

# Math 375: Statistics

Winter 2020 Lecture Notes

Shereen Elaidi

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Properties of Random Samples</b>	<b>2</b>
2.1	Sampling from the Normal Population . . . . .	6
2.2	Standardisation . . . . .	9
2.3	F-Test Basis . . . . .	10
2.4	Limiting Sample Distributions (Asymptotics) . . . . .	11
2.4.1	Question 1 . . . . .	12
2.4.2	Question 2 . . . . .	13
2.5	Order Statistics . . . . .	15
<b>3</b>	<b>Point Estimation</b>	<b>16</b>
3.1	Ways We Can Construct Estimators . . . . .	16
3.1.1	Method of moments . . . . .	16
3.1.2	Method of Maximum Likelihood . . . . .	19
3.1.3	Invariance of the MLE Theorem . . . . .	22
3.1.4	Bayesian Method . . . . .	24
3.2	Method of Evaluating Estimators . . . . .	25
3.3	Best Unbiased Estimators . . . . .	28
3.3.1	Geometric Interpretation . . . . .	30
<b>4</b>	<b>Sufficiency and Completeness</b>	<b>37</b>
4.1	Proving the Sufficiency of Statistics . . . . .	38
4.2	Sufficient Statistics . . . . .	40
4.3	Applications of Sufficiency in Point Estimation . . . . .	42

## 1 Introduction

Data consists of observations  $x_1, \dots, x_n$ . These are regarded as realisations of random phenomena modelled by the random variables  $X_1, \dots, X_n$ . In this course, the  $X_i$ 's will be random variables in  $\mathbb{R}^d$ , usually with  $d = 1$ .

**Definition 1** (Random Sample). The random variables  $X_1, \dots, X_n$  are called a random sample from a distribution  $F$  for  $i = 1, \dots, n$  (or, if we want, we write  $X \sim F$ ).

In this course, the data will be a realisation of the random sample  $F$ . The basic issue is that  $F$  is unknown, so our task is to learn  $F$  from the realisations  $x_1, \dots, x_n$ . A model for  $F$  is  $\mathcal{F}$ , which is a collection of (certain) probability distributions such that  $F \in \mathcal{F}$ . It is always an (artificial) approximation to reality.

**Example 1** (U.S. 2016 Election Poll). .  $n = 2,000$ , and let  $x_1, \dots, x_n$  be the realisations of  $X_1, \dots, X_n$ , which are iid. Then, assuming there are only two candidates:

$$F \in \mathcal{F} = \{ \text{Bernoulli}(p) \mid p \in ]0, 1[ \}$$

This is an example of a **parametric model**, since  $p$ , which is the probability of success, is an *unknown* parameter. We know that for each  $x_i$ , we have  $x_i \in \{0, 1\}$ . To estimate  $p$ , we have:

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

i.e., the sample mean. By the Weak Law of Large Numbers, which we can use thanks to the iid assumption, we can see how good the estimator is by observing that it will  $\hat{p}$  will converge in probability to  $p$  as  $n \rightarrow \infty$ :

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i \xrightarrow{p} p \quad (1)$$

## 2 Properties of Random Samples

We first remark that the assumption that the  $X_1, \dots, X_N$  are iid is also called **sampling from an infinite population**. To see why, consider that our population were finite, say  $\{x_1, \dots, x_N\}$ , and we sample the  $X_1, \dots, X_N$  with replacement. Then, the probability of choosing a specific  $x_k$  would be:

$$\begin{aligned} \mathbb{P}[X_1 = x_k] &= \frac{1}{N} \quad \forall k \in \{1, \dots, N\} \\ \mathbb{P}[X_2 = x_k \mid X_1 = x_j] &= \begin{cases} 0 & \text{if } k = j \\ \frac{1}{N-1} & \text{if } k \neq j \end{cases} \end{aligned}$$

Using the law of total probability, we can obtain  $\mathbb{P}[X_2 = x_k]$ :

$$\begin{aligned} \mathbb{P}[X_2 = x_k] &= \sum_{j=1}^N \mathbb{P}[x_2 = x_k \mid X_1 = x_j] \mathbb{P}[X_1 = x_j] \\ &= \sum_{j=1, j \neq k}^N \frac{1}{N-1} \frac{1}{N} = \frac{1}{N} \end{aligned}$$

These samples are identically distributed but not independent. When  $N \gg n$ , the dependence between  $X_1, \dots, X_n$  plays *essentially no role*. The following example illustrates this:

**Example 2.** Let  $N = 2,000$  and  $n = 10$ . If we assume that they are not independent:

$$\mathbb{P}[X_1 > 2,000, \dots, X_n > 2,000] = \frac{\binom{800}{10} \binom{200}{0}}{\binom{1000}{0}} \approx 0.106164$$

If do assume they are independent:

$$\mathbb{P}[X_1 > 2,000, \dots, X_n > 2,000] = [\mathbb{P}[X_1 > 200]]^{10} = \left( \frac{800}{1000} \right)^{10} \approx 0.107374$$

**Definition 2** (Statistic). Let  $X_1, \dots, X_n$  be a random sample from  $F$ , a distribution on  $\mathbb{R}^d$ . For a measurable function:

$$T : \underbrace{\mathbb{R}^d \times \dots \times \mathbb{R}^d}_{n \text{ times}} \rightarrow \mathbb{R}^k \quad (2)$$

the random variable  $T(X_1, \dots, X_n)$  is called a **statistic**. The distribution of the statistic,  $T(X_1, \dots, X_n)$  is called a **sampling distribution** of the statistic  $T(X_1, \dots, X_n)$ .

**CAUTION:**  $T(X_1, \dots, X_n)$  is a function of  $X_1, \dots, X_n$  ONLY. That is,  $T(X_1, \dots, X_n)$  must be a *vector of numbers*.

**Example 3.** Let's go back to the opinion poll example. The estimator

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

is a statistic. A realisation is  $0.47 = T(X_1, \dots, X_n)$ . Note that the following is *not* a statistic:

$$\left( \frac{1}{n} \sum_{i=1}^n x_i - p \right)^2 \quad (3)$$

since there is a dependence on a parameter  $p$ .

**Definition 3** (Sample Mean, Sample Variance, Sample Standard Deviation). The average:

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i \quad (4)$$

is called the sample mean. The quantity:

$$s^2 := \frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (5)$$

is called the sample variation. The statistic  $s := \sqrt{s^2}$  is called the sample standard deviation. When these values are realised, they are denoted  $\bar{x}$ ,  $s^2$ , and  $s$ .  $\bar{x}$  and  $s^2$  are measures of central tendency and variability, respectively.

**Theorem 1.** Suppose that  $x_1, \dots, x_n \in \mathbb{R}$  and let  $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$ . Then:

1.

$$\min_{\alpha} \sum_{i=1}^n (x_i - \alpha)^2 = \sum_{i=1}^n (x_i - \bar{x})^2 \quad (6)$$

2.

$$(n-1)s^2 = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 \quad (7)$$

*Proof.* 1.

$$\begin{aligned} \sum_{i=1}^n (x_i - \alpha)^2 &= \sum_{i=1}^n (x_i \pm \bar{x} - \alpha)^2 \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + (\bar{x} - \alpha)n + 2(\bar{x} - \alpha) \underbrace{\sum_{i=1}^n (x_i - \bar{x})}_{:= n\bar{x} - n\bar{x} = 0} \\ &\geq \sum_{i=1}^n (x_i - \bar{x})^2 \end{aligned}$$

2. Set  $\alpha = 0$  in the preceding calculation. Then:

$$\sum_{i=1}^n x_i^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + n\bar{x}^2$$

□

**Lemma 2.** Let  $X_1, \dots, X_n$  be iid from  $F$  on  $\mathbb{R}$ , and let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be measurable, and  $X \sim F$ . Suppose that  $\text{Var}[g(X)] < \infty$ . Then:

1.

$$\mathbb{E} \left[ \sum_{i=1}^n g(X_i) \right] = n\mathbb{E}[g(X)] \quad (8)$$

2.

$$\text{Var} \left[ \sum_{i=1}^n g(X_i) \right] = n\text{Var}[g(X)] \quad (9)$$

**Theorem 3.** Let  $X_1, \dots, X_n$  be a random sample from  $F$ , and assume that  $X \sim F$ . Suppose that  $\text{Var}[X] = \sigma^2 < \infty$  and let  $\mu := \mathbb{E}[X]$ . Then:

1.  $\mathbb{E}[\bar{X}] = \mu$
2.  $\text{Var}[\bar{X}] = \sigma^2/n$ .
3.  $\mathbb{E}[s^2] = \sigma^2$ .

**Remark.** The reason why we divided by  $(n-1)$  and not  $n$  in  $s^2$  was because we wanted property (3) to be true.

*Proof.* 1. Obvious.

2.

$$\text{Var}[\bar{X}] = \text{Var} \left[ \frac{1}{n} \sum_{i=1}^n x_i \right] = \frac{1}{n^2} \text{Var} \left[ \sum_{i=1}^n x_i \right] = \frac{1}{n^2} n\sigma^2 = \sigma^2/n$$

3.

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2 \right] &\stackrel{\text{Thm 1.2}}{=} \mathbb{E} \left[ \frac{1}{(n-1)} \left[ \sum_{i=1}^n x_i^2 - n\bar{x} \right] \right] \\ &= \frac{1}{(n-1)} \left[ n\mathbb{E}[X^2] - n\mathbb{E}[\bar{X}]^2 \right] \\ &= \frac{1}{(n-1)} \left[ n(\sigma^2 + \mu^2) - n \left( \frac{\sigma^2}{n} + \mu^2 \right) \right] \\ &= \sigma^2 \end{aligned}$$

□

**Example 4** (Simon Newcomb trying to measure the speed of light, 1835-1909). In 1882, he attempted to carry out measurements to determine the speed of light. He collected 66 data points. *Question: given the data, what would be the model for the distribution of the speed of light? How would you describe the stochastic mechanism?*

We have  $n$  random variables,  $X_1, \dots, X_n$ , where  $n = 66$ . The parametric statistical model is:

$$\begin{aligned} \mathcal{F} &= \{ N(\mu, \sigma^2), \mu \in \mathbb{R}^2, \sigma^2 > 0 \} \\ T_i &= (24800 + X_i) \cdot 10^{-9} \end{aligned}$$

We think that a normal distribution is reasonable. That is,  $X_i \sim N(\mu, \sigma^2)$ . We will write each  $X_i$  as the sum of the mean and an error term,  $\varepsilon$ :

$$X_i = \mu + \varepsilon_i$$

where we assume  $\varepsilon_i \sim N(0, \sigma^2)$ . The  $\varepsilon_i$  has some distribution function with CDF  $F_0$  with  $\mathbb{E}[\varepsilon_i] = 0$ . We can express the distribution  $\mathcal{F}$  in an alternative way in terms of shifts:

$$\mathcal{F} = \{F_0(\cdot - \mu), \mu \in \mathbb{R}, F_0 \text{ is a CDF with expectation zero} \}$$

This gives us something that we call a **semi-parametric model**. We then obtain the  $\bar{x}$ , sample mean. We can use this to estimate  $\mu$ .

*Question: How confident are we in the obtained sample mean?* In order to answer this question, we need to know something about the sampling distribution of the statistics. We can estimate the variance of the sample mean using the previous theorem and the assumption about the errors being normally distributed as follows. Recall that the sample variance is given by:

$$\frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2$$

So, the expected value of  $s^2$  is:

$$\mathbb{E}[s^2] = \text{Var}[X] = \sigma^2 \text{ (if the errors are normally distributed)}$$

and,

$$\text{Var}[\bar{X}] = \sigma^2/n$$

This means that the uncertainty of  $\bar{X}$  depends on the underlying distribution since the  $\sigma^2$ 's are the variances of the  $X_i$ 's.

**Theorem 4.** Suppose that  $X_1, \dots, X_n$  is a random sample from an underlying distribution  $F$ . Let  $X \sim F$ . Suppose also that  $X$  has a moment generating function  $M_X(t)$  for  $t \in I$ . Then, the moment generating function of  $\bar{X}$  is:

$$M_{\bar{X}}(t) = \{M_X(t/n)\}^n, \quad t/n \in I \quad (10)$$

*Proof.* The proof follows from the IID property of the random variables  $X_1, \dots, X_n$ .

$$\begin{aligned} \text{MGF} &= \mathbb{E}[e^{t\bar{X}}] = \mathbb{E}[e^{t/n(X_1 + \dots + X_n)}] \\ &= \mathbb{E}\left[\prod_{i=1}^n e^{t/n X_i}\right] \\ &= \prod_{i=1}^n \mathbb{E}[e^{t/n X_i}] \\ &= \{M_X(t/n)\}^n \end{aligned}$$

□

The next example will give us some concrete examples of applying the previous theorem.

**Example 5.** 1. Let  $F = N(\mu, \sigma^2)$ . So, we know that  $X \sim N(\mu, \sigma^2)$ . From Math 356 we know the moment generating function is:

$$M_X(t) = e^{t\mu + \frac{1}{2}\sigma^2 t^2} \quad t \in \mathbb{R}$$

Invoking the theorem and simplifying:

$$\begin{aligned} M_{\bar{X}}(t) &= \left(e^{\frac{t}{n}\mu + \frac{1}{2}\sigma^2 \frac{t^2}{n^2}}\right)^n \quad t \in \mathbb{R} \\ &= e^{t\mu + \frac{1}{2}\sigma^2 \frac{t^2}{n}} \end{aligned}$$

since the MGF uniquely determines the underlying distribution, this gives us that:

$$\bar{X} \sim N(\mu, \sigma^2/n)$$

2. If  $X \sim \text{Binomial}(m, p)$ :

$$M_X(t) = (1 - p + pe^t)^m$$

$$M_{\bar{X}}(t) = (1 - p + pe^{t/n})^{m \cdot n}$$

a) A modification of the  $\bar{X}$  will be distributed binomially. Namely,  $n \times \bar{X}$  will get rid of the  $n$  in the denominator:

$$M_{n \times \bar{X}}(t) = \mathbb{E} \left[ e^{nt\bar{X}} \right] = (1 - p + pe^t)^{m \cdot n}$$

$$\Rightarrow n \times \bar{X} \sim \text{binomial}(m \cdot n, p).$$

3. If  $X \sim \text{Gamma}(\alpha, \beta)$ :

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$$

$$M_X(t) = (1 - t\beta)^{-\alpha}, \quad t < \frac{1}{\beta}$$

$$M_{\bar{X}}(t) = \left(1 - \frac{t}{n}\beta\right)^{-\alpha n}, \quad t < \frac{n}{\beta}$$

$$\Rightarrow \bar{X} \sim \text{Gamma}\left(\alpha \cdot n, \frac{\beta}{n}\right)$$

## 2.1 Sampling from the Normal Population

The setup for this section will be as follows: let  $X_1, \dots, X_n$  be iid from  $N(\mu, \sigma^2)$ .

**Theorem 5.** Let  $X_1, \dots, X_n$  be iid from  $N(\mu, \sigma^2)$ . Then:

1.  $\bar{X} \sim N(\mu, \sigma^2)$ . (Shown using the MGF).
2.  $\bar{X}$  and  $s^2$  are independent.

*Proof.* WLOG for (a), we can assume that  $\mu = 0$  and  $\sigma^2 = 1$ , since we standardise random variables. Why? The standardisation of a random variable  $X_i$ , denoted by  $X_i^*$ , is an affine transformation given by:

$$X_i^* = \frac{X_i - \mu}{\sigma}$$

$X_i^*$  is still normally distributed, which can be shown using an MGF argument. Thus:

$$\bar{X}_i^* = \frac{\bar{X} - \mu}{\sigma}$$

$(\bar{X} = \sigma \bar{X}^* + \mu)$ . Moreover, for  $(s^*)^2$ :

$$(s^*)^2 = \frac{1}{(n-1)} \sum_{i=1}^n (X_i^* - \bar{X}^*)^2 = \frac{1}{(n-1)} \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma^2} = \frac{s^2}{\sigma^2}$$

This justifies why we can say WLOG. To show that (b) holds when  $\mu = 0$  and  $\sigma^2 = 1$ , we will play with  $s^2$ :

$$s^2 = \frac{1}{(n-1)} \left[ \sum_{i=2}^n (X_i - \bar{X})^2 + (X_1 - \bar{X})^2 \right]$$

because we also know that

$$\sum_{i=1}^n (X_i - \bar{X}) = 0 \Rightarrow (X_1 - \bar{X}) = -\sum_{i=2}^n (X_i - \bar{X})$$

This implies that we can re-write  $s^2$  as:

$$s^2 = \frac{1}{(n-1)} \left[ \sum_{i=2}^n (X_i - \bar{X})^2 + \left( \sum_{i=2}^n (X_i - \bar{X}) \right)^2 \right] := h(X_2 - \bar{X}, \dots, X_n - \bar{X})$$

$h$  is clearly measurable. We now have a function of only  $n-1$  random variables, which is why we normalise  $s^2$  by  $(n-1)$ .

**Lemma 6** (Core of the argument). If  $X_1, \dots, X_n$  are iid  $N(0, 1)$ , then

$$\bar{X} \perp (X_2 - \bar{X}, \dots, X_n - \bar{X}) \quad (11)$$

From this lemma 1.14, we can then immediately conclude that  $\bar{X} \perp s^2$ .

*Proof.* This is where we use the normality assumption. Define a one-to-one function  $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$ :

$$g(x_1, \dots, x_n) \mapsto (\bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x})$$

Since  $g$  is one-to-one, we can invert it:

$$g^{-1}(y_1, \dots, y_n) \mapsto \left( g_n - \sum_{i=2}^n y_i, y_n + y_1, \dots, y_n + y_1 \right)$$

Need to calculate the Jacobian of the transformation:

$$\text{Jac}(g) = \begin{bmatrix} 1 & -1 & \cdots & -1 \\ 1 & 1 & \cdots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 1 & 0 & \cdots & 1 \end{bmatrix}$$

$\det(\text{Jac}(g)) = n$ . So, by the transformation laws:

$$f_{Y_1, \dots, Y_n}(y_1, \dots, y_n) = f_{(X_1, \dots, X_n)}(g^{-1}(y_1, \dots, y_n)) |\text{Jac}(g)|$$

Here,  $Y_1 = \bar{X}$ ,  $Y_2 = X_2 - \bar{X}$ , ...,  $Y_n = X_n - \bar{X}$ , and so by the IID of the  $X_i$ :

$$\begin{aligned} f_{(X_1, \dots, X_n)}(x_1, \dots, x_n) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-x_i^2/2} \end{aligned}$$

where the second inequality follows by the standard normal assumption.

$$\begin{aligned} f_{(Y_1, \dots, Y_n)}(y_1, \dots, y_n) &= \left( \frac{1}{\sqrt{2\pi}} \right)^n \exp \left\{ -\frac{1}{2} \left( y_1 - \sum_{i=2}^n y_i \right)^2 + (y_2 + y_1)^2 + \dots + (y_n + y_1)^2 \right\} \\ &= \left( \frac{1}{\sqrt{2\pi}} \right)^n \exp \left\{ \frac{-ny_1^2}{2} \right\} \prod_{i=2}^n \exp \left\{ -\frac{1}{2} \left( \sum_{i=2}^n y_i^2 - \left( \sum_{i=2}^n y_i \right)^2 \right) \right\} \end{aligned}$$

(the cross terms will drop out). Thus, we have factored the densities of the  $Y$ 's into the product of two functions with difference dependencies ( $y_1$  vs.  $y_2, \dots, y_n$ ). Thus,  $Y_1 \perp (Y_2, \dots, Y_n)$   $\square$

By the following theorem from Chapter 4 of the textbook, we obtain the desired result:

**Theorem 7** (Generalisation of Theorem 4.3.2). Let  $X_1, \dots, X_N$  be independent random vectors. Let  $g_i(x_i)$  be only a factor of  $x_i$ ,  $i = 1, \dots, n$ . Then, the random variables  $U_i := g_i(X_i)$ ,  $i = 1, \dots, n$  are mutually independent.

□

**Definition 4** (Chi Squared Distribution with  $v$  degrees of freedom).

$$f_v(x) := \frac{1}{2^{v/2}\Gamma(v/2)} x^{v/2-1} e^{-x/2} \quad (12)$$

$x > 0$ .

**Remark:**  $\chi_v^2$  is a Gamma( $v/2, 2$ ) distribution (special case of the gamma distribution). Recall from earlier that the MGF of  $\chi_v^2$  is  $(1 - 2t)^{-v/2}$  for  $t < 1/2$ .

**Lemma 8** (Facts about  $\chi^2$ ). 1. If  $X \sim \chi_v^2$ , then  $\mathbb{E}[X] = v$  and  $\text{Var}[X] = 2v$   
 2. If  $X_1 \sim \chi_{v_1}^2$  and  $X_2 \sim \chi_{v_2}^2$ ,  $X_1 \perp X_2$ , then  $X_1 + X_2 \sim \chi_{(v_1+v_2)}^2$   
 3. If  $X \sim N(0, 1)$  then  $X^2 \sim \chi_1^2$ . (Proof of this one is on assignment 1).

The following theorem is very important since it leads to the chi squared test.

**Theorem 9.** Let  $X_1, \dots, X_n$  be a random sample from  $N(\mu, \sigma^2)$ . Then:

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi_{(n-1)}^2 \quad (13)$$

Reality check:

$$\mathbb{E}\left[\frac{(n-1)s^2}{\sigma^2}\right] = (n-1)$$

which implies that:

$$\frac{(n-1)}{\sigma^2} \mathbb{E}[s^2] = (n-1) \Rightarrow \mathbb{E}[s^2] = \sigma^2$$

We can elegantly prove this using moment generating functions:

*Proof.* From the preceding Lemma and the first theorem of the section, we have that we can standardise:

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

Squaring this, we obtain:

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

So, summing the random variables gives:

$$\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 \sim \chi_n^2$$



Therefore, by adding and subtracting the sample mean and simplifying:

$$\begin{aligned}\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 &= \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} + \frac{n(\bar{X} - \mu)^2}{\sigma^2} \\ &= \frac{(n-1)s^2}{\sigma^2} + \underbrace{\frac{n(\bar{X} - \mu)^2}{\sigma^2}}_{\sim \chi_1^2}\end{aligned}$$

From the first theorem of the section, these are independent. We cannot subtract to solve for the quantity that we are interested in, so we will work with moment generating functions. The MGF of the left hand side is:

$$(1 - 2t)^{n/2}, \quad t < 1/2$$

and the MGF of the right hand side is:

$$M_{\frac{(n-1)s^2}{\sigma^2}}(t) \cdot (1 - 2t)^{-1/2}, \quad t < 1/2$$

equating these, we obtain:

$$M_{\frac{(n-1)s^2}{\sigma^2}}(t) = (1 - 2t)^{-(n-1)/2}, \quad t < 1/2$$

However, this is the MGF of  $\chi_{(n-1)}^2$ , and since the MGF uniquely determines the distribution, we obtain that  $\frac{(n-1)s^2}{\sigma^2} \sim \chi_{(n-1)}^2$ , which is what we wanted to show.  $\square$

**Observe:** this proof heavily relies on the normality of the distribution, independence, and the chi squared.

## 2.2 Standardisation

**Motivation:** quality control problem. Let  $\bar{X}$  be the random variable estimating the underlying mean. Say that we set a quality cutoff,  $\mu_0$ , and we want to answer the question:

$$\text{does } \mu = \mu_0?$$

What if we try to standardise:

$$\frac{\bar{X} - \mu_0}{\sigma/n} \sim N(0, 1)$$

and then compare the quantities? The problem here is that this is not a statistic. We do not know what  $\sigma$  is. However, if you don't know something, estimate it:

$$\sqrt{n} \frac{\bar{X} - \mu_0}{\sqrt{s^2}}$$

The problem with this approach is that it is not distributed  $N(0, 1)$ . Especially when the sample is not too large. This motivates the following definition:

**Definition 5** (Student-t distribution with  $v$  degrees of freedom). has the density:

$$f_v(x) := \frac{\Gamma\left(\frac{v+1}{2}\right)}{\sqrt{v\pi}\Gamma\left(\frac{v}{2}\right)} \left(1 + \frac{x^2}{v}\right)^{-(v+1)/2}$$

where  $x \in \mathbb{R}$  and  $v > 0$ .

**Remark:** if you set  $v = 1$ , you obtain the Cauchy Lorentz distribution.  $\mathbb{E}[X] = \infty$ . Also observe that the student-t is a heavy-tailed distribution.

**Lemma 10.** Let  $X \sim t_v$ . Then:

1.  $\mathbb{E}[X] = 0$ , provided that  $v > 1$  (otherwise it does not exist).
2.  $\text{Var}[X] = v/(v-2)$  when  $v > 2$  (otherwise,  $\text{var}[X]$  does not exist).
3. (Assignment): if  $Z \sim N(0, 1)$ ,  $V \sim \chi_v^2$ ,  $Z \perp V$ , then:

$$\frac{Z}{\sqrt{v/v}} \sim T_v \quad (14)$$

This is the **most important part of the distribution**. You can prove it using the transformation theorem for densities.

**Theorem 11.** Let  $X_1, \dots, X_n$  be a random sample from  $N(\mu, \sigma^2)$ . Then:

$$\frac{\bar{X} - \mu}{\sqrt{s^2/n}} \sim t_{(n-1)}$$

The proof is pretty obvious using the lemma:

*Proof.* We will express the ratio as a standard normal  $Z$ . We already have that:

$$\sqrt{n} \frac{\bar{X} - \mu}{\sigma} \sim N(0, 1) \perp \frac{s^2(n-1)}{\sigma^2} \sim \chi_{(n-1)}^2$$

So, taking the ratio in the form of what is given to us in the previous lemma gives:

$$\frac{\sqrt{n} \frac{\bar{X} - \mu}{\sigma}}{\sqrt{\frac{s^2(n-1)}{\sigma^2}}} = \frac{\sqrt{n}(\bar{X} - \mu)}{\sqrt{s^2}} \sim t_{(n-1)}$$

where we obtain the distribution from the previous lemma. □

The student-t model forms a nice statistical model for certain types of data.

## 2.3 F-Test Basis

If we want to compare quality, we need to standardise by variance. The following section helps answer the question, **is the variance between two samples the same?**

**Definition 6** (The Fischer-Snedecor's F Distribution with Two Parameters  $v_1$  and  $v_2$  degrees of freedom). This distribution is denoted by  $F_{v_1, v_2}$ . This is the distribution of the following:

$$\frac{V_1/v_1}{V_2/v_2}$$

where  $V_1 \sim \chi_{v_1}^2$ ,  $V_2 \sim \chi_{v_2}^2$ , and  $V_1 \perp V_2$ .

This will lead to the  $F$ -test.

**Theorem 12.** If  $X_1, \dots, X_n$  is a random sample from  $N(\mu_x, \sigma_x^2)$ , and  $Y_1, \dots, Y_m$  is a random sample from  $N(\mu_y, \sigma_y^2)$ , and the two samples are *independent*,  $s_x^2$  and  $s_y^2$  are the sample variances, then:

$$\frac{s_x^2/\sigma_x^2}{s_y^2/\sigma_y^2} \sim F_{n-1, m-1}$$

*Proof.* From the previous result and independence, we have:

$$\frac{(n-1)s_x^2}{\sigma_x^2} \sim \chi_{(n-1)}^2$$

$$\frac{(m-1)s_y^2}{\sigma_y^2} \sim \chi_{(m-1)}^2$$

are independent. If we divide by the degrees of freedom and invoke the definition of the F distribution, we obtain the desired result.

$$\frac{s_x^2/\sigma_x^2}{s_y^2/\sigma_y^2} \sim F_{(n-1),(m-1)} \quad (15)$$

□

This forms the basis of the F-test. If we want them to have the same variances, then  $s_y^2$  and  $s_x^2$  had better be close. The way to assess this is to look at ratios, and then the  $\sigma^2$ 's will drop out and we can then test hypotheses. Moreover, if  $\sigma_x^2 = \sigma_y^2$ , then:

$$\frac{s_x^2}{s_y^2} \sim F_{(n-1),(m-1)}$$

we will later use this to construct the so-called **F-test**.

## 2.4 Limiting Sample Distributions (Asymptotics)

What happens as  $n \rightarrow \infty$ ? These questions are answered by **Weak Law of Large Numbers** and **Convergence in Distribution**.

More precisely, let  $X_1, \dots, X_n$  be iid from F. Then, define  $T_n := T(X_1, \dots, X_n)$  be a real-valued statistic. Then:

**Q1:** Does  $T_n$  converge in probability to an estimator  $\theta \in \mathbb{R}$ ?

$$\forall \varepsilon > 0 \quad \mathbb{P}[|T_n - \theta| > \varepsilon] \rightarrow 0 \quad \text{as } n \rightarrow \infty \quad (16)$$

**Q2:** What happens to the distribution as  $n \rightarrow \infty$ ? In other words, if  $(r_n)$  is a sequence of real numbers, typically such that  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ , does

$$r_n(T_n - \theta) \xrightarrow{d} T \quad (17)$$

What the  $(r_n)$  is doing is that it is “zooming” into  $T_n \rightarrow \theta$ . **Idea:** when  $n$  is “large enough”:

$$r_n(T_n - \theta) \xrightarrow{d} T_n \frac{T}{r_n} + \theta$$

where the final term is called the **location and scale model**. This may have a nice distribution. The distribution of  $\frac{T}{r_n} - \theta$  is often called the **large-sample** or **(asymptotic) distribution of  $T_n$**  or the **limiting distribution** of  $T_n$ . This is fundamental for quantifying uncertainty and for hypothesis testing.

### 2.4.1 Question 1

First we will have a refresher from Math 356. The prime tool for convergence in probability is the Weak Law of Large Numbers.

**Theorem 13** (Weak Law of Large Numbers). Let  $X_1, X_2, \dots$  be iid random variables with  $\mathbb{E}[X_i] = \mu$  and  $\text{Var}[X_i] = \sigma^2 < \infty$ . Define

$$\overline{X}_n := \frac{1}{n} \sum_{i=1}^n X_i \quad (18)$$

then,  $\overline{X}_n$  converges in probability to the expected value of  $X$ .

**Theorem 14** (Continuous Mapping Theorem). If  $T_n \xrightarrow{P} T$  and  $g$  is continuous on the set  $C$  such that  $\mathbb{P}[T \in C] = 1$ , then  $g(T_n) \xrightarrow{P} g(T)$ .

In particular, if  $T_n \xrightarrow{P} \theta$  and if  $g$  is continuous at  $\theta$ , then  $g(T_n) \xrightarrow{P} g(\theta)$ .

**Example 6** (Another justification for the sample variance). Let  $X_1, \dots, X_n$  be a random sample from  $X$  (i.e., from  $F$  with  $X \sim F$ ), and assume that  $\mathbb{E}[X^2] < \infty$ . Then, by the WLLN,  $\overline{X} \xrightarrow{P} \mathbb{E}[X]$ . For the sample variance:

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

Applying the weak law of large numbers to the first term and the continuous mapping theorem to the second term gives us that the difference will converge in probability to:

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

since the square root is continuous:

$$s = \sqrt{s^2} \xrightarrow{P} \sqrt{\text{Var}[X]}$$

**Example 7.** Let  $X_1, \dots, X_n$  be a random sample from  $X$ . Suppose that we are after  $\mathbb{P}[X \in A]$ . This can be estimated by the **empirical probability**, which counts how many  $x$ 's fall into  $A$ . Set  $Z_i := \chi_{X_i \in A}$ . Then, the  $Z_i$  are iid Bernoulli random variables. Using the Weak Law of Large Numbers:

$$\mathbb{P}_n \xrightarrow{P} \mathbb{P}[X \in A]$$

which is the expected value of a Bernoulli random variable. By the Strong Law of Large numbers:

$$\mathbb{P}_n \rightarrow \mathbb{P}[X \in A] \text{ a.s.}$$

Since  $F$  is the CDF of  $X$ , we can try to learn it from the data. For  $x \in \mathbb{R}$ ,  $F(x) = \mathbb{P}[X \leq x]$ . The **empirical distribution function**  $F_n$  is given by:

$$F_n(x) := \frac{1}{n} \sum_{i=1}^n \chi_{X_i \leq x} \quad (19)$$

This is a central object in mathematical statistics.

Given an observed sample  $x_1, \dots, x_n$ , a sample empirical CDF is given by:

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \chi_{x_i \leq x} \quad (20)$$

This is a CDF itself. It is a CDF with a discrete distribution with support  $\{x_1, \dots, x_n\}$  and  $\mathbb{P}[x_i] = 1/n$ . By the Weak Law of Large Numbers:

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \chi_{X_i \leq x} \xrightarrow{P} F(x)$$

for  $\forall x \in \mathbb{R}$ . Actually, by the strong law of large numbers:

$$F_n(x) \rightarrow F(x) \text{ a.s.}$$

for any  $x \in \mathbb{R}$ .

**Theorem 15** (Glivenlco-Cantelli Theorem). We have uniform convergence:

$$\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \rightarrow 0 \text{ as } n \rightarrow \infty \quad (21)$$

Notice that:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \int x dF_n(x) \quad (22)$$

## 2.4.2 Question 2

We will need the following tools from probability: the Central Limit Theorem and Slutsky's theorem.

**Theorem 16.** Let  $Z_1, \dots, Z_n$  be iid with  $\mathbb{E}[Z_i] < \infty$ ,  $\mathbb{E}[X_1] = \mu$ ,  $\text{var}[Z_1] < \infty$ , and  $\text{var}[Z_i] = \sigma^2$ . Then, the **Central Limit Theorem** states:

$$\sqrt{n} \frac{\bar{Z} - \mu}{\sigma} \xrightarrow{CLT} N(0, 1) \quad (23)$$

or

$$\sqrt{n}(\bar{Z} - \mu) \xrightarrow{CLT} N(0, \sigma^2) \quad (24)$$

**Theorem 17** (Slutsky's Theorem). Assume that  $T_n \xrightarrow{d} T$ ,  $Y_n \xrightarrow{p} c$  where  $c \in \mathbb{R}$ . Then:

1.  $T_n + Y_n \xrightarrow{d} T + c$
2.  $T_n \cdot Y_n \xrightarrow{d} T \cdot c$
3.  $T_n/Y_n \xrightarrow{d} T/c$  if  $c \neq 0$ .

**Remarks:** If  $T_n \xrightarrow{d} c$  where  $c \in \mathbb{R}$ , then  $T_n \xrightarrow{p} c$ . Moreover, we can say something with the CMT. If  $T_n \xrightarrow{d} T$ ,  $g$  is continuous on  $C$  with  $\mathbb{P}[T \in C] = 1$ . Then,  $g(T_n) \xrightarrow{d} g(T)$ .

**Important remark:** if  $(r_n)$  is a sequence of numbers,  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Then, if  $r_n(T_n - \theta) \xrightarrow{d} T$ , then we have:

$$T_n - \theta = \underbrace{r_n(T_n - \theta)}_{\xrightarrow{d} T} \underbrace{\frac{1}{r_n}}_{\rightarrow 0} \xrightarrow{\text{slutsky's}} 0$$

So,  $T_n - \theta \xrightarrow{p} 0$ , or  $T_n \xrightarrow{d} \theta$ . This means that the second statement implies the first statement.

**Example 8** (Alternative proof of the CLT). Let  $X_1, \dots, X_N$  be iid from  $F$ ,  $\mathbb{E}[X] = \mu$ ,  $\text{Var}[X] = \sigma^2 < \infty$  (this assumption is very important!)

$$\sqrt{n} \frac{\bar{X} - \mu}{\sqrt{s^2}} = \underbrace{\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}}_{\text{By vanilla CLT, this goes to } N(0,1)} \cdot \underbrace{\frac{\sigma}{\sqrt{s^2}}}_{\xrightarrow{p} 1} \quad (25)$$

and so by Slutsky's theorem

$$\xrightarrow{d} N(0, 1)$$

**Corollary 1.** Assume that each  $T_n \sim t_n$  (student-t with  $n$  degrees of freedom). Then  $T_n \xrightarrow{d} N(0, 1)$ .

**Example 9.** Suppose  $X_1, \dots, X_n$  are iid from Bernoulli( $p$ ). Suppose that we want to estimate  $\text{var}[X_1] = p(1 - p)$ . **Q:** could we estimate it by:

$$\bar{X}(1 - \bar{X}) \xrightarrow[\text{P}]{\text{CMT}} p(1 - p)$$

More precisely, can we find a sequence  $(r_n)$  such that

$$(r_n)(\bar{X}(1 - \bar{X}) - p(1 - p)) \xrightarrow{d} ??$$

This is not exactly the CMT; this example provides the motivation for developing the **delta method**. Suppose that we know that  $r_n(T_n - \theta) \xrightarrow{d} T$  and we are interested in  $g(T_n) - g(\theta)$ . What could we do? We could use the Taylor Expansion:

$$\begin{aligned} g(T_n) - g(\theta) &\approx g'(\theta)(T_n - \theta) \\ (r_n)(g(T_n) - g(\theta)) &\approx g'(\theta)\{r_n(T_n - \theta)\} \end{aligned}$$

By Slutsky's theorem,  $r_n(T_n - \theta)$  converges in distribution to  $T \cdot g'(\theta)$ .

**Theorem 18** (Delta Method). Suppose  $(r_n)(T_n - \theta) \xrightarrow{d} T$ , where  $(r_n)$  is a real sequence with  $r_n \rightarrow \infty$ . Let  $g$  be a function and  $T_n$  takes values in the domain of  $g$ , and assume that  $g$  is *differentiable* at  $\theta$ . Then:

$$r_n(g(T_n) - g(\theta)) \xrightarrow{d} T \cdot g'(\theta) \quad (26)$$

Proof will be postponed. Back to the example. Here:

$$\begin{aligned} g(x) &= x(1 - x) \\ g'(x) &= 1 - 2x \end{aligned}$$

By the central limit theorem

$$(\bar{X} - p)\sqrt{n} \xrightarrow{d} N(0, p(1 - p))$$

Let  $\sqrt{n} := (r_n)$ . Then, by the delta method

$$\begin{aligned} \sqrt{n}(\bar{X}(1 - \bar{X}) - p(1 - p)) &\xrightarrow{d} N(0, p(1 - p)) \cdot (1 - 2p) \\ &= N(0, (1 - 2p)^2 p(1 - p)) \end{aligned}$$

When  $p = 1/2$ , the statement is uninteresting since the derivative is zero. We will have both convergence in probability and distribution to zero, which doesn't give us any information really.

We are now ready to prove the delta method.

*Proof.* This proof uses several common arguments that you should learn for stats :-). Observe that by the continuous mapping theorem,  $T_n \xrightarrow{p} \theta$ , also  $g(T_n) \xrightarrow{p} g(\theta)$ . Define the following function:

$$f : \mathbb{R} \rightarrow \mathbb{R}$$

$$h \mapsto \begin{cases} \frac{g(\theta+h)-g(\theta)}{h} - g'(\theta) & \text{if } h \neq 0 \\ 0 & \text{if } h = 0 \end{cases}$$

Here, the interesting domain is small and is about zero. Observe that  $f$  is continuous at zero. Now we can use the continuous mapping theorem:

$$f(T_n - \theta) \xrightarrow{p} f(0) = 0$$

since the first term is equal to

$$\frac{g(T_n) - g(\theta)}{T_n - \theta} - g'(\theta)$$

By Slutsky's theorem:

$$\begin{aligned} & r_n(g(T_n) - g(\theta)) - g'(\theta)r_n[T_n - \theta] \\ &= \underbrace{r_n[T_n - \theta]}_{\xrightarrow{d} T} \underbrace{f(T_n - \theta)}_{\xrightarrow{p} 0} \xrightarrow{d} T \cdot 0 = 0 \end{aligned}$$

where the last convergence follows from Slutsky's theorem. Therefore:

$$r_n(g(T_n) - g(\theta)) - g'(\theta)r_n[T_n - \theta] \xrightarrow{p} 0$$

Now:

$$r_n[g(T_n) - g(\theta)] = \underbrace{r_n[g(T_n) - g(\theta)] - g'(\theta)r_n(T_n - \theta)}_{\xrightarrow{p} 0} + \underbrace{g'(\theta)(T_n - \theta) \cdot r_n}_{\xrightarrow{d} g'(\theta)T \text{ (CMT)}}$$

and so by Slutsky's theorem, we obtain convergence in distribution to  $g'(\theta)T$ , which completes the proof.  $\square$

## 2.5 Order Statistics

**Definition 7** (Order Statistics). The order statistics of a random sample  $X_1, \dots, X_n$  are the sample values placed in ascending order. They are denoted by  $X_{(1)}, \dots, X_{(n)}$ . In other words, they are random variables that satisfy  $X_{(1)} \leq \dots \leq X_{(n)}$ . In particular:

$$\begin{aligned} X_{(1)} &= \min_{1 \leq i \leq n} X_i \\ X_{(2)} &= \text{second smallest } X_i \\ &\vdots \\ X_{(n)} &= \max_{1 \leq i \leq n} X_i \end{aligned}$$

**Definition 8** (Sample Range). The sample range is defined as  $R := X_{(n)} - X_{(1)}$ .

**Definition 9** (Sample Median). The sample median, denoted by  $M$ , is the number such that approx one half of the observations are less than  $M$  and one half are greater. It can be defined in terms of the order statistics as:

$$M := \begin{cases} X_{(n+1)/2} & \text{if } n \text{ is odd} \\ \frac{X_{(n/2)} + X_{(n/2+1)}}{2} & \text{if } n \text{ is even} \end{cases}$$

### 3 Point Estimation

Let  $X_1, \dots, X_n$  be a random sample from  $F$ , and we will *assume*:

$$F \in \mathcal{F} = \{F_\theta, \theta \in \Theta\}$$

This means that our distribution is known up to a parameter  $\theta$ .  $\theta$  is an unknown vector of parameters.  $\Theta$  is called the **parameter space**,  $\Theta \subseteq \mathbb{R}^k$ . **Question:** where does  $\Theta$  come from? This is the “art” of statistical modelling. The objective in this section is to learn  $\theta$  from the data.

**Notation:**  $F_\theta$  is a CDF.  $P_\theta$  is the corresponding probability measure.  $F_\theta$  has a density / probability mass function, which we will denote by  $f(x, \theta)$ . This notation emphasises the dependence on  $\theta$ .

**Example 10.** Assume that we have the normal model. Then:

$$\mathcal{F} = \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$$

Here,  $\theta = (\mu, \sigma^2)$  and the parameter space  $\Theta = \mathbb{R} \times ]0, \infty[$ .

**Definition 10** (Point Estimator). A **point estimator** is any statistic  $T(X_1, \dots, X_n)$  (function of only the data) taking values in  $\mathbb{R}^k$  that has been constructed with the aim of estimating  $\theta$ . The observed value  $T(x_1, \dots, x_n)$  is called an **estimate** of  $\theta$ , which is a concrete number.

This definition is vague since we do not want to eliminate important estimates a priori. Often, estimators are denoted with  $\hat{\theta}, \tilde{\theta}, \theta_n$ , or  $\hat{\theta}_n$ .

#### 3.1 Ways We Can Construct Estimators

We will explore two of these methods in depth in this class. The third we will briefly mention. The first two are frequentist methods and the final one is a bayesian method.

1. Method of moments
2. Method of maximum likelihood
3. Bayesian estimation method

##### 3.1.1 Method of moments

This is the oldest method. It was developed in the 19th century by Karl Pearson. It has been what we heuristically have been doing all along. Suppose that  $X \sim F_\theta$ . Then, as an idea we could first calculate the  $j$ th moment  $\mu_j$  of  $X$  for  $j = 1, \dots, k$ , that is:

$$\mu_j = \mathbb{E}_\theta [X^j] = \int x^j dF_\theta(x) = \int x^j dP_\theta = \int x^j f(x; \theta) dx$$

**How would we estimate  $\mu_j$  from the data?** By the weak law of large numbers,  $\mu_j$  can be estimated by the  $j$ th sample moment:

$$\mu_j = \frac{1}{n} \sum_{i=1}^n X_i^j$$

observe that this  $\mu_j$  is a function of  $\theta$ . Thus, the general idea of the method of moments is:

Compute the  $j$ th moment of the parametric model. Then, solve for  $\theta$ :

1. Calculate the population moments  $\mu_j$ ,  $j = 1, \dots, k$  (the first  $k$  moments), and observe that  $\mu_j$  is a function of  $\theta$ .
2. Set  $\mu_j = m_j$  for  $j = 1, \dots, k$  and solve for  $\theta$ . Start with  $j = 1$  and see where it takes us.



Note that sometimes you may need more moments than the first  $k_1$  since some moments may not actually depend on  $\theta$  (they could be zero, for example).

**Example 11** (Method of moments for the normal distribution). Let  $\mathcal{F} = \{N(\mu, \sigma^2) \mid \mu \in \mathbb{R}, \sigma^2 > 0\}$ . Here, since we have two parameters to estimate, we need at least two moments. We have:

$$\mathbb{E}[X] = \mu \quad \mathbb{E}[X^2] = \text{var}[X] + (\mathbb{E}[X])^2 = \sigma^2 + \mu^2$$

Now equate the population moments and solve for  $\theta$ . That is, we need to solve:

$$m_1 = \frac{1}{n} \sum_{i=1}^n X_i = \mu$$

$$m_2 = \frac{1}{n} \sum_{i=1}^n X_i^2 = \sigma^2 + \mu^2$$

and this gives us

$$\hat{\mu} = m_1$$

$$\hat{\sigma}^2 = m_2 - \hat{\mu}^2$$

Substituting in what all of this means gives us:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \left( \frac{1}{n} \sum_{i=1}^n X_i \right)^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{(n-1)}{n} s^2$$

Hence, the estimator (**method of moments estimator**) for  $(\mu, \sigma^2)$  is:

$$(\hat{\mu}, \hat{\sigma}^2) = \left( \bar{X}, \frac{(n-1)}{n} s^2 \right)$$

**Q: What are some advantages and disadvantages of the method of moments?**

Some advantages include:

1. It is intuitive.
2. Population moments are easy to estimate and we know how they behave asymptotically.
  - a) Weak law of large numbers and the CLT.
3. Asymptotic properties of  $\hat{\theta}$  can often be derived using the CMT and the Delta Method.

Some disadvantages include (“the bitter reality”):

1. There is a systematic bias for  $n$  small, which is encoded by the  $(n-1)/n$  term.
2. These estimators are often not optimal.

**Example 12** (Method of moments for binomial). Let our model be  $\mathcal{F} = \{B(N, p), p \in [0, 1], N = \{1, \dots\}\}$ . There are two cases:  $N$  can either be known or unknown. If  $N$  is known, then it is easy. In this case,  $\theta = p$  and  $\Theta = [0, 1]$ . In this case:

$$\mu_1 = Np$$

$$m_1 = Np \Rightarrow \hat{p} = \frac{m_1}{N} = \frac{\bar{x}}{N}$$

The not so easy case is when both  $p$  and  $N$  are unknown. We have the following system of equations:

$$\mu_1 = Np$$

$$\mu_2 = Np(1-p) + (Np)^2$$

We thus have the following system of equations that we need to solve:

$$\begin{aligned}m_1 &= Np \\m_2 &= Np - Np^2 + N^2p^2\end{aligned}$$

Solving for  $\hat{p}$  gives:

$$\hat{p} = \frac{m_1}{N}$$

Substituting this estimation into the second equation and simplifying gives:

$$\begin{aligned}m_2 &= m_1 - \frac{m_1^2}{N^2}N(n-1) \\ \Rightarrow m_2 &= m_1 + \frac{m_1^2}{N}(N-1) \\ \Rightarrow Nm_2 &= Nm_1 + m_1^2(N-1) \\ \Rightarrow N(m_2 - m_1 - m_1^2) &= -m_1^2\end{aligned}$$

What if  $X_1 = \dots = X_n = 0$ ? This event has non-zero probability, and so we must account for it. In that case, there is nothing more we can do. Otherwise, we obtain:

$$\begin{aligned}\hat{N} &= \begin{cases} \frac{-m_1^2}{m_2 - m_1 - m_1^2} \\ \text{undefined if } X_1 = \dots = X_n = 0 \end{cases} \\ \hat{p} &= \frac{-m_2 - m_2 - m_1^2}{m_1} \text{ or undefined in the case } X_1 = \dots = X_n = 0\end{aligned}$$

So, our method of moments (MoM) estimator of  $(p, N)$  is:

$$\left( \frac{-m_2 + m_1 + m_1^2}{m_1}, \frac{m_1^2}{-m_2 + m_1 + m_1^2} \right) = \left( \frac{\bar{X} - \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}{\bar{X}}, \frac{(\bar{X})^2}{\bar{X} - \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2} (*) \right)$$

Observations from this example:

- Estimating  $N$  was not very intuitive.
- There is no guarantee that the second component of  $(*)$  is an integer or non-negative. This is why we chose such a vague definition of an estimator; because we don't want to exclude this method.
- Estimating  $N$  for the binomial is quite difficult.

Here is an example where  $\theta$  determines the support of the random variable.

**Example 13.** Let  $\mathcal{F} = \{U(\cdot - \theta, \theta], \theta \in [0, \infty[ \}$ . Since the expected value is zero, we will need more moments. Recall the following from probability theory: if  $X \sim U[a, b]$ , then:

$$\mathbb{E}[X] = \frac{a+b}{2} \text{ and } \text{var}[X] = \frac{1}{12}(b-a)^2$$

then the moments are:

$$\begin{aligned}\mu_1 &= 0 \\ \mu_2 &= \frac{1}{2}(2\theta)^2 = 2\theta\end{aligned}$$

this is good since it is a genuine function of  $\theta$ . Set:

$$\mu_2 = m_2$$

and solve the following to obtain the method of moment estimator:

$$2\theta^2 = m_2 \Rightarrow \hat{\theta} = \sqrt{\frac{m_2}{2}} = \sqrt{\frac{\frac{1}{n} \sum_{i=1}^n X_i^2}{2}}$$

Observe that  $\theta$  could theoretically be zero. But, since this is a continuous distribution, the probability of  $\theta$  actually being zero is zero.

### 3.1.2 Method of Maximum Likelihood

This is more optimal, but it can be harder to obtain. Assume that  $F_\theta$  has a PDF or PMF for any  $\theta \in \Theta$ . Recall that we denoted this PMF by  $f(x; \theta)$  to emphasise the dependence on  $\theta$ . Before proceeding, a word of caution:

**The likelihood function is NOT RANDOM!!!**

**Definition 11** (Likelihood Function). Given that  $(X_1, \dots, X_n) = (x_1, \dots, x_n)$  have been observed, the function of  $\theta$  defined by  $L : \Theta \rightarrow [0, \infty[$ ,

$$L(\theta) := \prod_{i=1}^n f(x_i; \theta) \quad (27)$$

is called the **likelihood function** for a fixed  $x = (x_1, \dots, x_n)$  as a function of  $\theta$ .

**Example 14.** Suppose that  $x_1 = 1$ ,  $x_2 = 2$ ,  $x_3 = 2$ , and  $x_4 = 5$  and that we assume the Poisson model:

$$\mathcal{F} = \{P(\lambda), \lambda \in ]0, \infty[ \}$$

then, the likelihood function is:

$$L(\lambda; x) = e^{-\lambda} \frac{\lambda^1}{1!} + e^{-\lambda} \frac{\lambda^2}{2!} + e^{-\lambda} \frac{\lambda^3}{2!} + e^{-\lambda} \frac{\lambda^4}{5!}$$

We can think of the Likelihood function as a sort of “function summary” whose domain is the parameter space.

#### Interpretation of the Likelihood Function:

1. If  $F_\theta$  is discrete, then:

$$L(\theta; x) = \mathbb{P}[X_1 = x_1, \dots, X_n = x_n]$$

It is simply the product of the PMFs. It thus represents the probability of observing the sample that we actually observed. Moreover, if for some  $\theta_1, \theta_2 \in \Theta$ , if  $L(\theta_1; x) > L(\theta_2; x)$ , then it means that we were more likely to observe that data if the parameter is  $\theta_1$  instead of  $\theta_2$ . Thus, *information is encoded in the function*.

2. If the distribution  $F_\theta$  is continuous, then this interpretation will obviously not work. However, we can do the usual “trick”:

$$\mathbb{P}_\theta[X \in ]x - \varepsilon, x + \varepsilon[ ] = \int_{x-\varepsilon}^{x+\varepsilon} f(t; \theta) dt \approx f(x; \theta) \cdot 2\varepsilon$$

It is thus proportional up to a constant depending on  $\varepsilon$ . Mathematically:

$$L(\theta; x) \propto \mathbb{P}[X_1 \in ]x_1 - \varepsilon, x_1 + \varepsilon[, \dots, X_n \in ]x_n - \varepsilon, x_n + \varepsilon[ ]$$

Similar to the discrete case, comparing  $L(\theta_1; x)$  with  $L(\theta_2; x)$  will give us a comparison of the probability of the actual observed sample.

So the aim is to find a way to maximise the likelihood function in  $\theta$ .

**Definition 12** (Maximum Likelihood Estimate). For an observed sample  $x = (x_1, \dots, x_n)$ , let  $\hat{\theta}(x)$  be the value for which  $L$  is maximised. In other words:

$$L(\hat{\theta}(x); x) := \sup_{\theta \in \Theta} L(\theta; x) \quad (28)$$

Then,  $\hat{\theta}$  is called the **maximum likelihood estimate**. (Grounding: it is a number, not an object). If  $\hat{\theta}(x)$  exists for almost all samples, and as a function:

$$\hat{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}^k, x \mapsto \hat{\theta}(x)$$

is a measurable, then the function  $\hat{\theta}(X)$  is called the **Maximum Likelihood Estimator (MLE)** of  $\theta$ .

**Example 15** (Continued from the Poisson Case). Let  $x_1, \dots, x_n$  be a sample from the Poisson distribution. Then, the aim is to maximise:

$$L(\theta; x) = \prod_{i=1}^n e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} = e^{-n\lambda} \frac{\lambda^{\sum_{i=1}^n x_i}}{x_1! \cdots x_n!}$$

You *could* differentiate this monstrosity, but differentiating products is in general not fun. However, sums are not too bad to differentiate. So, observe that this function is:

1. Positive
2. Wherever  $\log(L)$  is maximised,  $L$  is maximised.

So we can equivalently maximise  $\log(L(\lambda; x))$ . Thus:

$$\begin{aligned} \log L(\lambda; x) &= \ell(\lambda; x) = -n\lambda + \sum_{i=1}^n x_i \log(\lambda) - \log(x_1! \cdots x_n!) \\ \Rightarrow \frac{\partial \ell(\lambda; x)}{\partial \lambda} &= -n + \left( \sum_{i=1}^n x_i \right) \frac{1}{\lambda} = 0 \iff n\lambda = \sum_{i=1}^n x_i \\ \Rightarrow \hat{\lambda}(x) &= \bar{x} \end{aligned}$$

Take the second partial derivative to verify that this is indeed a maximum.

$$\frac{\partial^2 \ell}{\partial \lambda^2} = -\frac{1}{\lambda^2} \sum_{i=1}^n x_i < 0$$

This, the MLE is  $\bar{x}$ .

Some observations for calculating maximum likelihood estimators:

1. It is often simpler to maximise the **log likelihood function**, which is defined as:

$$\ell(\theta; x) := \log L(\theta; x) \quad (29)$$

2. If  $\ell$  is differentiable, we can look for the maximum by solving the so-called **likelihood equations**:

$$\frac{\partial \ell}{\partial \theta_j} = 0 \quad j = 1, \dots, k \quad (30)$$

**Example 16.** Let  $x_1, \dots, x_n$  be a random sample from  $B(N, p)$ , where  $N$  is *known*. Then,  $\Theta = [0, 1]$  and

$$L(p; x) = \prod_{i=1}^n \binom{N}{x_i} p^{x_i} (1-p)^{N-x_i}$$

thus, the log likelihood function  $\ell$  is:

$$\ell(p; x) = \sum_{i=1}^n \left[ x_i \log(p) + (N - x_i) \log(1-p) + \underbrace{\log \binom{N}{x_i}}_{\text{if } N \text{ were unknown, this would be a big mess}} \right]$$

We are searching for critical points. Let's first assume that  $\bar{x} \notin \{0, N\}$ . We obtain:

$$\frac{\partial \ell}{\partial p} = \left( \sum_{i=1}^n x_i \right) \frac{1}{p} + \left( nN - \sum_{i=1}^n x_i \right) \left( \frac{-1}{1-p} \right)$$

setting the partial derivative equal to zero and solving, we obtain:

$$\begin{aligned}\frac{n\bar{x}}{p} - \frac{n(N - \bar{x})}{1 - p} &= 0 \\ \Rightarrow \frac{\bar{x}(1 - p) - (N - \bar{x})}{p(1 - p)} &= 0 \\ \Rightarrow \bar{x}(1 - p) - (N - \bar{x})p &= 0 \\ \Rightarrow \bar{x} - \bar{x}p - Np + \bar{x}p &= 0 \\ \Rightarrow p &= \frac{\bar{x}}{N}\end{aligned}$$

Thus,  $p = \bar{x}/N$  is the method of moments estimator. We now need to verify that this is indeed a maximum with the second derivative test:

$$\frac{\partial^2 \ell}{\partial^2 p} = -n\bar{x} \frac{1}{p^2} + \underbrace{(nN - n\bar{x})}_{\geq 0} \frac{-1}{(1 - p)^2} < 0 \text{ (concave in } p\text{)}$$

We now need to check the boundary cases. If  $\bar{x} = 0$ , then

$$\ell(p; x) = nN \log(1 - p)$$

This function is decreasing in  $p$ , which implies that it is maximised at  $p = 0 = \frac{\bar{x}}{N}$ . If  $\bar{x} = N$ , then  $x_i = N \forall i \in \{1, \dots, N\}$ . This is an increasing function in  $p$ , which implies that it is maximised at  $p = 1$ , which implies that  $\hat{p} = 1 = \bar{x}/N$ .

For the past two examples, we've seen that the MLE estimator agrees with the method of moments examples. The following example will illustrate that this is not always the case.

**Example 17.** Let  $x_1, \dots, x_n$  be an observation of a random sample taken from the model:

$$\mathcal{F} = \{\mathcal{U}[0, \theta[ \mid \theta \geq 0\}$$

If  $X \sim \mathcal{U}[0, \theta[$ , then using the method of moments we have that:

$$\begin{aligned}\mathbb{E}[X] &= \theta/2 \\ \frac{\theta}{2} &= \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow \hat{\theta} = 2\bar{x}\end{aligned}$$

In contrast, using the method of maximum likelihood gives:

$$L(\theta; x) = \prod_{i=1}^n \frac{1}{\theta} \chi_{x_i \in [0, \theta]}$$

since  $\chi$  depends on  $\theta$ , taking the log will not be a good idea. Instead, this becomes:

$$= \left(\frac{1}{\theta}\right)^n (\chi_{\min_{1 \leq i \leq n}(x_i) \geq 0})(\chi_{\max_{1 \leq i \leq n}(x_i) \leq \theta})$$

For almost all samples, we have:

$$\min_{1 \leq i \leq n} x_i \geq 0$$

but for  $\theta < \max(x_i)$  we have that  $L = 0$  but for  $\theta \geq \max x_i$ ,  $L$  is decreasing. Thus, the MLE is maximised at  $\max_{1 \leq i \leq n} x_i$ . Thus,

$$\hat{\theta}(x) = \max_{1 \leq i \leq n} X_i = X_{(n)}$$

is the MLE (to do: insert graph).

**Example 18.** Let  $x_1, \dots, x_n$  be a random sample from  $N(\mu, \sigma^2)$  from the model

$$\mathcal{F} = \{N(\mu, \sigma^2) \mid \mu \in \mathbb{R}, \sigma^2 \in ]0, \infty[ \}$$

Then:

$$L(\mu, \sigma^2; x) = \prod_{i=1}^n \frac{1}{\sqrt{2\sigma}} \frac{1}{\sqrt{\sigma^2}} \exp \left\{ -\frac{(x_i - \mu)^2}{2\sigma} \right\}$$

and so the log likelihood is:

$$\ell(\mu, \sigma^2; x) = \sum_{i=1}^n \left[ \log \frac{1}{\sqrt{2\pi}} - \frac{1}{2} \log \sigma^2 - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

taking the partials, we obtain:

$$\begin{aligned} \frac{\partial \ell}{\partial \mu} &= - \sum_{i=1}^n \frac{2(x_i - \mu)(-1)}{2\sigma^2} = \frac{n\bar{x} - n\mu}{\sigma^2} \\ \frac{\partial \ell}{\partial \sigma^2} &= \sum_{i=1}^n \left[ -\frac{1}{2} \frac{1}{\sigma^2} + \frac{(x_i - \mu)^2}{2} \frac{1}{\sigma^4} \right] \end{aligned}$$

setting  $\frac{\partial \ell}{\partial \mu} = 0$ , we obtain that  $\hat{\mu} = \bar{x}$ . For the other one:

$$\begin{aligned} -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \bar{x})^2 &= 0 \\ n\sigma^2 &= \sum_{i=1}^n (x_i - \bar{x})^2 \\ \Rightarrow \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \end{aligned}$$

**How do we prove that our values  $\hat{\mu}$  and  $\hat{\sigma}^2$  are indeed maximums?** The trick is to use something called the profile likelihood. Define:

$$\ell^*(\sigma^2) := \sup_{\mu \in \mathbb{R}} \ell(\mu, \sigma^2) \quad (31)$$

we observe that if we maximise  $\ell^*$ , then this amounts to maximising the whole thing, since  $\bar{x}$  minimises the sum  $\sum_{i=1}^n (x_i - \mu)^2$ . So all we need to do is maximise this:

$$\ell^*(\sigma^2) = \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}} - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2$$

the above is called the profile likelihood univariate problem. Solving the above gives you the  $\hat{\sigma}^2$  that is the solution.

### 3.1.3 Invariance of the MLE Theorem

**Motivation:** suppose that  $X_1, \dots, X_n$  are iid taken from Bernoulli( $p$ ). In this case, we have that the method of moments and MLE method agree:

$$\hat{p} = \bar{x} = \hat{p}_{MLE}$$

The standard deviation of  $X_1 = \sqrt{p(1-p)}$ . We can estimate this by  $\sqrt{\hat{p}(1-\hat{p})}$ . This is a non-bijective function of  $p$ . In this sense, this is the MLE of the standard deviation, which motivates the idea of an invariance principle.

**Theorem 19** (Invariance of the MLE). Consider a statistical model  $\{F_\theta \mid \theta \in \Theta\}$  and suppose that  $g : \Theta \rightarrow \mathbb{R}^m$  is an *arbitrary* function. Set  $\Gamma := g(\Theta)$  and  $\gamma := g(\theta)$ . Then, if  $\hat{\theta}(x_1, \dots, x_n)$  is the maximum likelihood estimate of  $\theta$ , then  $g(\hat{\theta}(x_1, \dots, x_n))$  is the ML estimate of  $\gamma$  in the following sense: if

$$L^*(\gamma; x_1, \dots, x_n) = \sup_{\theta \in \Theta \mid g(\theta) = \gamma} \{L(\theta; x_1, \dots, x_n)\} \quad (32)$$

then:

$$L^*(g(\hat{\theta}(x_1, \dots, x_n)); x_1, \dots, x_n) = \sup_{\gamma \in \Gamma} L^*(\gamma; x_1, \dots, x_n) \quad (33)$$

*Proof.*

$$\begin{aligned} L^*(g(\hat{\theta}(x_1, \dots, x_n)); x_1, \dots, x_n) &= \sup_{\theta \mid g(\theta) = g(\hat{\theta})} \{L(\theta; x_1, \dots, x_n)\} \\ &= L(\hat{\theta}; x_1, \dots, x_n) \text{ (since } \hat{\theta} \text{ is the maximiser)} \\ &= \sup_{\theta \in \Theta} L(\theta; x_1, \dots, x_n) \end{aligned}$$

Sup over the range of  $g$ , and then through the pre-image (just running through the  $\theta$ 's in a systematic way):

$$= \sup_{\gamma \in \Gamma} \underbrace{\sup_{\theta \in \Theta \mid g(\theta) = \gamma} L(\theta; x_1, \dots, x_n)}_{= L^*(\gamma; x_1, \dots, x_n)}$$

□

**Example 19.** Let  $X_1, \dots, X_n$  be iid from Bernoulli( $p$ ). If  $\hat{p} = \bar{X} = p_{\text{MLE}}$ , then  $\sqrt{\hat{p}(1 - \hat{p})}$  is the MLE of the standard deviation  $\sqrt{p(1 - p)}$ .

The next example will be a nice philosophical example.

**Example 20.** Among 20 tosses, assume that there were 7 heads. Suppose that we want to estimate  $p$  in two ways:

- Experiment # 1: the coin was tossed 20 times and it had 7 heads. However, we don't know *when* those 7 heads happened. In this case, think of likelihood as the probability of observing what you have observed:

$$\binom{20}{7} p^7 (1 - p)^{13}$$

implies that the log likelihood is

$$\log \binom{20}{7} + \underbrace{7 \log(p) + 13 \log(1 - p)}_{\text{this is called the **kernel** of the log-likelihood}}$$

- Experiment # 2: suppose you waited until 7 heads were tossed, and you are told it took 20 tosses to get there. Therefore, the difference here is that *you know that the 20th toss was a head*. This probability is modelled by the **negative binomial**:

$$\binom{19}{6} p^7 (1 - p)^{13}$$

and so the log likelihood is:

$$\log \binom{19}{6} + \underbrace{7 \log(p) + 13 \log(1 - p)}_{\text{kernel}}$$

Even though the likelihoods are not the same, the MLE between the two will be the same! We have:

$$L_1 \propto p^7(1-p)^{13} \quad (34)$$

that is,  $L_1(p) = c_1 \cdot p^7(1-p)^{13}$ . Similarly,  $L_2(p) = c_2 \cdot p^7(1-p)^{13} \propto p^7(1-p)^{13}$ . Therefore, the MLE estimates are the same in the two experiments, since  $L_1(p) = cL_2(p)$ . This is called the **likelihood principle**: when two likelihood functions of two experiments are equal up to a multiplicative constant, they contain the same information about the unknown parameter.

### 3.1.4 Bayesian Method

In the **Bayesian philosophy**, we quantify the lack of knowledge of a parameter  $\theta$  with a probability distribution or density  $\pi(\theta)$ . This is called a **prior**. A prior distribution is *your* belief of what the probability is. Once the data has been collected—say,  $x_1, x_2, \dots, x_n$ , we can update the prior to incorporate this information, and this leads to the **posterior density**.

The **posterior density function** is given by:

$$\pi(\theta \mid x_1, \dots, x_n) := \frac{f(\theta, x_1, \dots, x_n)}{f(x_1, \dots, x_n)} = \frac{f(x_1, \dots, x_n; \theta)\pi(\theta)}{\int_{\Theta} f(x_1, \dots, x_n; \theta)\pi(\theta)d\theta} \quad (35)$$

**Example 21.** Consider  $X_1, \dots, X_n$  taken from a Bernoulli( $p$ ) distribution with a prior from the Beta distribution. Recall that the **beta function** is given by:

$$\int_0^1 x^{\alpha-1}(1-x)^{\beta-1}dx$$

and so the **beta distribution**  $B(\alpha, \beta)$  has the density:

$$f(x; \alpha, \beta) = x^{\alpha-1}(1-x)^{(\beta-1)} \cdot \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \chi_{x \in ]0,1[} \quad (36)$$

The set of beta distributions has a lot of different shapes as you vary the parameters, which makes it a nice distribution to choose priors from. For example, the uniform is a special case of a beta distribution. Given the data  $x_1, \dots, x_n$ , we obtain:

$$\begin{aligned} f(x_1, \dots, x_n; p) &= \prod_{i=1}^n p^{x_i}(1-p)^{1-x_i} \\ &= p^{n\bar{x}}(1-p)^{n-n\bar{x}} \end{aligned}$$

We thus obtain the posterior:

$$\begin{aligned} \pi(p|x) &= \frac{p^{n\bar{x}}(1-p)^{n-n\bar{x}}p^{\alpha-1}(1-p)^{\beta-1} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}}{\int_0^1 p^{n\bar{x}}(1-p)^{n-n\bar{x}}p^{\alpha-1}(1-p)^{\beta-1} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} dp} \\ &\propto c(x_1, \dots, x_n, \alpha, \beta)p^{n\bar{x}+\alpha-1}(1-p)^{n-n\bar{x}+\beta-1} \end{aligned}$$

The posterior is a beta density:

$$\pi(p|x_1, \dots, x_n) \sim \text{Beta}(n\bar{x} + \alpha, n - n\bar{x} + \beta)$$

here, the Beta is the so-called **conjugate prior**; that means, the posterior belongs to the same class of densities as the prior. Observe that we no longer have a point estimate of  $p$ ; we have a distribution. To obtain a point estimate, we can take the expected value of the density:

$$\begin{aligned} \hat{p}_B &= \text{expected value of the posterior.} \\ \Rightarrow \hat{p}_B &:= \mathbb{E}[Y], \quad Y \sim \text{Beta}(n\bar{x} + \alpha, n - n\bar{x} + \beta) \end{aligned}$$



Note that the expected value of  $B(\alpha, \beta)$  is  $\frac{\alpha}{\alpha+\beta}$  and:

$$\hat{p}_B = \frac{n\bar{x} + \alpha}{n + \alpha + \beta} = \frac{n}{n + \alpha + \beta} \underbrace{\bar{x}}_{\text{MLE}} + \frac{\alpha + \beta}{n + \alpha + \beta} \underbrace{\frac{\alpha}{\alpha + \beta}}_{\text{expectation of prior}}$$

Notice the weights. The more data that you have, the more the weights favour the data. Thus, a small sample means that the prior has a larger effect.

**Example 22** (Poisson Likelihood Example - Continued). Consider the  $\text{Poisson}(\lambda)$  model with parameter  $\lambda > 0$  and data  $x_1, \dots, x_n$ . In order to go Bayesian, we need to multiply the prior estimate of  $\lambda$ , and we want it to be a conjugate prior. *What do we multiply it by?*  $\pi(\lambda) \sim \text{Gamma}(\alpha, \beta)$  works. We obtain:

$$f(x_1, \dots, x_n; \lambda)\pi(\lambda) = \frac{1}{\Gamma(\alpha)} \frac{\lambda^{n\bar{x}} e^{-n\lambda}}{x_1! \cdots x_n!} \underbrace{\lambda^{\alpha-1} e^{-\lambda\beta} \beta^\alpha}_{\pi(\lambda) \sim \text{Gamma}(\alpha, \beta)} \propto \alpha \lambda^{\alpha+n\bar{x}-1} e^{-\lambda(n+\beta)} \quad (37)$$

The posterior is again a Gamma, so the Gamma prior is a conjugate prior. Thus:

$$\pi(\lambda|x) \sim \text{Gamma}(\alpha + n\bar{x}, n + \beta)$$

which gives:

$$\hat{\lambda}_B = \frac{\alpha + n\bar{x}}{n + \beta} = \frac{n}{n + \beta} \bar{x} + \frac{\beta}{n + \beta} \frac{\alpha}{\beta}$$

where the final term comes from the fact that the expected value of  $\text{Gamma}(\alpha, \beta)$  is  $\alpha/\beta$ .

### 3.2 Method of Evaluating Estimators

**Definition 13** (Unbiased, Consistent, and MSE). Let  $\mathcal{F} = \{F_\theta \mid \theta \in \Theta\}$  be our statistical model, let  $\gamma : \Theta \rightarrow \mathbb{R}^n$ , and also let  $T(x_1, \dots, x_n)$  be an estimator of  $\gamma(\theta)$ . Then:

1. An estimator  $T(X_1, \dots, X_n)$  is called **unbiased** if:

$$\mathbb{E}_\theta [T(X_1, \dots, X_n)] = \gamma(\theta) \quad (38)$$

*On average, does it get the right thing?*

2.  $T$  is called **consistent** if

$$T(X_1, \dots, X_n) \xrightarrow{P} \gamma(\theta) \text{ as } n \rightarrow \infty \quad (39)$$

*Some kind of limiting statement – generally uses WLLN techniques*

3. The **mean squared error (MSE)** of  $T(X_1, \dots, X_n)$  is:

$$\text{MSE}(T) := \mathbb{E}_\theta [\{T(X_1, \dots, X_n) - \gamma(\theta)\}^2] \quad (40)$$

Before going on, need to make some important remarks:

1. If  $X_1$  has an expectation  $\mu$  and variance  $\sigma^2$ , then we already had that the sample mean is an unbiased estimator of  $\mu$  and  $s^2$  is an unbiased estimator of  $\sigma^2$  (Theorem 1.9).
2. Sometimes the consistency as defined here is called **weak consistency** to differentiate it from **strong consistency**, which is the case when  $T(X_1, \dots, X_n) \rightarrow \gamma(\theta)$ . In general, it is hard to prove strong consistency. It generally requires the use of the strong law of large numbers.

- a) *Caution!*: beware of asymptotics! You can have quite stupid examples. For example, consider  $N(\mu, \sigma^2)$ . Suppose that we want to estimate  $\mu$ . We can construct this silly unbiased estimator:

$$T(X_1, \dots, X_n) = X_1$$

And we can construct this silly consistent estimator:

$$T(X_1, \dots, X_n) = \begin{cases} X_1 & \text{if } n \leq 10^6 \\ \bar{X} & \text{if } n > 10^6 \end{cases}$$

- b) (Deriving bias variance decomposition of MSE). Write  $T = T(X_1, \dots, X_n)$ . Then, for any  $\theta$ :

$$\begin{aligned} \text{MSE}(T) &= \mathbb{E}_\theta [T(X_1, \dots, X_n) + \mathbb{E}_\theta [T] - \mathbb{E}_\theta [T] + \gamma(\theta)]^2 \\ &= \underbrace{\mathbb{E}_\theta [T - \mathbb{E}_\theta [T]]^2}_{\text{variance of } T} + \mathbb{E}_\theta [\mathbb{E}_\theta [T] - \gamma(\theta)]^2 \end{aligned}$$

$T$  is the only random part.

$$\begin{aligned} &= \text{var}_\theta [T] + (\mathbb{E}_\theta [T] - \gamma(\theta))^2 + 2(\mathbb{E}_\theta [T] - \gamma(\theta)) \underbrace{\mathbb{E} [T - \mathbb{E}_\theta [T]]}_{=0} \\ &= \text{var}_\theta [T] + [\mathbb{E}_\theta [T] - \gamma(\theta)]^2 \end{aligned}$$

We call  $\mathbb{E}_\theta [T] - \gamma(\theta)$  the **bias** of  $T$ , and so:

$$\text{MSE}_\theta = \text{var}_\theta [T] + (\text{Bias}_\theta(T))^2 \quad (41)$$

**Example 23.** Let  $X_1, \dots, X_n$  be a random sample from  $N(\mu, \sigma^2)$ . Then:

$$\text{MSE}(\bar{X}) = \text{var}_\theta [\bar{X}] = \sigma^2/n$$

this converges to zero as  $n \rightarrow \infty$ . Moreover:

$$\text{MSE}[s^2] = \text{var}[s^2] = \text{var} \left[ \underbrace{\frac{n-1}{\sigma^2} s^2}_{\sim \chi_{n-1}^2} \cdot \frac{\sigma^2}{(n-1)} \right] = \frac{\sigma^4}{(n-1)^2} \underbrace{2(n-1)}_{\text{variance of } \chi_{(n-1)}^2} = \frac{2\sigma^4}{(n-1)}$$

This converges at the same rate as the sample mean. Recall the MLE of  $\sigma^2$ :

$$\hat{\sigma}^2 = \frac{(n-1)}{n} s^2$$

and so the bias is:

$$\begin{aligned} \text{Bias}(\hat{\sigma}^2) &= \mathbb{E} \left[ \frac{n-1}{n} s^2 \right] - \sigma^2 \\ &= \frac{n-1}{n} \sigma^2 - \sigma^2 \\ &= \frac{-\sigma^2}{n} \end{aligned}$$

Thus, the bias is always negative and so it will systematically underestimate  $\sigma^2$ . Thus, the  $\text{MSE}(\hat{\sigma}^2)$  is:

$$\begin{aligned}
 \text{MSE}(\hat{\sigma}^2) &= \mathbb{E}[(\hat{\sigma}^2 - \sigma^2)^2] \\
 &= \text{var}[\hat{\sigma}^2] + [\text{Bias}(\hat{\sigma}^2)]^2 \\
 &= \left(\frac{n-1}{n}\right)^2 \text{var}[s^2] + \frac{\sigma^4}{n^2} \\
 &= \frac{(n-1)^2}{n^2} \frac{2\sigma^4}{(n-1)} + \frac{\sigma^4}{n^2} \\
 &= \frac{\sigma^4(2(n-1) + 1)}{n^4} \\
 &= \frac{\sigma^4}{n^2}(2n-1) \\
 &= \frac{2\sigma^4}{n-1} \cdot \underbrace{\frac{2n^2-3n+1}{2n^2}}_{<1} \\
 &< \frac{2\sigma^4}{n-1}
 \end{aligned}$$

So,  $\hat{\sigma}^2$  is closer, on average, to the true value of  $\theta$  than  $s^2$  is. Thus, the MLE is preferable. This illustrates the idea of a tradeoff between bias and variability.

To do: insert figure

**Example 24.** Let  $X_1, \dots, X_n$  be taken from  $\text{Bernoulli}(p)$ . Recall that  $\hat{p}_{\text{MLE}} = \bar{x}$  and  $\text{MSE}(\hat{p}_{\text{MLE}}) = \frac{p(1-p)}{n}$ . In contrast, a Bayesian estimator with a  $\text{Beta}(\alpha, \beta)$  prior, we have:

$$\hat{p}_B = \frac{n}{n + \alpha + \beta} \bar{x} + \frac{\alpha}{n + \alpha + \beta} = \frac{n\bar{x} + \alpha}{n + \alpha + \beta}$$

Calculating the MSE gives:

$$\begin{aligned}
 \text{MSE}(\hat{p}_B) &= \text{var} \left[ \frac{n\bar{x} + \alpha}{n + \alpha + \beta} \right] + \left( \text{Bias} \left[ \frac{n\bar{x} + \alpha}{n + \alpha + \beta} \right] \right)^2 \\
 &= \underbrace{\frac{np(1-p)}{(n + \alpha + \beta)^2}}_{\text{variance}} + \left( \frac{np + \alpha}{n + \alpha + \beta} - p \right)^2 \\
 &= \frac{np(1-p)}{(n + \alpha + \beta)^2} + \frac{(\alpha - p\alpha - p\beta)^2}{(n + \alpha + \beta)^2}
 \end{aligned}$$

To compare, we need to choose  $\alpha$  and  $\beta$  to be independent of  $p$ .

$$\text{MSE}(\hat{p}_B) = \dots = \frac{\alpha^2 + p(n - 2\alpha^2 - 2\alpha\beta) + p^2(-n + \alpha^2 + \beta^2 + 2\alpha\beta)}{(n + \alpha + \beta)^2}$$

Thus, if we want the above equation to be independent of  $p$ , set the coefficients of  $p$  and  $p^2$  to equal zero. We thus must solve:

$$\begin{aligned}
 n - 2\alpha^2 - 2\alpha\beta &= 0 \\
 -n + \alpha^2 + \beta^2 + 2\alpha\beta &= 0
 \end{aligned}$$

which gives:

$$\alpha = \sqrt{\frac{n}{4}} \qquad \beta = \sqrt{\frac{n}{4}}$$

and therefore:

$$\text{MSE}(\hat{p}_B) = \frac{n}{4(n + \sqrt{n})^2} \quad \hat{p}_B = \frac{n\bar{x} + \sqrt{n/4}}{n + \sqrt{n}}$$

To do: include figures. Thus, for small sample sizes,  $\text{MSE}(\hat{p}_B) \ll \text{MLE}(\hat{p}_{\text{MLE}})$ . The advantage shrinks as sample size increases.

**Theorem 20.** Suppose  $\mathbb{E}_\theta [T(X_1, \dots, X_n)] \rightarrow \gamma(\theta)$  as  $n \rightarrow \infty$  (that is,  $T(X_1, \dots, X_n)$  is asymptotically unbiased) and  $\text{var}_\theta[T] \rightarrow 0$  as  $n \rightarrow \infty$ . Then,  $T$  is a consistent estimator of  $\gamma(\theta)$ .

*Proof.* Application of the Markov Inequality. Fix  $\varepsilon > 0$ . Then:

$$\mathbb{P}[|T(X_1, \dots, X_n) - \gamma(\theta)| > \varepsilon] \leq \frac{\mathbb{E}[(T(X_1, \dots, X_n) - \gamma(\theta))^2]}{\varepsilon^2} = \frac{\text{var}_\theta[T] - (\text{Bias}_\theta[T])^2}{\varepsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty$$

□

**Q:** Does there exist a strategy of choosing the best estimator? **A:** No! Consider the other “stupid” estimator: Suppose we want to estimate  $\theta$  and  $17 \in \Theta$ . Then the statistic:

$$T(X_1, \dots, X_n) = 17$$

will always have

$$\text{MSE}_{17}(T) = 0$$

You can never beat that silly statistic. More generally, you can never *systematically* minimise the MSE uniformly over all values of  $\theta$ . Thus, we need to restrict the class of estimators that we consider, which leads us to the next section.

### 3.3 Best Unbiased Estimators

**Definition 14** (Uniform Minimum Variance Unbiased Estimator). An estimator  $T^*(X_1, \dots, X_n)$  is called a uniform minimum variance unbiased estimator (UMVUE) of  $\gamma(\theta)$  if

1. For all  $\theta \in \Theta$ ,  $\mathbb{E}_\theta [T^*] = \gamma(\theta)$ . It must be unbiased for every  $\theta$ .
2. For every other unbiased estimator  $T$ :

$$\text{var}_\theta[T^*] \leq \text{var}_\theta[T]$$

for all  $\theta \in \Theta$ . It beats it for any  $\theta \Rightarrow$  uniform control.

This is an estimator that systematically reduces the bias.

**Example 25.** Assume that we have a random variable distributed as Poisson. Then,  $\hat{\lambda} = \bar{x}$ . Forget about trying to do a Bayesian estimate on this. However, for the variance we have:

$$\text{Var}[\hat{\lambda}] = \lambda/n$$

This seems like it could be a sensible estimator. **Q:** Is it the UMVUE? This question motivates the following theorem:

**Theorem 21** (Crammer-Rao Inequality). Suppose  $X_1, \dots, X_n$  form a random sample. Let  $F_\theta, \theta \in \Theta, \Theta = ]a, b[$  where  $-\infty \leq a < b \leq \infty$  form the statistical model. Let  $T(X_1, \dots, X_n)$  be an unbiased estimator of  $\gamma(\theta)$ . Assume:

1.  $\forall \theta \in \Theta$ ,  $F_\theta$  has a density/PMF  $f(\cdot; \theta)$ , and  $\partial f(x; \theta)/\partial \theta$  exists for all  $\theta \in \Theta$ .

2. You can differentiate the log:

$$\mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta} \log f(x; \theta) \right] = 0 \quad \forall \theta \in \Theta$$

and

$$\mathbb{E}_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f(x; \theta) \right)^2 \right] = I(\theta) < \infty$$

Remark the  $I(\theta)$  is called the **Fisher information**

3.  $\text{var}_\theta[T] < \infty$  and:

$$\sum_{i=1}^n \mathbb{E}_\theta \left[ T(X_1, \dots, X_n) \frac{\partial}{\partial \theta} \log f(X_i; \theta) \right] = \gamma'(\theta)$$

then:

$$\text{var}_\theta T(X_1, \dots, X_n) \geq \frac{\{\gamma'(\theta)\}^2}{nI(\theta)}$$

*Remark:* this is called the **Crammer-Rao lower bound**

*Proof.* The theorem follows from the Cauchy - Schwartz inequality. Recall from probability:

$$\{\text{cov}(W, Z)\}^2 \leq \text{var}[W] \text{var}[Z]$$

We need to choose our quantities for Cauchy-Schwartz inequality:

1.  $\text{var}[W]$  is:

$$\text{var}[W] = \sum_{i=1}^n \text{var} \left[ \frac{\partial}{\partial \theta} \log f(x_i; \theta) \right] = nI(\theta) < \infty \quad (\text{by assumption 2})$$

2. By assumption 3:

$$\text{var}_\theta T < \infty$$

3. Covariance:

$$\text{cov}_\theta[W, Z] = \mathbb{E}_\theta[W \cdot Z] = \mathbb{E}_\theta \left[ \sum_{i=1}^n T \frac{\partial}{\partial \theta} \log f(x_i; \theta) \right] = \sum_{i=1}^n \mathbb{E}_\theta[T] \cdot \frac{\partial}{\partial \theta} \log f(x_i; \theta) = \gamma'(\theta) \quad (\text{assumption 3})$$

Now it follows from Cauchy Schwartz:

$$[\gamma'(\theta)]^2 \leq nI(\theta) \cdot \text{var}_\theta T$$

which happens  $\iff$  :

$$\text{var}_\theta[T] \geq \frac{\{\gamma'(\theta)\}^2}{nI(\theta)}$$

Which is what we wanted to show. □

Some remarks from the proof:

1. There are a lot of versions of the Crammer-Rao inequality.

- a) This version has the cleanest possible assumptions.
- b) We did not need to differentiate between the continuous and discrete cases.

2. Observe the set  $\{x \mid f(x; \theta) > 0\}$ . We define  $N_\theta$  as the set such that  $\mathbb{P}[x \in \{x \mid f(x; \theta) > 0\}] = 1$ . On here,  $X \sim F_\theta$ . So, you can restrict your attention to this set which makes the log well-defined. Thus, on  $N_\theta$ ,  $\partial/\partial\theta \log f(x, \theta)$  exists.
3. **Q:** What do assumptions (ii) and (iii) mean?
- a) This means nothing more than interchanging differentiation and integration. In particular, if  $\{x \mid f(x; \theta) > 0\}$  does NOT depend on  $\theta$ , then, assuming that the  $X$ 's have a density, then these assumptions amount to interchanging differentiation and integration.
- b) From Assumption 2, we have:

$$\int_{\mathbb{R}} f(x; \theta) dx = 1$$

and so differentiating both sides with respect to  $\theta$  gives

$$\int_{\mathbb{R}} f(x; \theta) dx = 0$$

multiplying and dividing by  $f$  gives:

$$\int_{\mathbb{R}} \left[ \frac{\partial}{\partial \theta} \log f(x; \theta) \right] = 0$$

so, assumption 2.1 is saying that you can interchange the derivatives and integrals.

Also,

$$\gamma(\theta) = \mathbb{E}_\theta [T] = \int \cdots \int T(x_1, \dots, x_n) \sum_{i=1}^n \frac{\partial}{\partial \theta} f(x_i; \theta) \prod_{i=1}^n f(x_i; \theta) dx_1 \cdots dx_n$$

We will play the same game again: differentiating and swapping derivatives and integrals gives us:

$$\gamma'(\theta) = \int \cdots \int T(x_1, \dots, x_n) \sum_{i=1}^n \frac{\partial}{\partial \theta} f(x_i; \theta) \prod_{j \neq i} f(x_j; \theta) dx_1 \cdots dx_n$$

artificially multiplying and dividing by  $f(x_i; \theta)$  gives:

$$\begin{aligned} &= \int \cdots \int T(x_1, \dots, x_n) \sum_{i=1}^n \underbrace{\frac{\partial}{\partial \theta} \frac{f(x_i; \theta)}{f(x_i; \theta)}}_{\frac{\partial}{\partial \theta} \log f(x_i; \theta)} \prod_{j=1}^n \underbrace{f(x_j; \theta)}_{\text{expected value}} \\ &= \sum_{i=1}^n \mathbb{E}_\theta \left[ T(x_1, \dots, x_n) \frac{\partial}{\partial \theta} \log f(x_i; \theta) \right] \end{aligned}$$

These are actually smoothness conditions. There is also a nice geometric interpretation.

### 3.3.1 Geometric Interpretation

Define the following space:

$$\mathcal{L}_2 := \{g : \mathbb{R}^n \rightarrow \mathbb{R}, \text{ borel measurable, } \mathbb{E}_\theta [g^2(x_1, \dots, x_n)] < \infty\}$$

On  $\mathcal{L}_2$  space, we can define a scalar product:  $g_1, g_2 \in \mathcal{L}_2$ , then:

$$\langle g_1, g_2 \rangle = \mathbb{E}_\theta [g_1(x_1, \dots, x_n) \cdot g_2(x_1, \dots, x_n)]$$

On  $\mathcal{L}_2$ , the norm for  $g \in \mathcal{L}_2$

$$\|g\|^2 = \langle g, g \rangle = \mathbb{E}_\theta [g^2] < \infty$$

$\mathcal{L}_2$  is a Hilbert space. In particular, it is a space of equivalence classes of functions. In this case:

$$g_1 \sim g_2 \text{ if } \mathbb{P}_\theta[g_1(x_1, \dots, x_n) = g_2(x_1, \dots, x_n)] = 1$$

We can consider projections since  $\mathcal{L}_2$  is a Hilbert space. Let  $T$  be an unbiased estimator of  $\gamma(\theta)$ . Then,  $T \in \mathcal{L}_2$ . Define  $\bar{T} := T - \gamma(\theta)$ . Then,  $\bar{T} \in \mathcal{L}_2$  and  $\langle \bar{T}, 1 \rangle = 0$  (since it's unbiased), since  $\langle \bar{T}, 1 \rangle = \mathbb{E}_\theta[T - \gamma(\theta)]$ . Thus,  $\bar{T}$  is orthogonal to the subspace of constant functions. Define the following function:

$$k(x_1, \dots, x_n) := \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i; \theta) \quad (42)$$

From assumption (ii) of the Crammer-Rao inequality,  $\|k\|^2 = nI(\theta) \in ]0, \infty[$  which implies that  $k \in \mathcal{L}_2$ . Since the expected value again is zero, we have that  $\langle k, 1 \rangle = 0$ . Moreover, from assumption (iii) of the Crammer-Rao inequality:

$$\langle k, T \rangle = \gamma'(\theta) = \langle k, \bar{T} \rangle$$

Therefore, we can re-frame the assumptions of the Crammer-Rao inequality in terms of assumptions on the norms and scalar products. Now define  $V := \text{span}(k) = \{c \cdot k \mid c \in \mathbb{R}\}$  and project  $\bar{T}$  onto  $V$ . Then:

$$\mathbf{T}_V \bar{T} = \left\langle \bar{T}, \frac{k}{\|k\|} \right\rangle \frac{k}{\|k\|} = \langle \bar{T}, k \rangle \frac{k}{\|k\|^2} = \frac{\gamma'(\theta)k}{nI(\theta)}$$

Now, by the pythagorean theorem:

$$\underbrace{\|\bar{T}\|^2}_{\text{variance}} = \|\mathcal{T}_V \bar{T}\|^2 + \|\bar{T} - \mathcal{T}_V \bar{T}\|^2$$

However, this is precisely the inequality:

$$\begin{aligned} \text{var}_\theta[T] &= \frac{\{\gamma'(\theta)\}^2}{nI(\theta)} + \underbrace{\|T - \mathcal{T}_V \bar{T}\|^2}_{\geq 0} \\ &\geq \frac{\{\gamma'(\theta)\}^2}{nI(\theta)} \end{aligned}$$

equality holds *only* when  $\bar{T} \in V$ . So, the UMVUE must be in  $V$ . That is:

$$T - \gamma(\theta) = ck = a(\theta)k$$

where  $a(\theta)$  is a constant that depends on  $\theta$ . So, the UMVUE is of the form:

$$T = a(\theta)k(x_1, \dots, x_n) + \gamma(\theta) \quad (43)$$

We will see later that the only estimators of this form are members of the exponential family.

**Example 26.** Consider the Poisson family.

$$\{\text{Poi}(\lambda), \lambda > 0\}$$

then, the densities are given by:

$$f(x; \lambda) = e^{-\lambda} \frac{\lambda^x}{x!} \chi_{x \in N_\lambda}$$

Again, we don't need to worry about the bad points. In this case,  $N_\lambda = \mathbb{N}_0$ , and so we can differentiate the density with respect to  $\lambda$  and so we have no problems. We will now check the conditions of the Crammer-Rao inequality:

$$\frac{\partial}{\partial \lambda} \log f(x, \lambda) = \frac{\partial}{\partial \lambda} [-\lambda + x \log(\lambda) - \log x!] = -1 + \frac{x}{\lambda}$$

and so the first condition is met. Now for the second assumption:

$$\mathbb{E}_\lambda \left[ \frac{\partial}{\partial \lambda} \log f(x; \lambda) \right] \underbrace{=}_{x \sim \text{Pos}(\lambda)} \mathbb{E}_\lambda \left[ -1 + \frac{x}{\lambda} \right]$$

and for the variance

$$\text{var} \left[ \frac{\partial}{\partial \lambda} \log f(x; \lambda) \right] = \mathbb{E}_\lambda \left[ \left( \frac{\partial}{\partial \lambda} \log f(x; \lambda) \right)^2 \right] = \text{var}_\lambda \left[ -1 + \frac{x}{\lambda} \right] = \frac{1}{\lambda^2} \text{var}_\lambda[x] = \frac{1}{\lambda} = I(\lambda)$$

so, by the Crammer-Rao inequality, the Crammer-Rao bound is:

$$\frac{1}{I(\lambda) \cdot n} = \frac{\lambda}{n}$$

So, if  $T^*(x_1, \dots, x_n) = \bar{x}$ , then  $\mathbb{E}[T^*] = \lambda$  and  $\text{var}[T^*] = \lambda/n$ . So, the sample mean attains the Crammer-Rao upper bound. Is  $\bar{x}$  the UMVUE, or is there something else to check? We need to check (iii). If  $T$  is unbiased for  $\lambda$ , and  $\text{var}_\lambda[T] < \infty$ , can we have the validity of (iii)?

$$1 = \sum_{i=1}^n \mathbb{E}_\theta[T] \frac{\partial}{\partial \lambda} \log f(x_i; \lambda)$$

This comes down to interchanging summation and differentiation. We know that  $T$  is unbiased:

$$\lambda = \mathbb{E}_\theta[T] = \sum_{x_1=0}^{\infty} \cdots \sum_{x_n=0}^{\infty} T(x_1, \dots, x_n) e^{-\lambda} \frac{\lambda^{\sum x_i}}{\prod_{i=1}^n x_i!}$$

Trick: move the  $e^{-n\lambda}$  out:

$$= e^{-n\lambda} \sum_{k=1}^{\infty} \underbrace{\left( \sum_{(x_1, \dots, x_n) \in \mathbb{N}_0^n \mid \sum_{i=1}^n x_i = k} T(x_1, \dots, x_n) \frac{1}{\prod_{i=1}^n x_i!} \right)}_{:=a_k} \lambda^k$$

This is an absolutely convergent power series. The variance is finite. In the radius of convergence, we can interchange summation and differentiation:

$$1 = \sum_{k=0}^{\infty} \frac{\partial}{\partial \lambda} \left[ e^{-n\lambda} a_k \lambda^k \right]$$

You can deduce that (iii) holds. So, Crammer-Rao is applicable and

$$\text{var}_\theta[T] \geq \frac{\lambda}{n} = \text{var}_\lambda T^*$$

and so the sample mean is the UMVUE:  $T^* = \bar{x}$

**Example 27.** Consider the so-called **one-parameter exponential family**. This is an umbrella term for the following families of distributions with distribution functions of the form:

$$f(x; \lambda) = \chi_{x \in A} \cdot \exp \{c(\theta) \cdot T(x) + d(\theta) + s(x)\}$$



important: the support of the characteristic function does NOT depend on  $\theta$ ! For example, the PMF of the Poisson( $\lambda$ ) is:

$$f(x; \lambda) = \chi_{x \in \mathbb{N}_0} \exp\{-\lambda + x \log(\lambda) - \log(x!)\}$$

here,

$$\begin{aligned} c(\lambda) &= \log(\lambda) \\ d(\lambda) &= -\lambda \\ T(x) &= x \\ s(x) &= -\log(x!) \end{aligned}$$

Or, the PDF of  $N(\mu, \sigma^2)$  with KNOWN  $\sigma^2$  is:

$$f(x, \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \cdot \left\{ \frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} - \frac{x^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right\}$$

can be written as:

$$= \chi_{x \in \mathbb{R}} \exp \cdot \left\{ \frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} - \frac{x^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right\}$$

Which means that it is in the one parameter exponential from the set:

$$\begin{aligned} c(\mu) &= \frac{\mu}{\sigma^2} \\ d(\mu) &= \frac{-\mu^2}{2\sigma^2} \\ T(x) &= x \\ S(x) &= \frac{-x^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \end{aligned}$$

However, if  $\Theta = ]a, b[$ ,  $-\infty \leq a < b \leq \infty$ , and  $c$  is continuously differentiable with derivative  $c'(\theta) > 0$  for all  $\theta \in \Theta$  (this is the regularity condition), then for any unbiased estimator  $T^*$  for which the assumptions of the Crammer-Rao inequality are met (i.e., we can swap differentiation and integration), we have that the only densities for which the Crammer-Rao inequality is attained is if the density belongs to an exponential family.

If  $\mathbb{E}[T(x)] = \theta$ , then:

$$I(\theta) = \frac{1}{\text{var}[T(x)]}$$

and

$$\frac{1}{n} \sum_{i=1}^n T(x_i)$$

is the UMVUE.

**Example 28.**  $c(\theta)T(x) + d(\theta) + s(x)(:= (*))$  Then:

$$f(x; \theta) = \chi_{x \in A} e^{(*)}$$

Then, the one-parameter exponential family is:

$$\theta \in \Theta, \text{ where } \Theta \in ]a, b[$$

These are the kind of distributions for which the theory works nicely. Moreover:

1. If  $c$  is continuously differentiable on  $\Theta$ , and  $c' > 0$ , and  $w(x_1, \dots, x_n)$  is an estimator with finite variance and unbiased, then the assumptions of the Crammer-Rao theorem holds.
2. If  $\mathbb{E}_\theta [T(x)] = \theta$ , then  $I(\theta) = \frac{1}{\text{var}_\theta T(x)}$  and  $\frac{1}{n} \sum_{i=1}^n T(x_i)$  is the UMVUE of  $\theta$ .

Interestingly, if  $n = 1$ , we know that the estimator that attains the Crammer-Rao lower bound must have the form:

$$T(x) = a(\theta) \frac{\partial}{\partial \theta} \log f(x; \theta) + \theta$$

for  $x$  fixed,  $\frac{\partial}{\partial \theta} \log f(x; \theta) = \frac{T(x) - \theta}{a(\theta)} = c'(\theta)T(x) + d'(\theta)$ , which implies that  $\log f(x; \theta) = c(\theta)T(x) + d(\theta) + s(x)$ .

Other examples of one-parameter exponential families with  $T(x) = x$  are  $\text{Poisson}(\lambda)$  or  $N(\mu, \sigma^2)$  with *known*  $\sigma^2$ . In this case,  $\bar{x}$  is the UMVUE for  $\lambda, \mu$  respectively.

**Lemma 22.** If  $f(x; \theta)$  is a PDF/PMF so that  $f$  is differentiable with respect to  $\theta \forall \theta, x \in N_\theta$ , and so that for  $X$  with a PDF/PMF  $f(x; \theta)$ :

1.  $\mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta} \log f(x; \theta) \right] = 0$  and
2.  $\frac{\partial}{\partial \theta} \mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta} \log f(x; \theta) \right] = 0 =$

$$\left\{ \begin{array}{l} \int \frac{\partial}{\partial \theta} \left( \left[ \frac{\partial}{\partial \theta} \log f(x; \theta) \right] f(x; \theta) \right) dx \\ \sum_x \frac{\partial}{\partial \theta} \left[ \frac{\partial}{\partial \theta} \log f(x; \theta) \right] \end{array} \right.$$

depending on if  $x$  is discrete or continuous, then:

$$I(\theta) = \mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta} \log f(x; \theta) \right]^2 = -\mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta^2} \log f(x; \theta) \right]$$

*Proof.* We will only prove this in the case that  $X$  is continuous. All this is doing is carrying out differentiation. The LHS is zero. The RHS is:

$$\begin{aligned} & \int \frac{\partial}{\partial \theta} \left( \frac{\partial / \partial \theta f(x; \theta)}{f(x; \theta)} \right) dx \\ &= \int \left[ \frac{\partial^2}{\partial \theta^2} \log f(x; \theta) \right] \cdot f(x; \theta) dx + \int \left( \frac{\partial}{\partial \theta} \log f(x; \theta) \right) \frac{\partial}{\partial \theta} f(x; \theta) dx \end{aligned}$$

Observe that the  $f(x; \theta)$  term is simply:

$$\begin{aligned} &= \frac{\partial / \partial \theta f(x; \theta)}{f(x; \theta)} \cdot f(x; \theta) \\ &= \frac{\partial}{\partial \theta} \log f(x; \theta) f(x; \theta) \end{aligned}$$

which implies that:

$$0 = \mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta^2} \log f(x; \theta) \right] + \mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta} \log f(x; \theta) \right]^2$$

which proves what we wanted to show.  $\square$

The next example highlights the limitations of the result:

**Example 29.** Let's consider  $N(\mu, \sigma^2)$  with  $\mu$  known. Suppose that we want to estimate  $\sigma^2$ . Then:

$$\log f(x; \mu, \sigma^2) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma^2) - \frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}$$

This is a one-parameter exponential family. The assumptions of CR hold here (can be checked). To calculate the CR lower bound, we need to calculate  $I(\sigma^2)$ :

$$\frac{\partial}{\partial^2 \sigma^2} \log f(\mu, \sigma^2) = -\frac{1}{2} \frac{1}{\sigma^2} - \frac{1}{2} \frac{(x - \mu)^2}{\sigma^4}$$

and

$$\mathbb{E} \left[ -\frac{1}{2} \frac{1}{\sigma^2} - \frac{1}{2} \frac{(x - \mu)^2}{\sigma^4} \right] = -\frac{1}{2} \frac{1}{\sigma^2} - \frac{1}{2} \frac{1}{\sigma^4} \sigma^2 = 0$$

By the Lemma:

$$\frac{\partial}{\partial^2 \sigma^2} \log f(x; \mu, \sigma^2) = \frac{1}{2} \frac{1}{\sigma^4} + \frac{1}{2} (x - \mu)^2 (-2) \frac{1}{\sigma^6}$$

and so

$$\begin{aligned} \mathbb{E} \left[ \frac{\partial}{\partial^2 \sigma^2} \log f(x; \mu, \sigma^2) \right] &= \frac{1}{2} \frac{1}{\sigma^4} - \frac{1}{\sigma^6} \mathbb{E} [x - \mu]^2 \\ &= \frac{1}{2\sigma^4} - \frac{1}{\sigma^4} \\ &= -\frac{1}{2} \frac{1}{\sigma^4} \end{aligned}$$

So:

$$I(\sigma^2) = \frac{1}{2} \frac{1}{\sigma^4}$$

So, the Crammer-Rao lower bound for  $\sigma^2$  is:

$$\frac{2\sigma^4}{n}$$

**Q:** Can we find an unbiased estimator with this variance? For the sample variance:

$$\mathbb{E} [s^2] = \sigma^2$$

but

$$\text{var}[s^2] = \frac{2\sigma^4}{(n-1)}$$

where the variance was computed in Example 2.16. So, we have

$$\text{var}[s^2] > \frac{2\sigma^2}{n}$$

∴. AN UMVUE would be:

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

but this requires that we KNOW  $\mu$ . Note that

$$\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2} \sim \chi_n^2$$

But, if  $\mu$  is unknown, then the above is not an estimator, since there is dependence on a parameter, and thus the CR lower bound cannot be attained. This shows a limitation of the CR method.

**Example 30.** Observe that the CR lower bound is:

$$\frac{1}{nI(\theta)}$$

if we want to estimate  $\theta$ , and if CR is applicable, then:

$$\text{Var}_\theta T \geq \frac{1}{nI(\theta)}$$

which means that the maximum speed at which you could get  $\text{var}[T]$  to go to zero is  $1/n$ . This means it's at best  $O(1/n)$ . Now consider the uniform model  $U[0, \theta]$ ,  $\theta > 0$ . Then:

$$f(x; \theta) = \frac{1}{\theta} \chi_{x \in ]0, \theta[}$$

We already have:

$$\begin{aligned} \hat{\theta}_{\text{MOM}} &= 2\bar{x} \\ \hat{\theta}_{\text{MLE}} &= \max\{x_1, \dots, x_n\} \end{aligned}$$

Here,  $\hat{\theta}_{\text{MOM}}$  is unbiased and  $\text{var}[\hat{\theta}_{\text{MOM}}] = \frac{4\text{var}[x]}{n}$ . The rate,  $4/n$ , is what's interesting. If  $X_{(n)} = \max\{X_1, \dots, X_n\}$ , then:

$$f_{X_{(n)}}(x; \theta) = \frac{n}{\theta} x^{n-1}$$

and so

$$\mathbb{E}[X_{(n)}] = \int_0^\theta \frac{x^n}{\theta^n} n dx = \theta \frac{n}{n+1}$$

which means that the MLE is not unbiased. So, to make it unbiased, set:

$$\tilde{\theta} := \hat{\theta}_{\text{MLE}} \frac{(n+1)}{n}$$

which means that

$$\text{var}[\theta] = \left[ \frac{(n+1)}{n} \right]^2 \int_0^\theta \left[ x - \frac{\theta n}{n+1} \right]^2 \frac{nx^{n-1}}{\theta^n} dx = \dots = \frac{\theta^2}{n(n+2)} \ll \frac{1}{nI(\theta)} = \frac{\theta^2}{n}$$

Something is wrong. Let's check the conditions of the CR. Set  $N_\theta := ]0, \theta[$ . Then,  $\partial/\partial\theta \log f(x; \theta)$  on  $]0, \theta[$  is:

$$\begin{aligned} \frac{\partial}{\partial\theta} \left( \frac{1}{\theta} \right) &= -\frac{1}{\theta^2} \\ \frac{\partial}{\partial\theta} \log f(x; \theta) &= -\frac{1}{\theta} \\ I(\theta) &= \mathbb{E} \left[ \left( \frac{\partial}{\partial\theta} \log f(x; \theta) \right)^2 \right] = \mathbb{E} \left[ \frac{1}{\theta^2} \right] = \frac{1}{\theta^2} \end{aligned}$$

as we can see, the assumptions of CR are not fulfilled on SEVERAL accounts. Also:

$$\mathbb{E} \left[ \frac{\partial}{\partial\theta} \log f(x; \theta) \right] = -\frac{1}{\theta} \neq 0$$

we have

$$\begin{aligned} \frac{\partial}{\partial\theta} \int_0^\theta t(x) \underbrace{f(x; \theta)}_{1/\theta} dx \\ = \frac{t(\theta)}{\theta} + \int_0^\theta t(x) \frac{\partial}{\partial\theta} \left( \frac{1}{\theta} \right) dx \neq \int_0^\theta t(x) \frac{\partial}{\partial\theta} \left( \frac{1}{\theta} \right) dx \end{aligned}$$

unless  $t(\theta) = 0 \forall \theta$ .

So, we need to temporarily suspend our UMVUE search and do more theory.

*Cutoff for Midterm Material*

## 4 Sufficiency and Completeness

Let  $X_1, \dots, X_n$  be taken from  $F_\theta$ ,  $\theta \in \Theta \subseteq \mathbb{R}^k$ . **Q:** How can we summarise the sample without losing information about  $\theta$ ? This question is closely related to UMVUEs, since statistics are summaries of data. They will be optimal if we do not lose a lot of information about data when making statistics. For example, imagine the Bernoulli experiment:

$$\{\text{Bernoulli Distributions}, \text{Bernoulli}(p), p \in [0, 1]\}$$

If  $X_1, \dots, X_n$  is a random sample from  $\text{Bernoulli}(p)$ , then let  $\mathbb{X}$  denote the tuple  $(X_1, \dots, X_n)$ . Then, the sample space consists of sequences of length  $n$  consisting of 0's and 1s. It is denoted by:

$$\mathcal{X} = \{0, 1\}^n$$

To make things more concrete, we can recall the U.S. election example. We can think about it as follows:

$$\mathbb{P}[\mathbb{X} = (x_1, \dots, x_n)] = p^{\sum_{i=1}^n x_i} (1-p)^{1-\sum_{i=1}^n x_i}$$

Here we have a function of  $p$  (the likelihood,  $L(p)$ ), and it only depends on the  $x_i$ 's through their sum. In other words, the probability depends on  $(x_1, \dots, x_n)$  only through the sum of the  $x_i$ 's (or equivalently, through their mean). *Perhaps the sum  $\sum_{i=1}^n x_i$  contains all the relevant information?* To check this guess, we will calculate the conditional distribution:

$$\begin{aligned} \mathbb{P}\left[X_1 = x_1, \dots, X_n = x_n \mid \sum_{i=1}^n x_i = t\right] &= \begin{cases} 0 & \text{if } \sum_{i=1}^n x_i \neq t \\ 1 & \text{if } \sum_{i=1}^n x_i = t \end{cases} \\ &= \frac{\mathbb{P}[X_1 = x_1, \dots, X_n = x_n, \sum_{i=1}^n x_i = t]}{\mathbb{P}[\sum_{i=1}^n x_i = t]} \\ &= \frac{p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}}{\binom{n}{\sum_{i=1}^n x_i} p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}} \\ &= \frac{1}{\binom{n}{\sum_{i=1}^n x_i}} \end{aligned}$$

Which does not depend on  $p$ ! Now consider a statistic  $T(X_1, \dots, X_n)$ . Denote by  $\mathcal{X}$  the sample space of the  $(X_1, \dots, X_n)$  and by  $\mathcal{T}$  the range of the statistic  $T(X_1, \dots, X_n)$ :

$$\mathcal{T} := \{t \in \mathbb{R}^k \mid \exists (x_1, \dots, x_n) \in \mathcal{X} \text{ so that } T(x_1, \dots, x_n) = t\}$$

We want to look at a statistic as something that partitions the sample space  $\mathcal{X}$ .  $\forall t \in \mathcal{T}$ , define:

$$A_t := \{(x_1, \dots, x_n) \in \mathcal{X} \mid T(x_1, \dots, x_n) = t\}$$

Concretely, for the Bernoulli experiment, we have that  $\mathcal{T} = \{0, \dots, n\}$  and the partitions are:

$$A_t := \left\{ (x_1, \dots, x_n) \in \{0, 1\}^n \mid \sum_{i=1}^n x_i = t \right\}$$

Now,  $\{A_t \mid t \in \mathcal{T}\}$  is a partition of  $\mathcal{X}$ . The  $A_t$ 's are disjoint and  $\mathcal{X} = \bigcup_{t \in \mathcal{T}} A_t$ .

Insert graphic here

**Definition 15** (Sufficient). Let  $X_1, \dots, X_n$  be a random sample from  $F_\theta$ ,  $\theta \in \Theta$ . A statistic  $T = T(X_1, \dots, X_n)$  is called **sufficient** for  $\theta$  if the conditional distribution of  $(X_1, \dots, X_n)$  given  $T$  doesn't depend on  $\theta$ , and this holds for all  $t \in \mathcal{T}$  so that  $\mathbb{P}_\theta[t \in \mathcal{T}] = 1 \forall \theta \in \Theta$ .

**Remark:** Sufficiency means that the distribution of the random vector, given  $(x_1, \dots, x_n) \in A_t$ , does not depend on  $\theta$ .

**Example 31** (Bernoulli). Let  $X_1, \dots, X_n \sim \text{bernoulli}(p)$ . Then we have that  $\sum_{i=1}^n x_i$  is sufficient for  $p$ . So is  $\bar{x}$ .

In general, a random vector is a statistic. Thus,  $(X_1, \dots, X_N)$  is always sufficient for  $\theta$ . We also have that the order statistics  $(X_{(1)}, \dots, X_{(n)})$  is sufficient for  $\theta$  (this will be proven in Assignment 3).

## 4.1 Proving the Sufficiency of Statistics

**Motivation:** Suppose a statistic  $T$  is sufficient for  $\theta$ . Assume that there are two experimenters. Experimenter # 1 has access to the full sample  $(x_1, \dots, x_n)$ . Experimenter # 2 only has access to the statistic  $T(x_1, \dots, x_n)$ . Because  $(X_1, \dots, X_n) \mid T = t$  does not depend on  $\theta$ , let's call it  $P_t^*$ . So, Experimenter # 2 can draw a sample  $(y_1, \dots, y_n)$  from  $P_{t=T(x_1, \dots, x_n)}^*$ . So,  $(y_1, \dots, y_n)$  is a realisation of  $(X_1, \dots, X_n)$ . We can show that Experimenter # 2 has lost *no* information:

$$\begin{aligned} \mathbb{P}[Y_1 = y_1, \dots, Y_n = y_n] &= \sum_{t \in \mathcal{T}} \mathbb{P}[Y_1 = y_1, \dots, Y_n = y_n, T(y_1, \dots, y_n) = t] \\ &= \mathbb{P}[Y_1 = y_1, \dots, Y_n = y_n, T(Y_1, \dots, Y_n) = T(y_1, \dots, y_n)] \\ &= \underbrace{\mathbb{P}[Y_1 = y_1, \dots, Y_n = y_n \mid T = T(y_1, \dots, y_n)]}_{P_{T(y_1, \dots, y_n)}^*} \cdot \mathbb{P}[T = T(y_1, \dots, y_n)] \\ &= \mathbb{P}[X_1 = y_1, \dots, X_n = y_n \mid T = T(y_1, \dots, y_n)] \cdot \mathbb{P}[T = T(y_1, \dots, y_n)] \\ &= \mathbb{P}[X_1 = y_1, \dots, X_n = y_n] \end{aligned}$$

Which shows that Experimenter # 2 has lost no information. Thus, *sufficiency is a useful concept when wanting to reduce data to statistics*. However, in general it is very hard to prove that a statistic is sufficient :(. So, we want a criteria to prove that a statistic is sufficient. This motivates the following “very handy” theorem.

**Theorem 23** (Neymen-Fisher Criterion). Let  $f(x_1, \dots, x_n; \theta)$  denote the joint PMF/PDF of  $(x_1, \dots, x_n)$ , a random sample from  $F_\theta$ ,  $\theta \in \Theta$ . A statistic  $T = T(x_1, \dots, x_n)$  is sufficient for  $\theta$  if and only if  $\forall \theta \in \Theta$ , there exist measurable functions  $g_\theta$  and  $h$  so that  $h$  does *NOT* depend on  $\theta$  for all  $(x_1, \dots, x_n) \in \mathcal{X}$ . That is, it can be decomposed in the form:

$$f(x_1, \dots, x_n; \theta) = h(x_1, \dots, x_n) g_\theta(T(x_1, \dots, x_n)) \quad (44)$$

*Proof.* We will only do the proof in the discrete case. “ $\Rightarrow$ ”: first assume that  $T$  is sufficient. Set

$$\begin{aligned} g_\theta(t) &:= \mathbb{P}_\theta[T(X_1, \dots, X_n) = t], \quad t \in \mathcal{T} \\ h(x_1, \dots, x_n) &:= \begin{cases} \mathbb{P}[X_1 = x_1, \dots, X_n = x_n \mid T = T(x_1, \dots, x_n)] & \text{if } \mathbb{P}[T = T(x_1, \dots, x_n)] > 0 \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Now apply the law of Total Probability and the definition of conditional probability:

$$\begin{aligned} \mathbb{P}[X_1 = x_1, \dots, X_n = x_n] &= \mathbb{P}[X_1 = x_1, \dots, X_n = x_n, T = T(x_1, \dots, x_n)] \\ &= h(x_1, \dots, x_n) \cdot g_\theta(T(x_1, \dots, x_n)) \end{aligned}$$

where  $h$  is  $\mathbb{P}[x \mid T = T(x_1, \dots, x_n)]$  and  $g_\theta$  is  $\mathbb{P}[T = T(x_1, \dots, x_n)]$ .

“ $\Leftarrow$ ”: Now assume that the PMF factors. Then:

$$\begin{aligned}
 \mathbb{P}[X_1 = x_1, \dots, X_n = x_n \mid T = t] &= \frac{\mathbb{P}[X_1 = x_1, \dots, X_n = x_n, T = t]}{\mathbb{P}[T = t]} \\
 &= \begin{cases} 0 & \text{if } t \neq T(x_1, \dots, x_n) \\ \frac{\mathbb{P}[X_1 = x_1, \dots, X_n = x_n, T = T(x_1, \dots, x_n)]}{\mathbb{P}[T = T(x_1, \dots, x_n)]} & \text{otherwise} \end{cases} \\
 &= \frac{\mathbb{P}[X_1 = x_1, \dots, X_n = x_n]}{\mathbb{P}[(x_1, \dots, x_n) \in A_{T(x_1, \dots, x_n)}]} \\
 &= \frac{\mathbb{P}[X_1 = x, \dots, X_n = x_n]}{\mathbb{P}[(X_1, \dots, X_n) \in A_{T(x_1, \dots, x_n)}]} \\
 &= \frac{h(x_1, \dots, x_n)g_\theta(T(x_1, \dots, x_n))}{\sum_{(y_1, \dots, y_n) \in A_{T(x_1, \dots, x_n)}} h(y_1, \dots, y_n)g_\theta(T(y_1, \dots, y_n))} \\
 &= T(x_1, \dots, x_n) \\
 &= \frac{h(x_1, \dots, x_n)}{\sum_{(y_1, \dots, y_n) \in A_{T(x_1, \dots, x_n)}} h(y_1, \dots, y_n)}
 \end{aligned}$$

which doesn't depend on  $\theta$ . This thus proves sufficiency.  $\square$

**Example 32.** Let  $X_1, \dots, X_n$  be a random sample from  $\mathcal{U}[0, \theta[$ , where  $\theta > 0$ . Then, the PMF is of the form:

$$f(x; \theta) = \begin{cases} \frac{1}{\theta} & x \in ]0, \theta[ \\ 0 & \text{else} \end{cases}$$

To write the joint PMF of the whole random sample, we need to account for the dependence of the interval on  $\theta$  (we need to write it out in terms of characteristic functions):

$$\begin{aligned}
 f(x_1, \dots, x_n; \theta) &= \prod_{i=1}^n \frac{1}{\theta} \chi_{x_i \in ]0, \theta[} \\
 &= \chi_{\min x_i > 0} \chi_{\max x_i < \theta} \left( \frac{1}{\theta} \right)^n
 \end{aligned}$$

Observe that we can set the first term to be  $h(x_1, \dots, x_n)$ , since there is no dependence on  $\theta$ , and that we can set the second term to be  $g_\theta(\max_{1 \leq i \leq n} x_i)$ . Since this is the likelihood function, we can conclude that the maximum order statistic,  $X_{(n)}$ , is sufficient for  $\theta$ .

**Example 33.** Let  $X_1, \dots, X_n$  be a random sample from  $\mathcal{N}(\mu, \sigma^2)$ . Then, the joint PDF is given by:

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

From the proof of Theorem 1.7, we know that:

$$\sum_{i=1}^n (x_i - \mu)^2 = (n-1)s^2 + n(\bar{x} - \mu)^2$$

and so we can re-write the above as:

$$f(x_1, \dots, x_n; \mu, \sigma^2) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \left( \frac{1}{\sigma} \right)^n \exp \left\{ \left( \frac{-(n-1)}{2\sigma^2} s^2 - \frac{n(\bar{x} - \mu)^2}{2\sigma^2} \right) \right\}$$

Here, we may choose:

$$\begin{aligned}
 h &:= \left( \frac{1}{\sqrt{2\pi}} \right)^n \\
 g_\theta(\bar{x}, s^2) &= \left( \frac{1}{\sigma} \right)^n \exp \left\{ \left( \frac{-(n-1)}{2\sigma^2} s^2 - \frac{n(\bar{x} - \mu)^2}{2\sigma^2} \right) \right\}
 \end{aligned}$$

So,  $(\bar{x}, s^2)$  is sufficient for  $(\mu, \sigma^2)$ . Now suppose that  $\sigma^2$  is known. Then we can set:

1.  $h, g$  as before, and obtain that  $(\bar{x}, s^2)$  is sufficient for  $\mu$ .
2. Or we set:

$$h(x_1, \dots, x_n) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \left( \frac{1}{\sigma} \right)^n \exp \left\{ \frac{-(n-1)s^2}{2\sigma^2} \right\}$$

which has no dependence on  $\mu$ , and so

$$g_\mu(\bar{x}) = \exp \left( \frac{-n(\bar{x} - \mu)^2}{2\sigma^2} \right)$$

and so  $\bar{x}$  is sufficient for  $\mu$ .

Note: there are some statistics that lead to greater data reduction than others, and this leads into an idea called “minimal sufficiency.” Also note that:

1. In (1), if we use  $(\bar{x}, s^2)$  and set  $\mathcal{X} = \mathbb{R}^n$ , then:

$$A_{(a,b)} = \{ (x_1, \dots, x_n) \in \mathbb{R}^n \mid \bar{x} = a, s^2 = b \}$$

In this case,  $a \in \mathbb{R}$  and  $b > 0$ .

2. In (2), if we only use  $\bar{x}$ , then  $\mathcal{T} = \mathbb{R}$  and:

$$A_a = \{ (x_1, \dots, x_n) \in \mathbb{R}^n \mid \bar{x} = a \}$$

If you compare these, observe that  $A_{(a,b)} \subseteq A_a$ . In a way,  $A_a$  is a projection, and forms a *coarser* partition than  $(\bar{x}, s^2)$  does.

## 4.2 Sufficient Statistics

To summarise what we have so far, we have that  $T(X_1, \dots, X_n)$  is a statistic taking values in  $\mathcal{T}$ . Then,  $\{A_t, t \in \mathcal{T}\}$  is a partition of the sample space  $\mathcal{X}$ , where

$$A_t = \{ (x_1, \dots, x_n) \in \mathcal{X} \mid T(x_1, \dots, x_n) = t \}$$

This means that the statistics induce a partition on  $\mathcal{X}$ .

Insert graphic here.

**Definition 16** (Sufficient Statistic). A sufficient statistic  $T$  is called minimally sufficient if for any other sufficient statistic  $T^*$ , we have the following:

$$\begin{aligned} \forall (x_1, \dots, x_n), (y_1, \dots, y_n) \in \mathcal{X} \\ T^*(x_1, \dots, x_n) = T^*(y_1, \dots, y_n) \Rightarrow T(x_1, \dots, x_n) = T(y_1, \dots, y_n) \end{aligned}$$

Note that if  $T$  is minimally sufficient, and if  $T^*$  is sufficient, then  $\{A_t \mid t \in \mathcal{T}\}$  (the partition of  $\mathcal{X}$  corresponding to  $T$ ) is coarser than the partition  $\{A_{t^*}^*, t^* \in \mathcal{T}^*\}$  corresponding to  $T^*$ . Mathematically:

$$\forall t^* \in \mathcal{T}, \exists t \in \mathcal{T} \text{ s.t. } A_{t^*}^* \subseteq A_t$$

Here,  $A_{t^*}^* = \{ (x_1, \dots, x_n) \in \mathcal{X} \mid T^*(x_1, \dots, x_n) = t^* \}$ .

The next theorem applies to statistical models where  $F_\theta$  has a density:

**Theorem 24** (Lehmann-Schafte Criteria). Suppose that  $(X_1, \dots, X_n)$  has a PDF/PMF  $f(x_1, \dots, x_n; \theta)$ ,  $\theta \in \Theta$ . Suppose that  $T$  is a statistic such that for any  $(x_1, \dots, x_n)$  and  $(y_1, \dots, y_n) \in \mathcal{X}$  such that the ratio:

$$\frac{f(x_1, \dots, x_n; \theta)}{f(y_1, \dots, y_n; \theta)}$$

is constant as a function of  $\theta \iff T(x_1, \dots, x_n) = T(y_1, \dots, y_n)$ . This property gives us that  $T$  is minimally sufficient.



**Example 34.** Let  $X_1, \dots, X_n$  be a random sample from  $\mathcal{N}(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ . From last class, the density can be broken down as:

$$\begin{aligned} f(x_1, \dots, x_n; \mu, \sigma^2) &= \left(\frac{1}{\sqrt{2\pi}}\right)^n \left(\frac{1}{\sigma}\right)^n \exp\left\{-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right\} \\ &= \left(\frac{1}{\sqrt{2\pi}}\right)^n \left(\frac{1}{\sigma^2}\right)^n \exp\left\{-\frac{(n-1)s^2}{2\sigma^2} - \frac{n(\bar{x} - \mu)^2}{2\sigma^2}\right\} \end{aligned}$$

**Q:** is it minimally sufficient? By the previous theorem:

$$\frac{f(x_1, \dots, x_n; \mu, \sigma^2)}{f(y_1, \dots, y_n; \mu, \sigma^2)} = \exp\left\{-\frac{(n-1)}{2\sigma^2}(s_x^2 - s_y^2) - \frac{n}{2\sigma^2}[(\bar{x} - \mu^2) - (\bar{y} - \mu^2)]\right\}$$

where

$$\begin{aligned} s_x^2 &= \frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2 \\ s_y^2 &= \frac{1}{(n-1)} \sum_{i=1}^n (y_i - \bar{y})^2 \end{aligned}$$

This ratio does not depend on  $\mu$  and  $\sigma^2$  if and only if  $\bar{x} = \bar{y}$  and  $s_x^2 = s_y^2$ , which implies that  $(\bar{x}, s^2)$  is minimally sufficient for  $(\mu, \sigma^2)$ .

*Proof.* Proof of the Theorem. We need to prove that  $T$  is sufficient and minimally sufficient.

1. ( $T$  is sufficient for  $\theta$ ): We will show this using the Neymann-Fischer criteria. As before, consider  $\{A_t, t \in \mathcal{T}\}$ . For any  $A_t$ , choose a point  $x_t \in A_t$ . Now, fix an arbitrary  $(x_1, \dots, x_n) \in \mathcal{X}$ . Then we have that  $(x_1, \dots, x_n) \in A_{T(x_1, \dots, x_n)}$ . Let  $x_{T(x_1, \dots, x_n)}$  be a fixed point in  $A_{T(x_1, \dots, x_n)}$ . This is the representative. We thus have that  $(x_1, \dots, x_n) \in A_{T(x_1, \dots, x_n)}$ . So, we have that

$$\frac{f(x_1, \dots, x_n; \theta)}{f(x_{T(x_1, \dots, x_n)}; \theta)}$$

does not depend on  $\theta$  by assumption. Now set:

$$\begin{aligned} \frac{f(x_1, \dots, x_n; \theta)}{f(x_{T(x_1, \dots, x_n)}; \theta)} &= h(x_1, \dots, x_n) \\ f(x_{T(x_1, \dots, x_n)}; \theta) &= g_\theta(T(x_1, \dots, x_n)) \end{aligned}$$

So, we have thus written

$$f(x_1, \dots, x_n; \theta) = h(x_1, \dots, x_n) g_\theta(T(x_1, \dots, x_n))$$

and so by the Neymann Fischer Criteria,  $T$  is sufficient.

2. ( $T$  is minimally sufficient). Let  $T^*$  be another sufficient statistic. Then, from the Neymann Fischer criteria, we know that  $\exists h^*, g_\theta^*$  for which  $\forall (x_1, \dots, x_n) \in \mathcal{X}$ :

$$f(x_1, \dots, x_n; \theta) = h^*(x_1, \dots, x_n) \cdot g_\theta^*(T^*(x_1, \dots, x_n))$$

Now pick an arbitrary  $(x_1, \dots, x_n), (y_1, \dots, y_n) \in \mathcal{X}$  so that  $T^*(x_1, \dots, x_n) = T^*(y_1, \dots, y_n)$ . We need to prove that then also  $T(x_1, \dots, x_n) = T(y_1, \dots, y_n)$ . But:

$$\frac{f(x_1, \dots, x_n; \theta)}{f(y_1, \dots, y_n; \theta)} = \frac{h^*(x_1, \dots, x_n) \cdot g_\theta^*(T^*(x_1, \dots, x_n))}{h^*(y_1, \dots, y_n) \cdot g_\theta^*(T^*(y_1, \dots, y_n))}$$

But we assumed that the  $g^*$  are the same and so those drop out and we obtain:

$$= \frac{h^*(x_1, \dots, x_n)}{h^*(y_1, \dots, y_n)}$$

which does not depend on  $\theta$ . This implies that  $T(x_1, \dots, x_n) = T(y_1, \dots, y_n)$ .

□

**Example 35.** Let  $X_1, \dots, X_n$  be a random sample taken from  $\mathcal{U}[\theta, \theta + 1[$  where  $\theta \in \mathbb{R}$ . Then, the probability density of the sample is:

$$f(x_1, \dots, x_n; \theta) = \begin{cases} 1 & \text{if } x_i \in ]\theta, \theta + 1[ \forall i \in \{1, \dots, n\} \\ 0 & \text{otherwise} \end{cases}$$

$$= \chi_{\min_{1 \leq i \leq n} x_i > \theta} \chi_{\max_{1 \leq i \leq n} x_i < \theta + 1}$$

for  $(x_1, \dots, x_n)$  and  $(y_1, \dots, y_n)$ . Thus,

$$\frac{f(x_1, \dots, x_n; \theta)}{f(y_1, \dots, y_n; \theta)}$$

does not depend on  $\theta$  if  $\min_{1 \leq i \leq n} x_i = \min_{1 \leq i \leq n} y_i$  and  $\max_{1 \leq i \leq n} x_i = \max_{1 \leq i \leq n} y_i$ . So,  $(\min_{1 \leq i \leq n}(x_i), \max_{1 \leq i \leq n}(x_i))$  is minimally sufficient for  $\theta$ .

### 4.3 Applications of Sufficiency in Point Estimation

Let  $X_1, \dots, X_n$  be a random sample from  $F_\theta$ ,  $\theta \in \Theta$ . Our goal is to find the uniformly minimum variance unbiased estimator of  $\gamma(\theta)$  (the UMVUE). We will now denote our estimators by  $W$ . We want to find  $W$  so that

$$\mathbb{E}_\theta[W] = \gamma(\theta) \quad \forall \theta \in \Theta$$

and for any other unbiased estimator  $W^*$  of  $\gamma(\theta)$ , we have:

$$\text{Var}_\theta[W] \leq \text{Var}_\theta[W^*] \quad \forall \theta \in \Theta$$

First recall three important facts from Math 356. For any two random variables  $Y, Z$ :

1.  $\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|Z]]$  (Iterative expectation)
2.  $\mathbb{E}[Y|Z] = h(Z)$  for some measurable function  $h$ . That is, this is a function of  $Z$ .
3.  $\text{Var}[Y] = \text{Var}[\mathbb{E}[Y|Z]] + \mathbb{E}[\text{Var}[Y|Z]] \geq 0$ .

**Theorem 25** (Rao-Blackwell). Let  $W$  be an unbiased estimator of  $\gamma(\theta)$  with finite variance and let  $T$  be a sufficient statistic for  $\theta$ . Then, if  $W^* := \mathbb{E}[W|T]$ , then:

1.  $W^*$  is an unbiased estimator of  $\gamma(\theta)$ .
2.  $\text{Var}_\theta[W^*] \leq \text{Var}_\theta[W]$  for all  $\theta \in \Theta$ .

*Proof.* 1.  $\mathbb{E}_\theta[W^*] = \mathbb{E}_\theta[\mathbb{E}_\theta[W|T]] = \mathbb{E}_\theta[W] = \gamma(\theta)$  for all  $\theta \in \Theta$ .

2.  $\text{Var}_\theta[W] = \text{Var}_\theta[\mathbb{E}_\theta[W|T]] + \mathbb{E}_\theta[\text{Var}_\theta[W|T]] = \text{Var}_\theta[W^*] + \mathbb{E}_\theta[\text{Var}_\theta[W|T]]$ . Since the  $\mathbb{E}_\theta[\text{Var}_\theta[W|T]]$  term is positive, this implies that  $\text{Var}_\theta[W] \geq \text{Var}_\theta[W^*]$  for  $\theta \in \Theta$ .

We need to use the sufficiency – finally, since  $T$  is sufficient, the distribution of  $(X_1, \dots, X_n)$  given  $T$  does not depend on  $\theta$ , which means that  $W^*$  does not depend on  $\theta$ , which gives us the fact that  $W^*$  is an estimator. □

**Example 36.** Assume that  $X_1, X_2 \sim \mathcal{N}(\theta, 1)$ . Let's investigate the sample mean,  $\bar{x} = (1/2)(x_1 + x_2)$ . Computing the quantities:

$$\mathbb{E}[\bar{x}] = \theta$$

$$\text{Var}[\bar{x}] = 1/2$$

$$\mathbb{E}[\bar{x}|x_1] = \mathbb{E}\left[\frac{1}{2}(x_1 + x_2)|x_1\right] = \frac{1}{2}\mathbb{E}[x_1|x_1] + \mathbb{E}[x_2|x_1] = \frac{1}{2}(x_1 + \theta) \equiv W^*$$

Note that this is not an estimator, since it involves the unknown parameter. So, we need to repair it.

$$\mathbb{E}[W^*] = \frac{1}{2}(\mathbb{E}[X_1] + \theta) = \frac{1}{2}(\theta + \theta) = \theta$$

and

$$\text{Var}[W^*] = \frac{1}{4}\text{Var}[X_1] = \frac{1}{4} < \frac{1}{2}$$

because  $X_1$  is not sufficient.

**Corollary 2.** Let  $W$  be any estimator of  $\gamma(\theta)$ , not necessarily unbiased. Then,  $W^* = \mathbb{E}[W|T]$  where  $T$  is a sufficient estimator of  $\gamma(\theta)$  so that,  $\forall \theta \in \Theta$ :

$$\text{MSE}_\theta[W^*] \leq \text{MSE}_\theta[W]$$

Note: if  $W, Z$  are arbitrary random variables. Then:

$$\mathbb{E}[W|Z] = h(Z)$$

for some measurable function  $Z$ . We have:

$$\mathbb{E}[h(Z)|Z] = h(Z)$$

**Intuition from Rao Blackwell:** we should condition on a minimally sufficient statistic to achieve the greatest reduction of the variance possible. **Q:** if we condition on a minimally sufficient statistic, will we get the UMVUE/an UMVUE? Unfortunately, no.

**Example 37.** (Nasty Way of Improving an Estimator): Suppose that  $W$  is some estimator of  $\gamma(\theta)$ , unbiased, and has finite variance. Suppose that  $U$  is a statistic,  $U(X_1, \dots, X_n)$  so that  $\mathbb{E}_\theta[U] = 0$  for all  $\theta \in \Theta$  (i.e.,  $U$  is an estimator of 0, which means that it is unbiased). Now let:

$$W_a^* := W + aU$$

where  $a \in \mathbb{R}$ . We get that  $\mathbb{E}[W_a^*] = \gamma(\theta)$  and  $\text{Var}[W_a^*] = \text{Var}_\theta[W] + a^2\text{Var}_\theta[U] + 2a\text{Cov}_\theta[W, U]$ . We have problems when  $\text{cov}_{\theta_0}[W, U] \neq 0$  for some  $\theta_0$ . If  $\text{cov}_{\theta_0}[W, U] < 0$ , then for  $a \in ]0, ((-2\text{cov}_{\theta_0}[W, U])/\text{var}_{\theta_0}[U])[$ , we have that  $\text{var}_{\theta_0}[W_a^*] < \text{var}_{\theta_0}[W]$ . This means that  $W$  is not an UMVUE since we found an improvement. We will prevent this from happening.

**Definition 17.** A family  $\{F_\theta \mid \theta \in \Theta\}$  of distributions is called **complete** if for any measurable function  $g$  and  $X \sim F_\theta$  such that

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

it holds that  $g(X)$  is almost surely zero, i.e.:

$$\mathbb{P}[g(X) = 0] = 1 \quad \forall \theta \in \Theta$$

A statistic  $T$  is called **complete** if for any measurable function  $g : \mathcal{T} \rightarrow \mathbb{R}$ ,

$$(\forall \theta \in \Theta, \mathbb{E}_\theta[g(T)] = 0) \Rightarrow (\mathbb{P}_\theta[g(T) = 0] = 1 \quad \forall \theta \in \Theta)$$

**Theorem 26** (Lehmann Scheffe Theorem). Consider a random variable  $X_1, \dots, X_n$  from  $F_\theta$ ,  $\theta \in \Theta$ . Let  $\gamma(\theta)$  be a real-valued parameter of interest and  $W$  an unbiased estimator for  $\gamma(\theta)$ . If  $T$  is a complete and sufficient statistic, then  $W^* = \mathbb{E}[W|T]$  is THE UMVUE of  $\gamma(\theta)$ .