as the sample size  $n$  increases and for large values of  $n$ , the efficiency tends to zero.

## 7.2 Consistency

**Definition** (Consistent Estimator). We say  $\hat{\theta}_n$  is a consistent estimator of  $\theta$  if  $\hat{\theta}_n \xrightarrow{p} \theta$  as  $n \rightarrow \infty$ ; i.e:

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| > \epsilon) = 0, \quad \forall \epsilon > 0 \quad (\dagger)$$

Consistency essentially means "being right-headed". It essentially says that if we have all the population, our procedure,  $\hat{\theta}_n$ , sizes the target.

Note that  $(\dagger)$  is equivalent to :

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| \leq \epsilon) = 1, \quad \forall \epsilon > 0$$

Now this definition can be compared with the notion of the limit of a sequence of real numbers.

$$\lim_{n \rightarrow \infty} a_n = a \quad \text{iff} \quad \forall \epsilon > 0 \exists N(\epsilon) \ni |a_n - a| < \epsilon \quad \text{if} \quad n \geq N(\epsilon)$$

Now since that  $\hat{\theta}_n$  is a random variable no matter how large  $n$  is, there is always a chance that  $|\hat{\theta}_n - \theta| > \epsilon$ . This chance, however, tends to zero as  $n \rightarrow \infty$ .

**Example.**  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$ ,

$$X_i = \begin{cases} 1 & p \\ 0 & 1-p \end{cases}, \quad i = 1, 2, \dots, n$$

$\hat{p}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . We want to show that:

$$\lim_{n \rightarrow \infty} P(|\hat{p}_n - p| > \epsilon) = 0, \quad \forall \epsilon > 0.$$

Compare  $P(|\hat{p}_n - p| > \epsilon)$  with Tchebyshev's Inequality:

$$P(|X - \mathbb{E}(X)| > K \underbrace{\sqrt{\text{Var}(X)}}_{\sigma}) \leq \frac{1}{K^2}$$



$X$  is replaced by  $\hat{p}_n$ ,  $\mu_X = \mathbb{E}(X)$  by  $p$  and  $K\sigma$  by  $\epsilon$ .

Note that  $\mathbb{E}(\hat{p}_n) = p$  so everything is in order for using Tchebyshev's Inequality.

Now  $\epsilon = K\sigma_X$  implies that  $K = (\frac{\sigma_X}{\epsilon})^{-1}$  and given that  $X$  is replaced by  $\hat{p}_n$ , we should have  $K = (\frac{\sigma_{\hat{p}_n}}{\epsilon})^{-1}$ . Thus:

$$\begin{aligned} P(|\hat{p}_n - p| > \epsilon) &\leq \frac{1}{(\frac{\sigma_{\hat{p}_n}}{\epsilon})^{-2}} = \frac{\sigma_{\hat{p}_n}^2}{\epsilon^2} \\ P(|\hat{p}_n - p| > \epsilon) &\leq \frac{Var(\hat{p}_n)}{\epsilon^2} \\ \Rightarrow Var(\hat{p}_n) &= Var(\frac{1}{n} \sum_{i=1}^n X_i) = \frac{p(1-p)}{n} \end{aligned}$$

$$\text{Therefore} \quad P(|\hat{p}_n - p| > \epsilon) \leq \frac{p(1-p)}{n\epsilon^2} \leq \frac{1}{4n\epsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty \quad \because p(1-p) \leq \frac{1}{4} \quad (\ddagger)$$

$$\text{Thus} \quad \lim_{n \rightarrow \infty} P(|\hat{p}_n - p| > \epsilon) = 0, \quad \forall \epsilon > 0$$

Note further that we can let  $\epsilon$  tend to zero as  $n \rightarrow \infty$ , i.e.  $\epsilon_n$  depend on  $n$  and  $\epsilon_n \rightarrow 0$  as  $n \rightarrow \infty$ . Using  $(\ddagger)$ :

$$\lim_{n \rightarrow \infty} P(|\hat{p}_n - p| > \epsilon_n) \leq \lim_{n \rightarrow \infty} \frac{1}{4n\epsilon_n^2}$$

Let  $\epsilon_n = \frac{\log(n)}{\sqrt{n}}$ , then:

$$\lim_{n \rightarrow \infty} P(|\hat{p}_n - p| > \frac{\log(n)}{\sqrt{n}}) \leq \lim_{n \rightarrow \infty} \frac{1}{4n(\frac{\log(n)}{\sqrt{n}})^2} = \lim_{n \rightarrow \infty} \frac{1}{4(\log(n))^2} = 0$$

This actually gives us an idea at what rate  $|\hat{p}_n - p| \xrightarrow{P} p$ . Note  $\epsilon_n = \frac{\log(n)}{\sqrt{n}}$  as long as  $\log(n) \rightarrow \infty$ , no matter how slow, we still have the same result. This then suggests that perhaps  $|\hat{p}_n - p|$  tends to zero in probability at the same rate as  $\frac{1}{\sqrt{n}}$ .

Suppose  $X_1, \dots, X_n$  have the same mean  $\mu$  and variance  $\sigma^2$ . Suppose further that  $Cov(X_i, X_j) = 0$   $i \neq j$ . Then  $\bar{X}_n \xrightarrow{P} \mu$ , i.e.  $\bar{X}_n$  is a consistent estimator of  $\mu$ , the population mean. Like the previous case:

$$P(|\bar{X}_n - \mu| > \epsilon) \leq \frac{1}{(\frac{\epsilon}{\frac{\sigma}{\sqrt{n}}})^2} = \frac{\sigma^2}{n\epsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty \quad (1)$$

Note that  $\epsilon = K \sqrt{\text{Var}(\bar{X}_n)} = K \sqrt{\frac{\sigma^2}{n}}$  and hence  $K = \frac{\epsilon}{\sqrt{\frac{\sigma^2}{n}}}$ . Then using Tchebyshev's Inequality we obtain (1) :

$$P(|\bar{X}_n - \mu| > \epsilon) \leq \frac{\sigma^2}{n\epsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty \quad \forall \epsilon > 0$$

Thus:

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| > \epsilon) = 0 \quad , \quad \forall \epsilon > 0$$

meaning that  $\bar{X}_n \xrightarrow{P} \mu$  , i.e.  $\bar{X}_n$  is a consistent estimator of  $\mu$  .

The same approach cannot be used to show that  $\delta_n^2 \xrightarrow{P} \sigma^2$  ( We need the law of large numbers(Kolmogorov's result)).

## 8 Lecture 8

### 8.1 Consistency

Consistency is the minimal property that an estimator is expected to possess. Consistency essentially means having right-headed ; in the sense that if "all" the population's information is available , the estimator produces the exact target. Recall once again:

**Definition:** Suppose  $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n)$  is an estimator of  $\theta$  . We say  $\hat{\theta}_n$  is a consistent estimator of  $\theta$  if  $\hat{\theta}_n \xrightarrow{P} \theta$  as  $n \rightarrow \infty$  , i.e:

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| > \epsilon) = 0 \quad , \quad \forall \epsilon > 0 \quad .$$

In lecture 7 we used Tchbyshev's inequality to establish consistency.

Markov's Inequality is an important tool in establishing consistency. In fact , Tchbyshev's inequality is a special case of Markov's inequality. It is often more straight forward to use Markov's inequality.

### 8.2 Markov's Inequality

Let  $X$  be a random variable and  $g$  a non-negative function. Then:

$$P(g(X) \geq \lambda) \leq \frac{\mathbb{E}[g(X)]}{\lambda} \quad , \quad \forall \lambda > 0 \quad .$$

Using Markov's inequality we have:

$$P(|\hat{\theta}_n - \theta| > \epsilon) \leq \frac{\mathbb{E}[|\hat{\theta}_n - \theta|]}{\epsilon} \quad (\dagger)$$

To establish consistency it then suffices to show that the upper bound of the above inequality tends to zero as  $n \rightarrow \infty$  .

Note that  $(\dagger)$  follows from Markov's inequality if we define  $g(x) = |x - \theta|$  . Note also that our random variable is  $\hat{\theta}_n$  .

To apply (†) , we need to find  $\mathbb{E}[|\hat{\theta}_n - \theta|]$  which is not always easy. We however have:

$$P(|\hat{\theta}_n - \theta| > \epsilon) = P(|\hat{\theta}_n - \theta|^2 > \epsilon^2) \quad \xrightarrow{\text{Markov's Ineq.}} \quad \frac{\mathbb{E}[|\hat{\theta}_n - \theta|^2]}{\epsilon^2}$$

and thus:

$$P(|\hat{\theta}_n - \theta| > \epsilon) \leq \frac{MSE(\hat{\theta}_n)}{\epsilon^2} = \frac{Var(\hat{\theta}_n) + Bias^2(\hat{\theta}_n)}{\epsilon^2} \quad (\ddagger)$$

where

$$MSE(\hat{\theta}_n) = \mathbb{E}[(\hat{\theta}_n - \theta)^2] = Var(\hat{\theta}_n) + \overbrace{[\mathbb{E}(\hat{\theta}_n) - \theta]^2}^{Bias(\hat{\theta}_n)}$$

Now (‡) is often much easier to use since Variance and Bias of an estimator are often hard to find.

**Theorem** (Slight Generalization of Theorem 9.1 , P450). *Suppose  $\hat{\theta}_n$  is an estimator of  $\theta$  . Then  $\hat{\theta}_n \xrightarrow{P} \theta$  if  $MSE(\hat{\theta}_n) \rightarrow 0$  as  $n \rightarrow \infty$ . In otherwords,  $\hat{\theta}_n$  is a consistent estimator of  $\theta$  if  $MSE(\hat{\theta}_n) \rightarrow 0$  as  $n \rightarrow \infty$  .*

*Proof.* Using (‡) we have:

$$P(|\hat{\theta}_n - \theta| > \epsilon) \leq \frac{MSE(\hat{\theta}_n)}{\epsilon^2} \rightarrow 0 \quad \text{as } n \rightarrow \infty \quad \forall \epsilon > 0 \quad \because MSE(\hat{\theta}_n) \rightarrow 0 \text{ as } n \rightarrow \infty$$

□

**Corollary 8.1.** *Let  $\hat{\theta}_n$  be an unbiased estimator of  $\theta$ . Suppose  $Var(\hat{\theta}_n) \rightarrow 0$  as  $n \rightarrow \infty$  . Then  $\hat{\theta}_n \xrightarrow{P} \theta$  , i.e.  $\hat{\theta}_n$  is a consistent estimator of  $\theta$  .*

*Proof.*

$$MSE(\hat{\theta}_n) = Var(\hat{\theta}_n) + Bias^2(\hat{\theta}_n) = Var(\hat{\theta}_n) \rightarrow 0 \text{ as } n \rightarrow \infty$$

Note that the  $Bias(\hat{\theta}_n) = \mathbb{E}(\hat{\theta}_n) - \theta = 0$  if  $\hat{\theta}_n$  is an unbiased estimator of  $\theta$  , i.e.  $\mathbb{E}(\hat{\theta}_n) = \theta$ . □

**Example 8.1.**  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$  ,  $i = 1, 2, \dots, n$

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

$$\mathbb{E}(\hat{p}_n) = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \frac{1}{n} \sum_{i=1}^n p = p$$

Then  $\mathbb{E}(\hat{p}_n) = p$  . i.e.  $\hat{p}_n$  is an unbiased estimator of  $p$  .

$$\text{Var}(\hat{p}_n) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} \sum_{i=1}^n p(1-p) = \frac{p(1-p)}{n}$$

$$\text{Now: } \text{Var}(\hat{p}_n) = \frac{p(1-p)}{n} \rightarrow 0 \text{ as } n \rightarrow \infty$$

Thus using corollary 8.1 ,  $\hat{p}_n \xrightarrow{P} p$  , i.e.  $\hat{p}_n$  is a consistent estimator of  $p$  .

**Example 8.2.** Suppose  $X_1, \dots, X_n$  are independent and identically distributed random variables with the common mean value  $\mu$  and common variance  $\sigma^2$ .

Then:

$$\mathbb{E}(\bar{X}_n) = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \frac{1}{n} \sum_{i=1}^n \mu = \mu$$

$$\text{Var}(\bar{X}_n) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}$$

$$\text{Thus: } \text{MSE}(\bar{X}_n) = \text{Var}(\bar{X}_n) + \overbrace{\text{Bias}^2(\bar{X}_n)}^0 = \frac{\sigma^2}{n} \rightarrow 0 \text{ as } n \rightarrow \infty$$

and hence  $\bar{X}_n \xrightarrow{\mu}$  using corollary 8.1 , i.e.  $\bar{X}_n$  is a consistent estimator of  $\mu$ .

**Remark 8.1.** The conclusion of 2<sup>nd</sup> example remains intact if the independence assumption is replaced by orthogonality , i.e.  $\text{Cov}(X_i, X_j) = 0$  ,  $i \neq j$ .

**Corollary 8.2.** Suppose  $\hat{\theta}_n$  is an asymptotically unbiased estimator of  $\theta$  , i.e.  $\lim_{n \rightarrow \infty} \mathbb{E}(\hat{\theta}_n) = \theta$  .

Suppose further that  $\text{Var}(\hat{\theta}_n) \rightarrow 0$  as  $n \rightarrow \infty$  . Then  $\hat{\theta}_n \xrightarrow{P} \theta$  , i.e.  $\hat{\theta}_n$  is a consistent estimator of  $\theta$

*Proof.*

$$\lim_{n \rightarrow \infty} \text{MSE}(\hat{\theta}_n) = \lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}_n) + \lim_{n \rightarrow \infty} \text{Bias}^2(\hat{\theta}_n) = 0 + \overbrace{[\lim_{n \rightarrow \infty} (\mathbb{E}(\hat{\theta}_n) - \theta)]^2}^0 = 0$$

The desired result then follows from the above theorem.  $\square$

**Remark.** The above result tells us that unbiasedness is NOT necessary for consistency.

**Example 8.3.** Suppose  $X_1, \dots, X_n$  from a random sample from a population with the mean  $\mu$  and variance  $\sigma^2$ . We want to estimate  $\sigma^2$ .

We showed that  $\mathbb{E}(S_n^2) = \sigma^2$  where :

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

We want to show that  $S_n^2 \xrightarrow{P} \sigma^2$ , i.e.  $S_n^2$  is a consistent estimator of  $\sigma^2$ . Note that using Markov's inequality we have:

$$P(|S_n^2 - \sigma^2| > \epsilon) \leq \frac{\text{Var}(S_n^2)}{\epsilon^2}$$

We need to show that  $\text{Var}(S_n^2) \rightarrow 0$  as  $n \rightarrow \infty$ . To do this, we need to find  $\text{Var}(S_n^2)$  in terms of the moments of the population. We therefore require conditions on the 4<sup>th</sup> moment of the population from which the samples were taken. Below we give a different approach that is much easier to apply and require lesser assumptions, but much more base.

### 8.3 Kolmogorov's Law of Large Numbers(LLN)

Suppose  $X_{n=1}^\infty$  is a sequence of *iid* random variables with common mean  $\mu$ .

Then:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mu = \mathbb{E}(X)$$

**Remark.** Kolmogorov's theorem is actually much stronger than this. It established a

stronger notion of convergence. The complete form has two sides. It also shows that if  $\bar{X}_n$  converges to a constant, since  $C$ , in that stronger notion of convergence, then:  $\mathbb{E}(|X|) < \infty$  &  $C = \mathbb{E}(X)$ .

**Corollary.** Suppose  $\{X_n\}_{n=1}^{\infty}$  is a sequence of iid random variables with the common  $K^{th}$ -moment  $\mu_K$ , i.e.  $\mathbb{E}(X^K) = \mu_K$ , for some  $K \in \mathbb{N}$ , then:

$$\frac{1}{n} \sum_{i=1}^n X_i^K \xrightarrow{P} \mu_K = \mathbb{E}(X^K)$$

This corollary follows from Kolmogorov's theorem immediately upon defining  $Y_i = X_i^K$ .

Note that if  $\mathbb{E}(X^K) < \infty$ , then  $\mathbb{E}(X^r) < \infty \quad \forall 0 \leq r \leq K$ . This then means that if  $X_{n=1}^{\infty}$  is a sequence iid random variables with the common  $K^{th}$ -moment  $\mu_K$ , then:

$$\frac{1}{n} \sum_{i=1}^n X_i^r \xrightarrow{P} \mu_K = \mathbb{E}(X^r) \quad \forall 0 \leq r \leq K$$

We also need the following theorem which is essentially theorem 9.2, page 451 of the textbook.

**Theorem 8.1** (Theorem 9.2, page 451). Suppose  $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n)$  and  $\hat{\mathcal{C}}_n = \hat{\mathcal{C}}(X_1, \dots, X_n)$  are consistent estimators of  $\theta$  and  $\mathcal{C}$ , respectively, i.e.  $\hat{\theta}_n \xrightarrow{P} \theta$  and  $\hat{\mathcal{C}}_n \xrightarrow{P} \mathcal{C}$ .

a)  $\hat{\theta}_n \hat{\mathcal{C}}_n \xrightarrow{P} \theta \mathcal{C}$

b)  $\hat{\theta}_n \pm \hat{\mathcal{C}}_n \rightarrow \theta \pm \mathcal{C}$

c)  $\frac{\hat{\theta}_n}{\hat{\mathcal{C}}_n} \rightarrow \frac{\theta}{\mathcal{C}} \quad \text{provided that } \mathcal{C} \neq 0, \mathcal{C} \neq 0$

d)  $g(\hat{\theta}_n) \rightarrow g(\theta) \quad \text{if } g(.) \text{ is a continuous function}$

Part(d) of the above theorem is called **Continuous Mapping Theorem**.

Now to establish consistency of  $S_n^2$ , we first establish  $S_{n,*}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \xrightarrow{P} \sigma^2$ .

**Step 1:**

$$\begin{aligned}
 S_{n,*}^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \left[ \sum_{i=1}^n X_i^2 - 2\bar{X}_n \sum_{i=1}^n X_i + n\bar{X}_n^2 \right] \\
 &= \frac{1}{n} \left[ \sum_{i=1}^n X_i^2 - 2\bar{X}_n(n\bar{X}_n) + n\bar{X}_n^2 \right] \\
 &= \frac{1}{n} \left[ \sum_{i=1}^n X_i^2 - 2n\bar{X}_n^2 + n\bar{X}_n^2 \right] \\
 &= \frac{1}{n} \left[ \sum_{i=1}^n X_i^2 - n\bar{X}_n^2 \right] \\
 &= \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2
 \end{aligned}$$

**Step 2:**

Using Kolmogorov's Theorem:

$$\begin{aligned}
 \text{a) } &\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{P} \mathbb{E}(X^2) \\
 \text{b) } &\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mathbb{E}(X)
 \end{aligned}$$

**Step 3:**

Using Step 2 and continuous mapping theorem (Theorem 9.2 (d))

$$\bar{X}_n^2 \rightarrow [\mathbb{E}(X)]^2$$

**Step 4:**

Using Step 1,2,3 and Theorem 9.2 (b) we have:

$$S_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 \xrightarrow{P} \mathbb{E}(X^2) - [\mathbb{E}(X)]^2 = \text{Var}(X) = \sigma^2$$

Thus  $S_{n,*}^2 \xrightarrow{P} \sigma^2$ , i.e.  $S_{n,*}^2$  is a consistent estimator of  $\sigma^2$ .

Next we note that:



$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \left(\frac{n}{n-1}\right) \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \left(\frac{n}{n-1}\right) S_{n,*}^2$$

Using Theorem 9.2 (a):

$$\lim_{n \rightarrow \infty} \left(\frac{n}{n-1}\right) = 1 \quad \& \quad S_{n,*}^2 \xrightarrow{P} \sigma^2$$

Then Theorem 9.2 (a) implies that:

$$S_n^2 = \left(\frac{n}{n-1}\right) S_{n,*}^2 \xrightarrow{P} 1 \cdot \sigma^2 = \sigma^2$$

i.e.  $S_n^2$  is a consistent estimator of  $\sigma^2$ .

**Remark.** Suppose  $P(X_n = C_n) = 1$ ,  $n = 1, 2, \dots$  where  $\{C_n\}_{n=1}^{\infty}$  is a sequence of real numbers such that  $\lim_{n \rightarrow \infty} C_n = C$ . Then  $X_n \xrightarrow{P} C$ . The proof of this result is as follow:

*Proof.*

$$\lim_{n \rightarrow \infty} C_n = C \text{ i.e. } \forall \epsilon > 0 \exists N(\epsilon) \in \mathbb{N} \ni |C_n - C| < \epsilon, \forall n \geq N(\epsilon)$$

Now suppose  $\epsilon > 0$  is given, then:

$$P(|X_n - C| > \epsilon) = P(|C_n - C| > \epsilon) = 0 \quad \text{if } n \geq N(\epsilon)$$

$$\text{Thus } \lim_{n \rightarrow \infty} P(|X_n - C| > \epsilon) = 0, \quad \forall \epsilon > 0 \quad \square$$

**Question:** Why couldn't we use Kolmogorov's theorem directly to establish consistency of  $S_{n,*}^2$ ? In other words, couldn't we define  $Z_i = (X_i - \bar{X}_n)^2$  and hence  $S_{n,*}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \sum_{i=1}^n Z_i$  and then apply Kolmogorov's theorem? The answer is that  $Z_i$ s are not independent. Note that  $\sum_{i=1}^n (X_i - \bar{X}_n) = 0$ .

## 8.4 Sufficiency

Sufficiency is essentially comparison. Sufficiency is one of the main pillars of the likelihood Inference.

As the following diagram shows the likelihood inference has three main components: the observable quantities, samples, the unobservable quantities, the

unknown parameters to be estimated , and a parametric distribution that links the observables to unobservables .

**Example:**  $f_{\theta}(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$  where  $\theta = (\mu, \sigma^2)$

#MISSING GRAPH Lecture 8 - Page 57

## 9 Lecture 9

### 9.1 Sufficiency

Suppose  $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m < n$  is a map from  $\mathbb{R}^n$  to  $\mathbb{R}^m$  and  $X = (X_1, X_2, \dots, X_n)$  and  $Y = (Y_1, Y_2, \dots, Y_n)$  are two  $n$ -dimensional random vectors. Then  $\tilde{X}$  &  $\tilde{Y}$  are called T-similar if:

$$P_{\tilde{X}|T=t}(u|t, \theta) = P_{\tilde{Y}|T=t}(u|t, \theta) \quad \forall u \text{ and } t$$

**Definition 9.1** (T-similar). A realization of  $\tilde{X}$ , say  $\tilde{x}$ , and a realization of  $\tilde{Y}$  say  $\tilde{y}$ , are called T-similar if:

1.  $\tilde{X}$  and  $\tilde{Y}$  are T-similar .
2.  $T(\tilde{x}) = T(\tilde{y})$

What do we expect from a good comparison?

Suppose  $\theta$  is the unknown parameter of interest. We want to estimate  $\theta$  using  $\tilde{X} = (X_1, \dots, X_n)$ . Now  $T_n = T(\tilde{X}) = T(X_1, \dots, X_n)$  is a good comparison if:

- a)  $T_n$  can preserve all the pertinent "information" in  $\tilde{X} = (X_1, \dots, X_n)$  to  $\theta$
- b) if  $\tilde{X}^* = (X_1^*, \dots, X_n^*)$  is the original sample, for any given value of  $T_n$ , say  $t$ , we can generate a  $T_n$ -similar sample of  $\tilde{x}$ s.

In order to generate a  $T_n$ -similar sample  $P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T_n = t)$  should be free from any unknown parameter.

As per retaining pertinent information in the data to the unknown parameter  $\theta$ , given that the link between the data and the unknown parameter(s) is the joint distribution:

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n; \theta) = P_{\theta}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

any possible information in the data about  $\theta$  should be in  $P_{\theta}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$

As such  $T_n = T(X_1, \dots, X_n)$  can preserve all the pertinent information if  $P_\theta(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$ . The Fisher-Neyman Factorization Theorem shows that this proportional is indeed a characterization of sufficient statistics. Let's dig into this a bit more. Note that:

$$P_\theta(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = P_\theta(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T_n = t) P_\theta(T_n = t) .$$

If  $P_\theta(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T_n = t)$  is actually free from  $\theta$ , then:

$$P_\theta(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) \propto P_\theta(T_n = t) \stackrel{\text{def.}}{=} g(t; \theta)$$

where the proportionality constant is a function of  $x = (x_1, x_2, \dots, x_n)$ , the observed sample. A formal definition then emerges.

**Definition.** Let  $X_1, \dots, X_n$  be a random sample from a distribution with an unknown parameter  $\theta$ . A statistic  $T_n = T(X_1, \dots, X_n)$  is called **sufficient** for  $\theta$  if the conditional distribution of  $(X_1, \dots, X_n)$  given  $T_n$  does not depend on  $\theta$ .

**Example.**  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$ ,  $i = 1, 2, \dots, n$ .

$$P_p(X = x) = p^x(1 - p)^{1-x}, \quad x = 0, 1$$

Consider  $T_n = \sum_{i=1}^n X_i$ . Note that  $T_n \sim \text{Bin}(n, p)$ . Then:

$$\begin{aligned}
 P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T_n = t) &= \begin{cases} \frac{P(X_1=x_1, X_2=x_2, \dots, X_n=x_n, T_n=t)}{P(T_n=t)} & \text{if } \sum_{i=1}^n x_i = t \\ 0 & \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{P(X_1=x_1, X_2=x_2, \dots, X_n=x_n)}{P(T_n=t)} & \text{if } \sum_{i=1}^n x_i = t \\ 0 & \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{\prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}}{\binom{n}{t} p^t (1-p)^{n-t}} & \text{if } \sum_{i=1}^n x_i = t \\ 0 & \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}}{\binom{n}{t} p^t (1-p)^{n-t}} & \text{if } \sum_{i=1}^n x_i = t \\ 0 & \text{otherwise} \end{cases} \\
 \text{thus } P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T_n = t) &= \begin{cases} \frac{1}{\binom{n}{t}} & \text{if } \sum_{i=1}^n x_i = t \\ 0 & \text{otherwise} \end{cases} \quad (1)
 \end{aligned}$$

and hence  $T_n$  is a sufficient statistic for  $p$ .

**Remark 9.1.** To generate a  $T_n$ -similar sample when  $T_n$  is given, say  $T_n = t$ , we define:

$$A_t = \left\{ (x_1, x_2, \dots, x_n) : \sum_{i=1}^n x_i = t \right\}$$

Note that  $\text{card}(A_t) = \binom{n}{t}$ . According to (1) we give equal weight, i.e. problem mass, to each element of  $A_t$ . We then choose one element of  $A_t$  randomly.

**Example 9.1.**  $X_i \stackrel{iid}{\sim} P(\lambda)$ ,  $i = 1, 2, \dots, n$

$$P_\lambda(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

Consider  $T_n = \sum_{i=1}^n x_i$ . Note that  $T_n \sim P_0(n\lambda)$  (Exercise)

$$\begin{aligned}
 P_\lambda(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T_n = t) &= \begin{cases} \frac{P_\lambda(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n, T_n = t)}{P(T_n = t)} & \text{if } \sum_{i=1}^n x_i = t \\ 0 & \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{\prod_{i=1}^n e^{-\lambda} \lambda^{x_i}}{\frac{e^{-n\lambda} (n\lambda)^t}{t!}} & \text{if } \sum_{i=1}^n x_i = t \\ 0 & \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{\prod_{i=1}^n \lambda^{x_i}}{\frac{(n\lambda)^t}{t!}} & \text{if } \sum_{i=1}^n x_i = t \\ 0 & \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{t!}{\prod_{i=1}^n x_i!} \cdot \frac{1}{n^t} & \text{if } \sum_{i=1}^n x_i = t \\ 0 & \text{otherwise} \end{cases} \quad (2)
 \end{aligned}$$

Thus:

$$\boxed{X |_{T_n=t} \sim \text{Multinomial}(t, p_i = \frac{1}{n}, i = 1, 2, \dots, n)}$$

Recall that:

$$(Y_1, Y_2, \dots, Y_k) \sim \text{Multinomial}(n, p_1, p_2, \dots, p_k) \text{ if } P(X_1 = x_1, \dots, X_k = x_k) = \binom{n}{x_1, \dots, x_k} \prod_{i=1}^k p_i^{x_i}$$

$$\text{where } \sum_{i=1}^n x_i = n \quad \& \quad \sum_{i=1}^k p_i = 1 \quad \text{and} \quad \binom{n}{x_1, x_2, \dots, x_k} = \frac{n!}{x_1! x_2! \dots x_k!}$$

Again, to generate a  $T_n$ -similar sample we define:

$$A_t = \left\{ (x_1, x_2, \dots, x_n) : \sum_{i=1}^n x_i = t \right\}$$

The probability mass associated to elements of  $A_t$  is given by (2). In other words, we choose an element of  $A_t$  using a multinomial distribution with  $n = t$ ,  $k = n$  and  $p_i = \frac{1}{n}$ ,  $i = 1, 2, \dots, n$ .

## 9.2 Likelihood

**Definition.** Let  $\tilde{Y} = (Y_1, Y_2, \dots, Y_n)$  be a random vector whose joint pdf or pmf depends on  $\tilde{\theta}$ , a vector of unknown parameters. **The Likelihood Function** is a function of  $\tilde{\theta}$ , for a realization  $\tilde{y} = (y_1, y_2, \dots, y_n)$  is defined to be:

$$\mathcal{L}(\tilde{\theta}; \tilde{y}) = \begin{cases} P_{\tilde{\theta}}(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n) & \text{if } y_i\text{'s are discrete random variables} \\ f_{\tilde{\theta}}(y_1, y_2, \dots, y_n) & \text{if } y_i\text{'s are continuous random variables} \end{cases}$$

**Example.**  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$ ,  $i = 1, 2, \dots, n$

$$P_p(X = x) = p^x(1-p)^{1-x}, \quad x = 0, 1$$

$$\mathcal{L}(p; x_1, x_2, \dots, x_n) = P_p(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

$$\implies = P_p(X_1 = x_1) \dots P_p(X_n = x_n) \quad \text{Independence and identically distributed}$$

$$= \prod_{i=1}^n p^{x_i}(1-p)^{1-x_i} = p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i} \quad (B)$$

**Example.**  $X_i \stackrel{iid}{\sim} P_0(\lambda)$ ,  $i = 1, 2, \dots, n$

$$P_\lambda(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

$$\mathcal{L}(\lambda; x_1, x_2, \dots, x_n) = P_\lambda(X_1 = x_1, \dots, X_n = x_n)$$

$$= \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \quad (P)$$

Independence and identically distributed

$$= \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}$$

**Example.**  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $i = 1, 2, \dots, n$

$$f_{\mu, \sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad -\infty < x < +\infty$$

$$\begin{aligned} \mathcal{L}(\mu, \sigma^2; x_1, x_2, \dots, x_n) &= f_{\mu, \sigma^2}(x_1, x_2, \dots, x_n) \\ &= \prod_{i=1}^n f_{\mu, \sigma^2}(x_i) \quad \text{Independent and identically distributed} \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right\} \quad (N) \end{aligned}$$

The examples we presented for sufficiency required first specifying a candidate statistic. The question then is how we come up with a sufficient statistic. The following theorem due to *Fisher & Neyman* tells us how to find sufficient statistic.

**Theorem** (Fisher-Neyman Factorization Theorem - Thm 9.4 , page 461). *A statistic  $T = T(Y_1, Y_2, \dots, Y_n)$  for  $\theta$  the parameter of the distribution of  $Y_1, Y_2, \dots, Y_n$  if and only if :  $\mathcal{L}(\theta; y_1, y_2, \dots, y_n) = g(t; \theta) h(y_1, y_2, \dots, y_n)$  For any realization  $(y_1, y_2, \dots, y_n)$ , where  $t = T(y_1, y_2, \dots, y_n)$ .*

**Example.**  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$  ,  $i = 1, 2, \dots, n$

Using (B) (the result above):

$$\mathcal{L}(p; x_1, x_2, \dots, x_n) = p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}$$

Then define  $T = \sum_{i=1}^n x_i$  ,  $g(t; p) = p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i} = p^t (1-p)^{n-t}$  and  $h(x_1, x_2, \dots, x_n) \equiv 1$ .

Using Fisher-Neyman theorem  $T = \sum_{i=1}^n X_i$  is a sufficient statistic for  $p$ .

**Example.**  $X_i \stackrel{iid}{\sim} P_0(\lambda)$  ,  $i = 1, 2, \dots, n$

Using (P) (the result above):  $\mathcal{L}(\lambda; x_1, \dots, x_n) = \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}$  . Define  $T = \sum_{i=1}^n X_i$

and  $g(t; \lambda) = e^{-n\lambda} \lambda^t = e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}$  and  $h(x_1, \dots, x_n) = \frac{1}{\prod_{i=1}^n x_i!}$  . Then using Fisher-



Neyman Theorem,  $T = \sum_{i=1}^n X_i$  is a sufficient statistic for  $\lambda$ .

**Example.**  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $i = 1, 2, \dots, n$ .

Note that now we have two unknown parameters,  $\mu$  and  $\sigma^2$ . Using (N) :

$$\mathcal{L}(\mu, \sigma^2; x_1, x_2, \dots, x_n) = \left(\frac{1}{\sqrt{2\pi\sigma}}\right)^n \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right\}$$

Now note that:  $\sum_{i=1}^n (x_i - \mu)^2 = \sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2$

Thus:  $\mathcal{L}(\mu, \sigma^2; x_1, x_2, \dots, x_n) = \left(\frac{1}{\sqrt{2\pi\sigma}}\right)^n \exp\left\{-\frac{\sum_{i=1}^n x_i^2}{2\sigma^2} + \frac{\mu \sum_{i=1}^n x_i}{\sigma^2} - \frac{n\mu^2}{2\sigma^2}\right\}$

Define:  $T = \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2\right)$ ,  $g(t; \theta) = \left(\frac{1}{\sqrt{2\pi\sigma}}\right)^n \exp\left\{-\frac{\sum_{i=1}^n x_i^2}{2\sigma^2} + \frac{\mu \sum_{i=1}^n x_i}{\sigma^2} - \frac{n\mu^2}{2\sigma^2}\right\}$  where

$\theta = (\mu, \sigma^2)$  and  $h(x_1, x_2, \dots, x_n) \equiv 1$ . Then using Fisher-Neyman Theorem

$T = \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2\right)$  is a sufficient statistic for  $\theta = (\mu, \sigma^2)$ .

**Remark.** Note that to identify the sufficient statistic using Fisher-Neyman Theorem, you only need the part of the likelihood in which you cannot separate the unknown parameters from observations. This part is called **kernel**. In other words, you can write a likelihood as the product of a function of observations alone, a function of parameters alone and the kernel. The sufficient statistic is in the kernel.

## 10 Lecture 10

### 10.1 The Rao-Blackwell Theorem

[Theorem 9.5 , page 464] An interesting and important application of sufficiency is in variance reduction. This application is formalized in a theorem due to Rao(C.R) and Blackwell (David) .

We first need to recall Theorem 5.14 (Page 286) and Theorem 5.15 (page 287):

Theorem 5.14 (page 286):

$$\mathbb{E}(X) = \mathbb{E} \left\{ \mathbb{E}(X|Y) \right\}$$

Theorem 5.15 (page 287):

$$\text{Var}(X) = \text{Var} \left\{ \mathbb{E}(X|Y) \right\} + \mathbb{E} \left\{ \text{Var}(X|Y) \right\}$$

**Theorem** (The Rao-Blackwell Theorem - Thm 9.5 , page 464). *Let  $\hat{\theta}$  be an unbiased estimator for  $\theta$  such that  $V(\hat{\theta}) < \infty$  . Suppose  $T$  is a sufficient statistic for  $\theta$  . Define  $\hat{\theta}^* = \mathbb{E}(\hat{\theta}|T)$ . Then, for all  $\theta$  :*

- (a)  $\mathbb{E}(\hat{\theta}^*) = \theta$
- (b)  $\text{Var}(\hat{\theta}^*) \leq \text{Var}(\hat{\theta})$  .

*Proof.* First note that  $T$  is a sufficient statistic for  $\theta$  , thus the distribution of  $\hat{\theta}$  given  $T$  does not depend on  $\theta$  . Therefore  $\mathbb{E}(\hat{\theta}|T)$  is a statistic . This is when sufficiency plays its role.

To prove part (a) of formula, we use Theorem 5.14, page 286:

$$\mathbb{E}(\hat{\theta}^*) = \mathbb{E} \left[ \mathbb{E}(\hat{\theta}|T) \right] = \mathbb{E}(\hat{\theta}) = \theta \quad \forall \theta.$$

To prove part (b), we use Theorem 5.15, page 287:

$$Var(\hat{\theta}^*) = Var\{\mathbb{E}(\hat{\theta}|T)\} \leq Var\{\mathbb{E}(\hat{\theta}|T)\} + \underbrace{\mathbb{E}\{Var(\hat{\theta}|T)\}}_{\geq 0}$$

□

**Remark (Completeness).** A statistic  $T$  or its family of distribution  $\{F_\theta: \theta \in \Theta\}$  where  $\Theta$  is the set of all admissible values of  $\theta$ , is called complete if for any reasonable  $g$ :

$$\mathbb{E}_\theta[g(T)] = 0 \quad , \quad \forall \theta \in \Theta$$

implies that  $g(t) = 0$  for all possible values of  $t$ . If a sufficient statistic  $T$  is also complete, then  $\hat{\theta}^* = \mathbb{E}(\hat{\theta}|T)$  will be the Minimum Variance Unbiased Estimator (MVUE). This often means that within the class of unbiased estimator  $\hat{\theta}^*$  is the least; i.e. the closest in the  $MSE$  sense, to the unknown parameter  $\theta$ . Recall that:

$$MSE_\theta(\hat{\theta}) = Var_\theta(\hat{\theta}) + Bias_\theta^2(\hat{\theta}) = Var_\theta(\hat{\theta})$$

if  $Bias_\theta(\hat{\theta}) = 0$ ; i.e. if  $\hat{\theta}$  is an unbiased estimator of  $\theta$ .

The notion of Completeness is due to **Lehmann(Eric Leo) and Scheffe'(Henry)**.

Then using the Rao-Blackwell and Lehmann-Schaffe' theorems we have an easy recipe for finding the MVUE.

**Step 1:** Using Fisher-Neyman theorem, find a sufficient statistic, say  $T$ , for  $\theta$ .

**Step 2:** Find an unbiased estimator  $\theta$ , say  $\hat{\theta}$ .

**Step 3:** Find  $\hat{\theta}^* = \mathbb{E}(\hat{\theta}|T)$ .

**Remark.** For the examples and exercises in the course, the sufficient statistic you find in **Step 1** using Fisher-Neyman Theorem is also complete.

**Example** (Ex. 9.6, page 466).  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$ ,  $i = 1, 2, \dots, n$

**Step 1:**

$$\mathcal{L}(p; x_1, \dots, x_n) = P_p(X_1 = x_1, \dots, X_n = x_n)$$

$$\text{Independence} \quad = \prod_{i=1}^n P_p(X_i = x_i)$$

$$\begin{aligned} \text{Identically distributed} \quad &= \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \\ &= p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i} \end{aligned}$$

Then  $T = \sum_{i=1}^n X_i$  is a sufficient statistic.

**Step 2:**

$\hat{p}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is an unbiased estimator of  $p$

$$\mathbb{E}(\hat{p}_n) = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{E}(X_i)}_p = \frac{1}{n} \cdot np = p$$

**Step 3:**

$$\text{Note that } \hat{p}_n = \frac{1}{n} \sum_{i=1}^n X_i = \frac{T}{n} \quad \text{thus} \quad \mathbb{E}(\hat{p}_n | T) = \frac{T}{n} = \hat{p}_n .$$

**Remark.** What we observed in *step 3* of the above example tells us that *step 3* of our recipe is redundant when  $\theta$  found in *step 2* is a function of the sufficient statistic found in *step 1*.

**Example** (Ex. 9.8 , page 467).  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$

**Step 1**

$$\begin{aligned}\mathcal{L}(\mu, \sigma^2; x_1, \dots, x_n) &= f_{\mu, \sigma^2}(x_1, \dots, x_n) \\ &= \prod_{i=1}^n f_{\mu, \sigma^2}(x_i) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}\end{aligned}$$

$$\text{hence: } \mathcal{L}(\mu, \sigma^2; x_1, \dots, x_n) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left\{-\frac{1}{2\sigma^2}\left[\sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2\right]\right\}$$

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

**Step 2:**  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is an unbiased estimator of  $\mu$  and  $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$  is an unbiased estimator of  $\sigma^2$ .

**Step 3:**

$$\mathbb{E}(\bar{X}_n | T) = \bar{X}_n$$

since  $\bar{X}_n$  is a

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

is also a function of  $T$

$$\text{Thus } \mathbb{E}(S_n^2 | T) = S_n^2$$

Thus  $\bar{X}_n$  is the *MVUE* of  $\mu$  and  $S_n^2$  is the *MVUE* of  $\sigma^2$ .

**Example** (Ex. 9.7, page 466-467).  $Y_i \stackrel{iid}{\sim} \text{Weibull}(m=2, \theta)$ ,  $i = 1, 2, \dots, n$

$$f_{\theta}(y) = \begin{cases} \left(\frac{2y}{\theta}\right) e^{-\frac{y^2}{\theta}} & y > 0 \\ 0 & \text{otherwise} \end{cases}$$

**Step 1:**

$$\begin{aligned}\mathcal{L}(\theta; y_1, y_2, \dots, y_n) &= \prod_{i=1}^n \left(\frac{2y_i}{\theta}\right) e^{-\frac{y_i^2}{\theta}} \\ &= \left(\frac{2}{\theta}\right)^n e^{-\frac{1}{\theta} \sum_{i=1}^n y_i^2} \prod_{i=1}^n y_i\end{aligned}$$

Thus  $T = \sum_{i=1}^n Y_i^2$  is a sufficient statistic for  $\theta$ .

Note that the Kernel is  $\exp\left\{-\frac{1}{\theta} \sum_{i=1}^n Y_i^2\right\}$ . We can also see this through Fisher-Neyman by choosing:

$$g(t; \theta) = \left(\frac{2}{\theta}\right)^n e^{-\frac{t}{\theta}} \quad \text{and} \quad h(y_1, \dots, y_n) = \prod_{i=1}^n y_i \quad \text{where} \quad t = \sum_{i=1}^n y_i^2$$

**Step 2:** Define  $W_i = Y_i^2$ ,  $i = 1, 2, \dots, n$ . Note that:

$$\begin{aligned} f_w(w) &= f_Y(\sqrt{w}) \left| \frac{d\sqrt{w}}{dw} \right| \quad \text{using Transformation Method} \\ &= \begin{cases} \left(\frac{2\sqrt{w}}{\theta}\right) e^{-\frac{(\sqrt{w})^2}{\theta}} \cdot \frac{1}{2\sqrt{w}} & \text{if } w > 0 \\ 0 & \text{otherwise} \end{cases} \\ &= \begin{cases} \frac{1}{\theta} e^{-\frac{w}{\theta}} & \text{if } w > 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

Thus  $W \sim \text{Exp}(\theta)$  and therefore

$$\mathbb{E}(T) = \mathbb{E}\left(\sum_{i=1}^n Y_i^2\right) = \sum_{i=1}^n \mathbb{E}(Y_i^2) = \sum_{i=1}^n \mathbb{E}(W_i) = n\theta \implies \mathbb{E}\left(\frac{T}{n}\right) = \theta$$

**Step 3**

$$\mathbb{E}\left(\frac{T}{n} \mid T\right) = \frac{T}{n} \quad \text{therefore} \quad \frac{T}{n} = \frac{1}{n} \sum_{i=1}^n Y_i^2 \quad \text{is the MVUE of } \theta.$$

**Remark.** Sufficient statistics can often be used to make a pivotal quantity, in the above example for instance,

$$\frac{2}{\theta} W \sim \chi_{(2)}^2 \quad (\text{Exercise})$$

and hence

$$\frac{2}{\theta} \sum_{i=1}^n W_i = \frac{2}{\theta} \sum_{i=1}^n Y_i^2 \sim \chi_{(2n)}^2 \quad \text{i.e.} \quad \frac{2T}{\theta} \sim \chi_{2n}^2$$

This pivotal quantity can therefore be used to make *exact confidence interval* for  $\theta$ . See example 9.10, page 468, confidence interval made using sufficient statistic based on pivotal quantities often have the shortest possible length for a given confidence level.

## 11 Lecture 11

**Methods of Estimation:**      A) Method of Maximum Likelihood (ML)  
   B) Method of Moments

### 11.1 Method of Maximum Likelihood (ML)

**Definition.** *The Maximum Likelihood Estimation (MLE) of a parameter  $\theta$  based on the realized values  $(y_1, y_2, \dots, y_n)$  of a sample  $Y_1, Y_2, \dots, Y_n$  is:*

$$\hat{\theta}_{ML} = \operatorname{argmax} \mathcal{L}(\theta; y_1, \dots, y_n)$$

Then we have a two step procedure for finding the  $\hat{\theta}_{ML}$ :

**Step 1:** Set up the likelihood  $\mathcal{L}(\theta; y_1, \dots, y_n)$

**Step 2:** Find the maximizer of the likelihood when we have a random sample, which is the case in this course:

$$\mathcal{L}(\theta; y_1, \dots, y_n) = \prod_{i=1}^n f_{\theta}(x_i)$$

To find  $\hat{\theta}_{ML}$  is often easier to work with  $\mathcal{L}(\theta; y_1, \dots, y_n) = \log(\mathcal{L}(\theta; y_1, \dots, y_n)) = \sum_{i=1}^n \log(f_{\theta}(x_i))$ .

Note that log is a monotone increasing function, Thus:

$$\operatorname{argmax} l(\theta; y_1, \dots, y_n) = \operatorname{argmax} \mathcal{L}(\theta; y_1, \dots, y_n)$$

**Example.**  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$  ,  $i = 1, 2, \dots, n$

$$P(X = x) = p^x(1-p)^{1-x} , \quad x = 0, 1$$

**Step 1:**  $\mathcal{L}(p; x_1, \dots, x_n) = p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}$



**Step 2:**

$$\begin{aligned}
 l(p; x_1, \dots, x_n) &= \log ( \mathcal{L}(p; x_1, \dots, x_n) ) \\
 &= \left( \sum_{i=1}^n x_i \right) \log p + \left( n - \sum_{i=1}^n x_i \right) \log(1-p) \\
 \frac{\partial}{\partial p} l(p; x_1, \dots, x_n) &= \frac{\sum_{i=1}^n x_i}{p} - \frac{n - \sum_{i=1}^n x_i}{1-p}
 \end{aligned}$$

Let  $t = \sum_{i=1}^n x_i$ . The  $\hat{p}_{ML}$  is then the solution to

$$\frac{t}{\hat{p}_{ML}} - \frac{n-t}{1-\hat{p}_{ML}} = 0 \iff \frac{t}{n-t} = \frac{\hat{p}_{ML}}{1-\hat{p}_{ML}} \iff \hat{p}_{ML} = \frac{t}{n} = \frac{1}{n} \sum_{i=1}^n X_i = \hat{p}_n$$

Note that

$$\frac{\partial^2 l}{\partial p^2} = -\frac{\sum_{i=1}^n x_i}{p^2} - \frac{n - \sum_{i=1}^n x_i}{(1-p)^2} < 0$$

Thus  $\hat{p}_{ML}$  is the maximizer of  $l(p; x_1, \dots, x_n)$ .

**Example.**  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $i = 1, 2, \dots, n$

**Step 1:**

$$\mathcal{L}(\mu, \sigma^2; x_1, \dots, x_n) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\}$$

$$\text{hence } l(\mu, \sigma^2, x_1, \dots, x_n) = -\frac{n}{2} \log \sqrt{2\pi} - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

**Step 2:**

$$\frac{\partial l(\mu, \sigma^2)}{\partial \mu} = -\frac{1}{2\sigma^2} \sum_{i=1}^n -2(x_i - \mu)$$

$$\frac{\partial l(\mu, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2$$

the MLE ,  $\hat{\mu}_{ML}$  &  $\hat{\sigma}_{ML}^2$  are therefore solutions to:

$$\begin{cases} -\frac{1}{2\hat{\sigma}_{ML}^2} \sum_{i=1}^n -2(x_i - \hat{\mu}_{ML}) = 0 \\ -\frac{n}{2\hat{\sigma}_{ML}^2} + \frac{1}{2\hat{\sigma}_{ML}^4} \sum_{i=1}^n (x_i - \hat{\mu}_{ML})^2 = 0 \end{cases}$$

From the 1<sup>st</sup> equation we find :

$$\sum_{i=1}^n (x_i - \hat{\mu}_{ML}) = 0 \implies \sum x_i = n\hat{\mu}_{ML} \implies \boxed{\hat{\mu}_{ML} = \bar{x}_n}$$

If we plug in for  $\hat{\mu}_{ML}$  in the second equation, we find

$$-\frac{n}{2\hat{\sigma}_{ML}^2} + \frac{1}{2\hat{\sigma}_{ML}^4} \sum_{i=1}^n (x_i - \bar{x}_n)^2 = 0$$

$$\implies \frac{1}{\hat{\sigma}_{ML}^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2 = n$$

$$\implies \boxed{\hat{\sigma}_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2}$$

To show that this is a maximizer we should check that:

$$\begin{bmatrix} \frac{\partial^2 l}{\partial \mu^2} & \frac{\partial^2 l}{\partial \mu \partial \sigma^2} \\ \frac{\partial^2 l}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 l}{\partial (\sigma^2)^2} \end{bmatrix} \text{ is a negative-definite matrix .}$$

This is not a hard task, but it is not required in this course. We only check the 2<sup>nd</sup> derive for cases that there is only one unknown parameter.

**Example.**  $X_i \stackrel{iid}{\sim} Unif(0, \theta)$  ,  $i = 1, 2, \dots, n$

$$f_{\theta}(x) = \begin{cases} \frac{1}{\theta} & 0 < x < \theta \\ 0 & \text{otherwise} \end{cases}$$

In other words,  $f_\theta(x) = \frac{1}{\theta} \cdot I_{[0,\theta]}(x)$  where  $I_A(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{if } x \notin A \end{cases}$  **Step 1:**

$$\mathcal{L}(\theta; x_1, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i) = \frac{1}{\theta^n} \prod_{i=1}^n I_{[0,\theta]}(x_i)$$

Now note that

$$\prod_{i=1}^n I_{[0,\theta]}(x_i) = I_{[0,\theta]}(\max_{1 \leq i \leq n} x_i)$$

Since  $0 \leq x_i \leq \theta$ ,  $i = 1, 2, \dots, n$  if and only if  $0 \leq \max_{1 \leq i \leq n} x_i \leq \theta$ . Therefore:

$$\mathcal{L}(\theta; x_1, \dots, x_n) = \frac{1}{\theta^n} I_{[0,\theta]}(\max_{1 \leq i \leq n} x_i)$$

**Step 2:** We note that the likelihood is a monotone decreasing function of  $\theta$ .

That is why, the max of  $\mathcal{L}(\theta; x_1, \dots, x_n)$  happens when  $\theta$  takes its smallest possible value. Since that  $0 \leq \max_{1 \leq i \leq n} x_i \leq \theta$ , the smallest value for  $\theta$  is  $\max_{1 \leq i \leq n} x_i$ , thus  $\hat{\theta}_{ML} = \max_{1 \leq i \leq n} x_i$ .

The method of *ML* has both the intuitive and theoretical appeal.

**Intuitive Appeal:** The *ML* method is based on the idea that "What I have observed is what should have expected to observe". In other words, we observe the most likely scenario. Now given a sample, we choose the unknown parameter such that what we have observed has its maximum possible chance.

**Theoretical Appeal:**

\* **Consistency:** *MLE* is a consistent estimator under rather mild conditions.

\* **Asymptotic Normality:**  $\sqrt{n}(\hat{\theta}_{ML} - \theta) \overset{app}{\sim} N(0, I^{-1}(\theta))$  for large  $n$  where  $I(\theta) = \mathbb{E}\left[\left[\frac{\partial}{\partial \theta} \log f_\theta(X)\right]^2\right]$ , the Fisher information.

We can therefore make confidence interval for  $\theta$  easily if we use  $\hat{\theta}_n = \hat{\theta}_{ML}$  as the estimator.

\* **Asymptotic Efficiency:** *MLE* is the most concentrated estimator about its estimand among a considerably large class of reasonable estimators called

"regular estimators".

\* **Invariance:** If  $\hat{\theta}_{ML}$  is the MLE of  $\theta$ , then  $g(\hat{\theta}_{ML})$  is the MLE of  $g(\theta)$ .

This property simplifies life a bit. Since if we find the MLE of  $\theta$ , then we have found MLE of any function of  $\theta$ .

**Example.**  $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$ ,  $i = 1, 2, \dots, n$

We showed that  $\hat{p}_{ML} = \hat{p}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is a MLE of  $p$ . Now suppose that we want to find the MLE of  $p(1-p)$ , which is the variance of  $\text{Bernoulli}(p)$ . Using the Invariance, the MLE of  $p(1-p)$  is  $\hat{p}_{ML}(1 - \hat{p}_{ML}) = (\frac{\sum_{i=1}^n X_i}{n})(1 - \frac{\sum_{i=1}^n X_i}{n})$

## 11.2 Method of Moments

Let  $\mu_K = \mathbb{E}(X^K)$  and  $m_K = \frac{1}{n} \sum_{i=1}^n X_i^K$ . Now suppose  $X_i \stackrel{iid}{\sim} f_{\tilde{\theta}}$ ,  $i = 1, 2, \dots, n$  where  $\tilde{\theta} = (\theta_1, \dots, \theta_r)$ . Clearly  $\mu_K$  is going to be a function of  $\tilde{\theta}$ . Suppose  $\mu_K$  exists. Then the method of moments estimators are the solution to the following equation:

$$\mu_K(\tilde{\theta}) = m_K, \quad K = 1, 2, \dots, r.$$

**Example.**  $X_i \stackrel{iid}{\sim} \Gamma(\alpha, \beta)$ ,  $i = 1, 2, \dots, n$

$$\mathbb{E}(X) = \alpha\beta, \quad \text{Var}(X) = \alpha\beta^2$$

Thus  $\mathbb{E}(X^2) = \text{Var}(X) + [\mathbb{E}(X)]^2 = \alpha\beta^2 + \alpha^2\beta^2 = \alpha\beta(\beta + \alpha\beta)$

$$\text{Now:} \quad \begin{cases} \alpha\beta = \bar{X}_n \\ \alpha\beta(\beta + \alpha\beta) = \frac{1}{n} \sum_{i=1}^n X_i^2 \end{cases}$$

Plugging-in from the 1<sup>st</sup> equation into the 2<sup>nd</sup> equation we find:

$$\bar{X}_n(\beta + \bar{X}_n) = \frac{1}{n} \sum_{i=1}^n X_i^2$$

hence

$$\beta = \frac{\frac{1}{n} \sum_{i=1}^n X_i^2}{\bar{X}_n} - \bar{X}_n$$

Using the first equation:

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

## 12 Lecture 12

### 12.1 Section 9.8 - Large Sample Property of the MLEs

Let  $X \sim f_\theta(x)$ . Suppose we have two observations from  $f_\theta$ ,  $x_1 = 2$  and  $x_2 = 5$ . We want to see how we can quantify the amount of information in each of these two observations about  $\theta$ . Note that our only link between the observations, i.e.  $X$ , and unobservable, i.e.  $\theta$ , is  $f_\theta(x)$ . So this is the channel through which information are transmitted. Now suppose the following figures depict the graph of  $f_\theta(2)$  and  $f_\theta(5)$ :

#MISSING GRAPHS - Lecture 12 P1

Since that  $f_\theta(2)$  is constant, it is not sensitive w.r.t changes in  $\theta$ . In contrast  $f_\theta(5)$  seems to respond to the changes in  $\theta$  values. As such 5 is much more information about  $\theta$  than 2. In fact, 2 does not have any information about  $\theta$ . Such sensitivity can be measured by derivative, i.e.  $\frac{\partial}{\partial \theta} f_\theta$ . Now Fisher would believe that information should increase linearly with the sample size  $n$ . As such, he would work with  $\log f_\theta$ , rather than  $f_\theta$ . This then leads us to  $\frac{\partial}{\partial \theta} \log f_\theta$ . Now to define the information in the random variable  $X$  about  $\theta$ , we can look at the (Euclidean) length of  $\frac{\partial}{\partial \theta} \log f_\theta(x)$  which is  $\mathbb{E}\left\{\left[\frac{\partial}{\partial \theta} \log f_\theta(x)\right]^2\right\}$

$$I(\theta) = \mathbb{E}\left\{\left[\frac{\partial}{\partial \theta} \log f_\theta(x)\right]^2\right\} = \int_{\theta} \left[\frac{\partial}{\partial \theta} \log f_\theta\right]^2 f_\theta(x) dx$$

is called the *Fisher Information Amount* . Note that :

$$\begin{aligned}
\mathbb{E}\left[\frac{\partial}{\partial\theta} \log f_{\theta}(x)\right] &= \int_x \frac{\partial}{\partial\theta} \log f_{\theta}(x) \cdot f_{\theta}(x) dx \\
&= \int_x \frac{\dot{f}_{\theta}(x)}{f_{\theta}(x)} \cdot \cancel{f_{\theta}(x)} dx \quad \text{s.t} \quad \dot{f}_{\theta}(x) = \frac{\partial}{\partial\theta} f_{\theta}(x) \\
&= \frac{\partial}{\partial\theta} \int_x \overbrace{f_{\theta}(x) dx}^1 = 0 \\
\text{Provided that } \frac{\partial}{\partial\theta} \int_x f_{\theta}(x) dx &= \int_x \frac{\partial}{\partial\theta} f_{\theta}(x) dx \quad (1) \\
\text{Thus: } I(\theta) &= \text{Var}\left\{\frac{\partial}{\partial\theta} \log f_{\theta}(x)\right\}
\end{aligned}$$

Taking the 2<sup>nd</sup> derivative we have:

$$\begin{aligned}
\frac{\partial^2}{\partial\theta^2} \log f_{\theta}(x) &= \frac{\partial}{\partial\theta} \left[ \frac{\dot{f}_{\theta}(x)}{f_{\theta}(x)} \right] = \frac{\dot{f}_{\theta}(x)f_{\theta}(x) - [\dot{f}_{\theta}(x)]^2}{[f_{\theta}(x)]^2} \\
\text{where } \ddot{f}_{\theta}(x) &= \frac{\partial^2}{\partial\theta^2} f_{\theta}(x) . \text{ Thus} \\
\frac{\partial^2}{\partial\theta^2} \log f_{\theta}(x) &= \frac{\ddot{f}_{\theta}(x)}{f_{\theta}(x)} - \left[ \frac{\dot{f}_{\theta}(x)}{f_{\theta}(x)} \right]^2 \\
&= \frac{\ddot{f}_{\theta}(x)}{f_{\theta}(x)} - \left[ \frac{\partial}{\partial\theta} \log f_{\theta}(x) \right]^2
\end{aligned}$$

Now taking  $\mathbb{E}$  from both sides, we have:

$$\mathbb{E}\left[\frac{\partial^2}{\partial\theta^2} \log f_{\theta}(X)\right] = \mathbb{E}\left[\frac{\ddot{f}_{\theta}(X)}{f_{\theta}(X)} - I(\theta)\right] \quad (\dagger)$$

Note that:

$$\begin{aligned}
\mathbb{E}\left[\frac{\ddot{f}_{\theta}(X)}{f_{\theta}(X)}\right] &= \int_x \frac{\ddot{f}_{\theta}(x)}{f_{\theta}(x)} \cdot \cancel{f_{\theta}(x)} dx = \int_x \ddot{f}_{\theta}(x) dx \\
&= \int_x \frac{\partial^2}{\partial\theta^2} f_{\theta}(x) dx \\
\text{Now if: } \frac{\partial^2}{\partial\theta^2} \int_x f_{\theta}(x) dx &= \int_x \frac{\partial^2}{\partial\theta^2} f_{\theta}(x) dx \quad (2)
\end{aligned}$$

$$\text{We obtain: } \mathbb{E}\left[\frac{\ddot{f}_{\theta}(x)}{f_{\theta}(x)}\right] = \int_x \frac{\partial^2}{\partial\theta^2} f_{\theta}(x) dx = \frac{\partial^2}{\partial\theta^2} \int_x \overbrace{f_{\theta}(x) dx}^1 = 0$$

Using (†) we therefore have:

$$\mathbb{E}\left[\frac{\partial^2}{\partial\theta^2} \log f_{\theta}(X)\right] = -I(\theta)$$

and hence:

$$I(\theta) = \mathbb{E}\left[\left\{\frac{\partial}{\partial\theta} \log f_{\theta}(X)\right\}^2\right] = -\mathbb{E}\left[\frac{\partial^2}{\partial\theta^2} \log f_{\theta}(X)\right]$$

Provided that (2) holds. Note that a necessary condition for (2) is that  $\text{supp } f_\theta(x) = \{x : f_\theta(x) > 0\}$  does not depend on  $\theta$ . This is the case for most distributions you have seen, i.e. Bernoulli, Binomial, Poisson, Exponential, Gamma,  $\chi^2$ , Normal. This is **is not**, however the case for **uniform distribution**.

## 12.2 Asymptotic Distribution of MLE and Approximate Pivotal Quantity

We discussed properties of the MLEs. We learned the invariance property of MLE which says that if  $\hat{\theta}_{ML}$  is the MLE of  $\theta$ , then  $\tau(\hat{\theta}_{ML})$  is the MLE of  $\tau(\theta)$ .

Under some mild conditions:

$$\sqrt{n}(\tau(\hat{\theta}_{ML}) - \tau(\theta)) \stackrel{app}{\approx} N\left(0, \left[\frac{\partial}{\partial \theta} \tau(\theta)\right]^2 I^{-1}(\theta)\right) \quad \text{for large } n \quad (*)$$

In particular when  $\tau(\theta) = \theta$ , we have:

$$\sqrt{n}(\tau(\hat{\theta}_{ML}) - \theta) \stackrel{app}{\approx} N(0, I^{-1}(\theta)) \quad (**)$$

Note that  $\left[\frac{\partial}{\partial \theta} \tau(\theta)\right]^2 I^{-1}(\theta)$  is the amount of information in variable  $X$  about  $\eta = \tau(\theta)$ . In fact:

$$\begin{aligned} I(\eta) &= \mathbb{E}\left[\left\{\frac{\partial}{\partial \theta} \log f_\theta(X)\right\}^2\right] = \mathbb{E}\left[\left\{\frac{\partial}{\partial \theta} \log f_\theta(X) \cdot \frac{\partial \theta}{\partial \eta}\right\}^2\right] \\ &= \left[\frac{\partial \theta}{\partial \eta}\right]^2 \cdot \mathbb{E}\left[\left\{\frac{\partial}{\partial \theta} \log f_\theta(X)\right\}^2\right] \\ &= \left[\frac{\partial \theta}{\partial \eta}\right]^2 \cdot I(\theta) \end{aligned}$$

Since  $\frac{\partial \theta}{\partial \eta}$  is not a function of  $X$ . Thus:

$$I^{-1}(\eta) = \left[\frac{\partial \eta}{\partial \theta}\right]^2 I^{-1}(\theta) = \left[\frac{\partial \tau(\theta)}{\partial \theta}\right]^2 I^{-1}(\theta)$$



using the above result:

$$\frac{\sqrt{n}(\tau(\hat{\theta}_{ML}) - \tau(\theta))}{\sqrt{\left[\frac{\partial}{\partial \theta} \tau(\theta)\right]^2 I^{-1}(\theta)}} \stackrel{app}{\sim} N(0, 1), (\mathcal{L}) \text{ for large } n$$

**For (\*) and (\*\*)** : When we develop theory for the maximum likelihood estimators, we show that under same mild condition:

$$\sqrt{n}(\hat{\theta}_{ML} - \theta) \stackrel{app}{\sim} N(0, I^{-1}(\theta))$$

Now for  $\eta = \tau(\theta)$  we have the same result, i.e.:

$$\sqrt{n}(\hat{\eta}_{ML} - \eta) \stackrel{app}{\sim} N(0, I^{-1}(\eta)) .$$

Then we use the invariance property to conclude that  $\hat{\eta}_{ML} = \tau(\hat{\theta}_{ML})$ , if further,  $\eta$  is continuously differentiable, then:

$$I^{-1}(\eta) = \left[\frac{\partial}{\partial \theta} \tau(\theta)\right]^2 I^{-1}(\theta)$$

Thus  $\tau(\hat{\theta}_{ML}) \pm \zeta_{\frac{\alpha}{2}} \sqrt{\frac{[\dot{\tau}(\theta)]^2 I^{-1}(\theta)}{n}}$  is a  $100(1 - \alpha)\%$  confidence interval for  $\tau(\theta)$ , where  $\dot{\tau}(\theta) = \frac{\partial}{\partial \theta} \tau(\theta)$ . For practical purposes the margin of error should be estimated.

As long as the estimator is consistent, the asymptotic confidence interval is still valid. Now if  $\tau$  is continuously differentiable,  $\dot{\tau}(\hat{\theta}_{ML})$  is a consistent estimator of  $\dot{\tau}(\theta)$  by continuous mapping theorem. Under mild conditions  $I(\theta)$  is a continuous function of  $\theta$ , and hence  $I(\hat{\theta}_{ML})$  is a consistent estimator of  $I(\theta)$  by the continuous mapping theorem. Thus:

$$\tau(\hat{\theta}_{ML}) \pm \zeta_{\frac{\alpha}{2}} \sqrt{\frac{\dot{\tau}(\hat{\theta}_{ML}) I^{-1}(\hat{\theta}_{ML})}{n}} \quad (\$)$$

is an approximate  $100(1 - \alpha)\%$  confidence interval for  $\tau(\theta)$ , where  $\zeta_{\frac{\alpha}{2}}$  is chosen such that  $P(Z > \zeta_{\frac{\alpha}{2}}) = \frac{\alpha}{2}$ ,  $Z \sim N(0, 1)$ .

**Example 12.1** (# 9.14 - page 484).

$Y_i \stackrel{iid}{\sim} \text{Ber}(p)$ ,  $i = 1, 2, \dots, n$ ,  $\tau(p) = p(1-p)$ ,  $\dot{\tau}(p) = \frac{\partial}{\partial p} \tau(p) = 1-2p$  and  $P_p(Y = y) = p^y(1-p)^{1-y}$ ,  $y = 0, 1$

$$\begin{aligned}\log P_p(Y = y) &= y \log p + (1 - y) \log(1 - p) \\ \frac{\partial}{\partial p} \log P_p(Y = y) &= \frac{y}{p} - \frac{1 - y}{1 - p} \\ \frac{\partial^2}{\partial p^2} \log P_p(Y = y) &= -\frac{y}{p^2} - \frac{1 - y}{(1 - p)^2}\end{aligned}$$

Now we have:

$$\begin{aligned}\mathbb{E}\left\{-\frac{\partial^2}{\partial p^2} \log P_p(Y = y)\right\} &= I(p) \\ I(p) &= \mathbb{E}\left\{\frac{Y}{p^2} + \frac{1 - Y}{(1 - p)^2}\right\} = \frac{\mathbb{E}(Y)}{p^2} + \frac{\mathbb{E}(1 - Y)}{(1 - p)^2} \\ &= \frac{p}{p^2} + \frac{1 - p}{(1 - p)^2} = \frac{1}{p} + \frac{1}{1 - p} = \frac{1}{p(1 - p)}\end{aligned}$$

Thus:  $I^{-1}(p) = p(1 - p)$ ,  $\hat{p}_{ML} = \frac{1}{n} \sum_{i=1}^n Y_i$

Using (\$), we have:

$$\hat{p}_{ML}(1 - \hat{p}_{ML}) \pm \zeta_{\frac{\alpha}{2}} \sqrt{\frac{(1 - 2\hat{p}_{ML})^2 \hat{p}_{ML}(1 - \hat{p}_{ML})}{n}}$$

is a  $100(1 - \alpha)\%$  confidence interval for  $p(1-p)$ .

## 12.3 Chapter 10 - Testing Statistical Hypothesis