

# GR5204: Statistical Inference\*

Johannes Wiesel  
Columbia University

April 11, 2023

## Contents

<b>1</b>	<b>Recap of Probability Theory</b>	<b>5</b>
<b>2</b>	<b>Statistical Inference: Introduction</b>	<b>10</b>
2.1	Statistical model . . . . .	10
2.2	Method of moments estimators . . . . .	12
<b>3</b>	<b>Method of Maximum Likelihood</b>	<b>16</b>
3.1	Properties of MLEs . . . . .	20
3.2	Computational methods for approximating MLEs . . . . .	20
<b>4</b>	<b>Principles of estimation</b>	<b>21</b>
4.1	Mean squared error . . . . .	21
4.2	Comparing estimators . . . . .	23
4.3	Unbiased estimators . . . . .	24
4.4	Sufficient Statistics . . . . .	26
<b>5</b>	<b>The sampling distribution of a statistic</b>	<b>30</b>
5.1	The gamma and the $\chi^2$ distributions . . . . .	30

---

\*adapted from Prof. Bodhisettva Sen's and Prof. Thibault Vatter's notes

5.1.1	The gamma distribution . . . . .	30
5.1.2	The Chi-squared distribution . . . . .	31
5.2	Sampling from a normal population . . . . .	32
5.3	The $t$ -distribution . . . . .	35
<b>6</b>	<b>Confidence intervals</b>	<b>36</b>
6.1	Construction of confidence interval using a pivot . . . . .	36
6.2	Asymptotic confidence intervals . . . . .	38
<b>7</b>	<b>The Cramér–Rao Information Inequality</b>	<b>41</b>
7.1	Information . . . . .	42
7.2	Examples . . . . .	45
7.3	Large sample properties of the MLE . . . . .	47
<b>8</b>	<b>Bayesian paradigm</b>	<b>52</b>
8.1	Prior distribution . . . . .	52
8.2	Posterior distribution . . . . .	53
8.3	Bayes Estimators . . . . .	54
8.4	Sampling from a normal distribution . . . . .	55
<b>9</b>	<b>Hypothesis Testing</b>	<b>57</b>
9.1	Principles of Hypothesis Testing . . . . .	57
9.2	Critical regions and test statistics . . . . .	58
9.3	Power function and types of error . . . . .	59
9.4	Significance level . . . . .	62
9.5	$P$ -value . . . . .	64
9.6	Testing simple hypotheses: optimal tests . . . . .	65
9.6.1	Minimizing the $\mathbb{P}$ (Type-II error) . . . . .	65
9.7	Uniformly most powerful (UMP) tests . . . . .	66
9.8	The $t$ -test . . . . .	67

9.8.1	Testing hypotheses about the mean with unknown variance . . .	67
9.8.2	One-sided alternatives . . . . .	70
9.9	Comparing the means of two normal distributions (two-sample $t$ test)	72
9.9.1	One-sided alternatives . . . . .	72
9.9.2	Two-sided alternatives . . . . .	73
9.10	Comparing the variances of two normal distributions ( $F$ -test) . . . .	73
9.10.1	One-sided alternatives . . . . .	74
9.10.2	Two-sided alternatives . . . . .	75
9.11	Likelihood ratio test . . . . .	75
9.12	Equivalence of hypothesis tests and confidence sets . . . . .	76
<b>10</b>	<b>Linear regression</b>	<b>80</b>
10.1	Simple linear regression . . . . .	80
10.1.1	Interpretation . . . . .	81
10.2	Method of least squares . . . . .	81
10.2.1	Normal equations . . . . .	82
10.2.2	Estimated regression function . . . . .	83
10.2.3	Properties . . . . .	83
10.2.4	Estimation of $\sigma^2$ . . . . .	84
10.2.5	Gauss-Markov theorem . . . . .	84
10.3	Normal simple linear regression . . . . .	85
10.3.1	Maximum likelihood estimation . . . . .	85
10.3.2	Inference . . . . .	86
10.3.3	Inference about $\beta_1$ . . . . .	86
10.3.4	Sampling distribution of $\hat{\beta}_0$ . . . . .	90
10.3.5	Mean response . . . . .	91
10.3.6	Prediction interval . . . . .	93
10.3.7	Inference about both $\beta_0$ and $\beta_1$ simultaneously . . . . .	94

10.3.8 Examples . . . . . 94

# 1 Recap of Probability Theory

**Definition 1** (Sample mean). *Suppose that  $X_1, X_2, \dots, X_n$  are i.i.d. random variables with (unknown) mean  $\mu \in \mathbb{R}$  (i.e.,  $\mathbb{E}(X_1) = \mu$ ) and variance  $\sigma^2 < \infty$ . A natural “estimator” of  $\mu$  is the **sample mean** (or **sample average**) defined as*

$$\bar{X}_n := \frac{1}{n}(X_1 + \dots + X_n) = \frac{1}{n} \sum_{i=1}^n X_i.$$

**Lemma 1.1.**  $\mathbb{E}(\bar{X}_n) = \mu$  and  $\text{Var}(\bar{X}_n) = \sigma^2/n$ .

*Proof.* Observe that

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

Also,

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

**Theorem 1.2** (Weak law of large numbers). *Suppose that  $X_1, X_2, \dots, X_n$  are  $n$  i.i.d. random variables with finite mean  $\mu$ . Then for any  $\epsilon > 0$ , we have*

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}(X) \right| > \epsilon \right) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

This says that if we take the sample average of  $n$  i.i.d random variables the sample average will be close to the true population average. Figure 1 illustrates the result: The left panel shows the density of the data generating distribution (in this example we took  $X_1, \dots, X_n$  i.i.d.  $\text{Exp}(10)$ ); the middle and right panels show the distribution (histogram obtained from 1000 replicates) of  $\bar{X}_n$  for  $n = 100$  and  $n = 1000$ , respectively. We see that as the sample size increases, the distribution of the sample mean concentrates around  $\mathbb{E}(X_1) = 1/10$  (i.e.,  $\bar{X}_n \xrightarrow{\mathbb{P}} 10^{-1}$  as  $n \rightarrow \infty$ ).

**Definition 2** (Convergence in probability). *In the above, we say that the sample mean  $\frac{1}{n} \sum_{i=1}^n X_i$  converges in probability to the true (population) mean.*

*More generally, we say that a sequence of random variables  $\{Z_n\}_{n=1}^\infty$  converges to  $Z$  in probability, and write*

$$Z_n \xrightarrow{\mathbb{P}} Z,$$

*if for every  $\epsilon > 0$ ,*

$$\mathbb{P}(|Z_n - Z| > \epsilon) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

*This is equivalent to saying that for every  $\epsilon > 0$ ,*

$$\lim_{n \rightarrow \infty} \mathbb{P}(|Z_n - Z| \leq \epsilon) = 1.$$

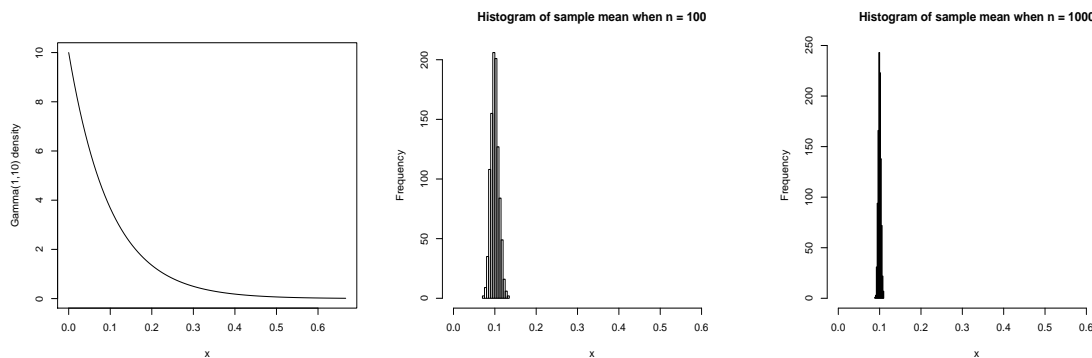


Figure 1: The plots illustrate the convergence (in probability) of the sample mean to the population mean.

**Definition 3** (Convergence in distribution). *We say a sequence of random variables  $\{Z_n\}_{i=1}^n$  with c.d.f's  $F_n(\cdot)$  **converges in distribution** to  $F$  if*

$$\lim_{n \rightarrow \infty} F_n(u) = F(u)$$

*for all  $u$  such that  $F$  is continuous<sup>1</sup> at  $u$  (here  $F$  is itself a c.d.f).*

The second fundamental result in probability theory, after the law of large numbers (LLN), is the Central limit theorem (CLT), stated below. The CLT gives us the approximate (asymptotic) distribution of  $\bar{X}_n$

**Theorem 1.3** (Central limit theorem). *If  $X_1, X_2, \dots$  are i.i.d with mean zero and variance 1, then*

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \xrightarrow{d} N(0, 1),$$

*where  $N(0, 1)$  is the standard normal distribution. More generally, the usual rescaling tell us that if  $X_1, X_2, \dots$  are i.i.d with mean  $\mu$  and variance  $\sigma^2 < \infty$ , then*

$$\sqrt{n}(\bar{X}_n - \mu) \equiv \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) \xrightarrow{d} N(0, \sigma^2).$$

---

<sup>1</sup>Explain why do we need to restrict our attention to continuity points of  $F$ . (Hint: think of the following sequence of distributions:  $F_n(u) = I(u \geq 1/n)$ , where the “indicator” function of a set  $A$  is one if  $x \in A$  and zero otherwise.) It’s worth emphasizing that convergence in distribution — because it only looks at the c.d.f. — is in fact **weaker** than convergence in probability. For example, if  $p_X$  is symmetric, then the sequence  $X, -X, X, -X, \dots$  trivially converges in distribution to  $X$ , but obviously doesn’t converge in probability. Also, if  $U \sim \text{Unif}(0, 1)$ , then the sequence

$$U, 1 - U, U, 1 - U, \dots$$

converge in distribution to a uniform distribution. But obviously they do not converge in probability.

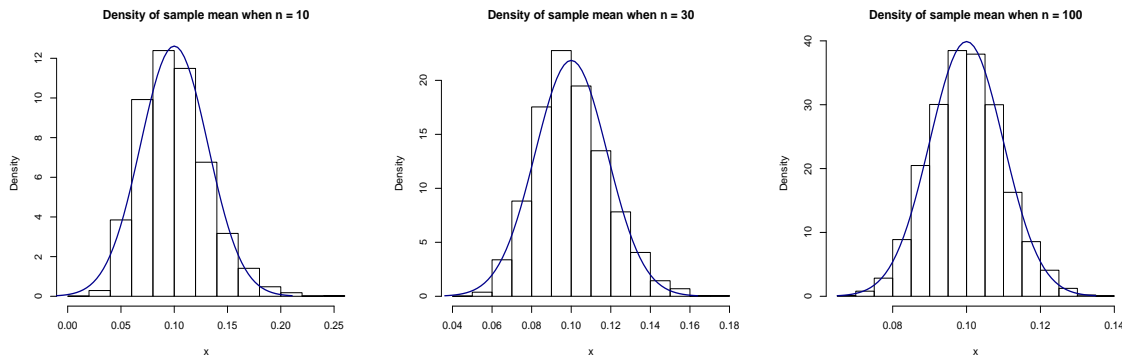


Figure 2: The plots illustrate the convergence (in distribution) of the sample mean to a normal distribution.

The following plots illustrate the CLT: The left, center and right panels of Figure 2 show the (scaled) histograms of  $\bar{X}_n$  when  $n = 10, 30$  and  $100$ , respectively (as before, in this example we took  $X_1, \dots, X_n$  i.i.d.  $\text{Exp}(10)$ ; the histograms are obtained from 5000 independent replicates). We also overplot the normal density with mean  $0.1$  and variance  $10^{-1}/\sqrt{n}$ . The remarkable agreement between the two densities illustrates the power of the CLT. Observe that the original distribution of the  $X_i$ 's,  $\text{Exp}(10)$ , is skewed and highly non-normal, but even for  $n = 10$ , the distribution of  $\bar{X}_{10}$  is quite close to being normal.

Another class of useful results we will use very much in this course go by the name “continuous mapping theorem”. Here are two such results.

**Theorem 1.4.** *If  $Z_n \xrightarrow{\mathbb{P}} b$  and if  $g(\cdot)$  is a function that is continuous at  $b$ , then*

$$g(Z_n) \xrightarrow{\mathbb{P}} g(b).$$

**Theorem 1.5.** *If  $Z_n \xrightarrow{d} Z$  and if  $g(\cdot)$  is a function that is continuous, then*

$$g(Z_n) \xrightarrow{d} g(Z).$$

The last result that we need from probability theory—and this may be new to you—is the so-called **delta method**. It allows us to find the asymptotic distribution of a *continuous transformation* of the rescaled sample mean. But let’s first state the abstract result.

**Theorem 1.6.** *Let  $Z_1, Z_2, \dots, Z_n$  be a sequence of random variables and let  $Z$  be a random variable with a continuous c.d.f  $F^*$ . Let  $\theta \in \mathbb{R}$ , and let  $a_1, a_2, \dots$ , be a sequence such that  $a_n \rightarrow \infty$ . Suppose that*

$$a_n(Z_n - \theta) \xrightarrow{d} F^*.$$

*Let  $g(\cdot)$  be a function with a continuous derivative such that  $g'(\theta) \neq 0$ . Then*

$$a_n \frac{g(Z_n) - g(\theta)}{g'(\theta)} \xrightarrow{d} F^*.$$

*Proof.* We will only give an outline of the proof (think  $a_n = n^{1/2}$ , if  $Z_n$  as the sample mean). As  $a_n \rightarrow \infty$ ,  $Z_n$  must get close to  $\theta$  with high probability as  $n \rightarrow \infty$ . As  $g(\cdot)$  is continuous,  $g(Z_n)$  will be close to  $g(\theta)$  with high probability. Let's say  $g(\cdot)$  has a Taylor expansion around  $\theta$ , i.e.,

$$g(Z_n) \approx g(\theta) + g'(\theta)(Z_n - \theta),$$

where we have ignored all terms involving  $(Z_n - \theta)^2$  and higher powers. Then if

$$a_n(Z_n - \theta) \xrightarrow{d} Z,$$

for some limit distribution  $F^*$  and a sequence of constants  $a_n \rightarrow \infty$ , then

$$a_n \frac{g(Z_n) - g(\theta)}{g'(\theta)} \approx a_n(Z_n - \theta) \xrightarrow{d} F^*.$$

□

In other words, limit distributions are passed through functions in a pretty simple way. We'll be using the delta method a lot. The main application is when we've already proven a CLT for  $Z_n$ , that is, when

$$\frac{\sqrt{n}(Z_n - \mu)}{\sigma} \xrightarrow{d} N(0, 1).$$

Then, by Theorem 1.6,

$$\sqrt{n}(g(Z_n) - g(\mu)) \xrightarrow{d} N(0, \sigma^2(g'(\mu))^2).$$

**Exercise 1:** Assume  $n^{1/2}Z_n \xrightarrow{d} N(0, 1)$ . What is the asymptotic distribution of

1.  $g(Z_n) = (Z_n - 1)^2$ ?
2. What about  $g(Z_n) = Z_n^2$ ? Does anything go wrong when applying the delta method in this case? Can you fix this problem?

Let's illustrate the theory through an example.

**Example 1.7.** A company sells a certain kind of electronic component. The company is interested in knowing about *how long* a component is likely to last on average. They can collect data on many such components that have been used under typical conditions. They choose to use the family of *exponential* distributions<sup>2</sup> to model the length of time (in years) from when a component is put into service until it fails.

The company believes that, if they knew the failure rate  $\theta$ , then  $\mathbf{X}_n = (X_1, X_2, \dots, X_n)$  would be  $n$  i.i.d random variables having the exponential distribution with rate  $\theta$ . Based on these hypotheses, we have the following results:

---

<sup>2</sup> $X$  has an exponential distribution with rate  $\theta > 0$  (in symbols,  $X \sim \text{Exp}(\theta)$ ), if the p.d.f of  $X$  is given by

$$f_\theta(x) = \theta e^{-\theta x} \mathbf{1}_{[0, \infty)}(x), \quad \text{for } x \in \mathbb{R}.$$

The mean of  $X$  is given by  $\mathbb{E}(X) = \theta^{-1}$ , and the variance of  $X$  is  $\text{Var}(X) = \theta^{-2}$ .



- by the LLN, the sample mean  $\bar{X}_n$  converges in probability to the expectation  $1/\theta$ , that is,

$$\bar{X}_n \xrightarrow{\mathbb{P}} \frac{1}{\theta};$$

- by the continuous mapping theorem (see Theorem 1.4)  $\bar{X}_n^{-1}$  converges in probability to  $\theta$ , i.e.,

$$\bar{X}_n^{-1} \xrightarrow{\mathbb{P}} \theta;$$

- by the CLT, we know that

$$\sqrt{n}(\bar{X}_n - \theta^{-1}) \xrightarrow{d} N(0, \theta^{-2})$$

where  $\text{Var}(X_1) = \theta^{-2}$ ;

- By the delta method, we can show that

$$\sqrt{n}(\bar{X}_n^{-1} - \theta) \xrightarrow{d} N(0, (\theta^2)^2 \theta^{-2}),$$

where we have considered  $g(x) = \frac{1}{x}$  with  $g'(x) = -\frac{1}{x^2}$  (observe that  $g$  is continuous on  $(0, \infty)$ ). Note that the variance of  $X_1$  is  $\text{Var}(X_1) = \theta^{-2}$ .

## 2 Statistical Inference: Introduction

### 2.1 Statistical model

**Definition 4** (Statistical model). *A statistical model is*

- *an identification of random variables of interest,*
- *a specification of a joint distribution or a family of possible joint distributions for the observable random variables,*
- *the identification of any parameters of those distributions that are assumed unknown,*
- *(Bayesian approach, if desired) a specification for a (joint) distribution for the unknown parameter(s).*

**Definition 5** (Statistical Inference). *Statistical inference is a procedure that produces a probabilistic statement about some or all parts of a statistical model.*

**Definition 6** (Parameter space). *The set  $\Omega$  of all possible values of a parameter  $\theta$  or of a vector of parameters  $\theta = (\theta_1, \dots, \theta_k)$  is called the parameter space.*

**Example 2.1.**

- The family of *binomial* distributions has parameters  $n$  and  $p$ .
- The family of *normal* distributions is parameterized by the mean  $\mu$  and variance  $\sigma^2$  of each distribution (so  $\theta = (\mu, \sigma^2)$  can be considered a pair of parameters, and  $\Omega = \mathbb{R} \times \mathbb{R}^+$ ).
- The family of *exponential* distributions is parameterized by the rate parameter  $\theta$  (the failure rate must be positive:  $\Omega$  will be the set of all positive numbers).

Note: The parameter space  $\Omega$  must contain all possible values of the parameters in a given problem.

**Example 2.2.** Suppose that  $n$  patients are going to be given a treatment for a condition and that we will observe for each patient whether or not they recover from the condition.

For each patient  $i = 1, 2, \dots$ , let  $X_i = 1$  if patient  $i$  recovers, and let  $X_i = 0$  if not. As a collection of possible distributions for  $X_1, X_2, \dots$ , we could choose to say that the  $X_i$ 's are i.i.d having the Bernoulli distribution with parameter  $p$ , for  $0 \leq p \leq 1$ .

In this case, the parameter  $p$  is known to lie in the closed interval  $[0, 1]$ , and this interval could be taken as the parameter space. Notice also that by the LLN,  $p$  is the limit as  $n \rightarrow \infty$  of the proportion of the first  $n$  patients who recover.

**Definition 7** (Statistic). *Suppose that the observable random variables of interest are  $X_1, \dots, X_n$ . Let  $\varphi$  be a real-valued function of  $n$  real variables. Then the random variable  $T = \varphi(X_1, \dots, X_n)$  is called a **statistic**.*

**Example 2.3.**

- the sample mean  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ ;
- the maximum  $X_{(n)}$  of the values  $X_1, \dots, X_n$ ;
- the sample variance  $s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  of the values  $X_1, \dots, X_n$ .

**Definition 8** (Estimator/Estimate). *Let  $X_1, \dots, X_n$  be observable data whose joint distribution is indexed by a parameter  $\theta$  taking values in a subset  $\Omega$  of the real line. An **estimator**  $\hat{\theta}_n$  of the parameter  $\theta$  is a real-valued function  $\hat{\theta}_n = \varphi(X_1, \dots, X_n)$ . If  $\{X_1 = x_1, \dots, X_n = x_n\}$  is observed, then  $\varphi(x_1, \dots, x_n)$  is called the **estimate** of  $\theta$ .*

*More generally, let  $X_1, \dots, X_n$  be observable data whose joint distribution is indexed by a parameter  $\theta$  taking values in a subset  $\Omega$  of  $d$ -dimensional space, i.e.,  $\Omega \subset \mathbb{R}^d$ . Let  $h : \Omega \rightarrow \mathbb{R}^d$ , be a function from  $\Omega$  into  $d$ -dimensional space. Define  $\psi = h(\theta)$ . An **estimator** of  $\psi$  is a function  $g(X_1, \dots, X_n)$  that takes values in  $d$ -dimensional space. If  $\{X_1 = x_1, \dots, X_n = x_n\}$  are observed, then  $g(x_1, \dots, x_n)$  is called the **estimate** of  $\psi$ .*

When  $h$  in Definition 8 is the identity function  $h(\theta) = \theta$ , then  $\psi = \theta$  and we are estimating the original parameter  $\theta$ . When  $g(\theta)$  is one coordinate of  $\theta$ , then the  $\psi$  that we are estimating is just that one coordinate.

Notice that an estimator need not be a “good” one. In fact, any transformation of the observations  $X_1, \dots, X_n$  is an estimator by Definition 8. For example, if  $\theta$  is the unknown mean of distribution,

$$\hat{\theta}_n = X_1 + \sum_{i=4}^n X_i - e^{X_2 + X_3}$$

is formally an estimator for  $\theta$  (but probably really bad one). So we need some sort of criteria to evaluate how good an estimator is. Here is one criterion that is often used.

**Definition 9** (Consistent estimator). *A sequence of estimators  $\hat{\theta}_n$  is said to be **consistent** for the unknown parameter  $\theta$  if  $\hat{\theta}_n$  converges in probability to  $\theta$ , that is, if for every  $\epsilon > 0$ ,*

$$\mathbb{P}(|\hat{\theta}_n - \theta| > \epsilon) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

In the following, we shall discuss three types of estimators:

- **Method of moments** estimators,
- **Maximum likelihood** estimators, and
- **Bayes** estimators.

## 2.2 Method of moments estimators

The *method of moments* (MOM) is an intuitive method for estimating parameters when other, more attractive, methods may be too difficult (to implement/compute).

**Definition 10** (Method of moments estimator). *Assume that  $X_1, \dots, X_n$  are observations from a distribution that is indexed by a  $k$ -dimensional parameter  $\theta$  and that has at least  $k$  finite moments. For  $j = 1, \dots, k$ , let*

$$\mu_j(\theta) := \mathbb{E}_\theta(X_1^j)$$

*be the  **$j$ th moment** of  $X_1, \dots, X_n$ . Suppose that the function*

$$\mu(\theta) = (\mu_1(\theta), \dots, \mu_k(\theta))$$

*is a one-to-one function of  $\theta$ . Let  $M(\mu_1, \dots, \mu_k)$  denote the inverse function, that is, for all  $\theta$ ,*

$$\theta = M(\mu_1, \dots, \mu_k).$$

*Define the  **$j$ th sample moment** as*

$$\hat{\mu}_j := \frac{1}{n} \sum_{i=1}^n X_i^j \quad \text{for } j = 1, \dots, k.$$

*The method of moments estimator of  $\theta$  is  $M(\hat{\mu}_1, \dots, \hat{\mu}_k)$ .*

Equivalently, the method of moments estimators can be obtained by setting up  $k$  equations

$$\hat{\mu}_j = \mu_j(\theta), \quad \text{for } j = 1, \dots, k,$$

and then solving for  $\theta$ .

**Theorem 2.4** (Consistency of the MOM estimator). *Suppose that  $X_1, X_2, \dots$  are i.i.d with a distribution indexed by a  $k$ -dimensional parameter vector  $\theta$ . If the first  $k$  moments of that distribution exist and are finite for all  $\theta$  and the inverse function  $M$  in Definition (10) is continuous, then the sequence of MOM estimators based on  $X_1, X_2, \dots$  is consistent for  $\theta$ .*

*Proof.* By the LLN, the sample moments converge in probability to the population moments  $\mu_1(\theta), \dots, \mu_k(\theta)$ . A generalization of the continuous mapping theorem to functions of  $k$  variables implies that  $M(\cdot)$  evaluated at the sample moments converges in probability to  $\theta$ , i.e., the MOM estimator converges in probability to  $\theta$ .  $\square$

**Example 2.5.** Let  $X_1, X_2, \dots, X_n$  be from a  $N(\mu, \sigma^2)$  distribution. Thus  $\theta = (\mu, \sigma^2)$ . What is the MOM estimator of  $\theta$ ?

**Solution:** Because  $\mu_1 = \mathbb{E}(X_1) = \mu$  and  $\mu_2 = \mathbb{E}(X_1^2) = \sigma^2 + \mu^2$ , it is easy to express the unknown parameters  $\mu$  and  $\sigma^2$  in terms of  $\mu_1$  and  $\mu_2$ :

$$\mu = \mu_1, \quad \sigma^2 = \mu_2 - \mu^2 = \mu_2 - \mu_1^2.$$

To get MOM estimates of  $\mu$  and  $\sigma^2$  we are going to plug in the sample moments. Thus

$$\hat{\mu} = \hat{\mu}_1 = \bar{X},$$

and

$$\hat{\sigma}^2 = \hat{\mu}_2 - \hat{\mu}_1^2 = \frac{1}{n} \sum_{j=1}^n X_j^2 - \bar{X}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

where we have used the fact that

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

**Example 2.6.** Suppose that  $X_1, X_2, \dots, X_n$  are i.i.d Gamma( $\alpha, \beta$ ),  $\alpha, \beta > 0$ . Thus,  $\theta = (\alpha, \beta) \in \Omega := \mathbb{R}_+ \times \mathbb{R}_+$ . The first two moments of this distribution are:

$$\mu_1(\theta) = \frac{\alpha}{\beta}, \quad \mu_2(\theta) = \frac{\alpha(\alpha + 1)}{\beta^2},$$

which implies that

$$\alpha = \frac{\mu_1^2}{\mu_2 - \mu_1^2}, \quad \beta = \frac{\mu_1}{\mu_2 - \mu_1^2}.$$

The MOM says that we replace the right-hand sides of these equations by the *sample moments*. In this case, we get

$$\hat{\alpha} = \frac{\hat{\mu}_1^2}{\hat{\mu}_2 - \hat{\mu}_1^2}, \quad \hat{\beta} = \frac{\hat{\mu}_1}{\hat{\mu}_2 - \hat{\mu}_1^2}.$$

**Remark:** MOM can thus be thought of as “plug-in” estimates; to get an estimate  $\hat{\theta}$  of  $\theta = M(\mu_1, \mu_2, \dots, \mu_k)$ , we plug-in estimates of the  $\mu_i$ ’s, which are the  $\hat{\mu}_i$ ’s, to get  $\hat{\theta}$ .

In general, we might be interested in estimating  $\Psi(\theta)$  where  $\Psi(\theta)$  is some (known) function of  $\theta$ ; in such a case, the MOM estimate of  $\Psi(\theta)$  is  $\Psi(\hat{\theta})$  where  $\hat{\theta}$  is the MOM estimate of  $\theta$ .

**Example 2.7.** Let  $X_1, X_2, \dots, X_n$  be the indicators of  $n$  Bernoulli trials with success probability  $\theta$ . We are going to find a MOM estimator of  $\theta$ .

**Solution:** Note that  $\theta$  is the probability of success and satisfies,

$$\theta = \mathbb{E}(X_1), \quad \theta = \mathbb{E}(X_1^2).$$

Thus we can get MOMs of  $\theta$  based on both the first and the second moments. Thus,

$$\hat{\theta} = \bar{X}$$

or

$$\hat{\theta} = \frac{1}{n} \sum_{j=1}^n X_j^2 = \frac{1}{n} \sum_{j=1}^n X_j = \bar{X}.$$

Here, the MOM estimate based on the second moment  $\mu_2$  coincides with the MOM estimate based on  $\mu_1$ . However, this is not necessarily the case; the MOM estimate of a certain parameter *may not be unique* as illustrated by the following example.

**Example 2.8.** Let  $X_1, X_2, \dots, X_n$  be i.i.d.  $\text{Poisson}(\lambda)$ ,  $\lambda > 0$ . Find the MOM estimator of  $\theta := \lambda + \lambda^2$ .

**Solution:** On the one hand, because  $\mu_1 = \mathbb{E}(X_1) = \lambda$ , the MOM estimate of  $\theta$  based on the first moment is

$$\hat{\theta} = \hat{\mu}_1 + \hat{\mu}_1^2 = \bar{X} + \bar{X}^2.$$

On the other hand,  $\mu_2 = \mathbb{E}(X_1^2) = \text{Var}(X_1) + \mathbb{E}(X_1)^2 = \lambda + \lambda^2 = \theta$ , so the MOM estimate of  $\theta$  based on the second moment is

$$\hat{\theta} = \frac{1}{n} \sum_{j=1}^n X_j^2.$$

However, these two estimates are not necessarily equal; in other words, it is not necessarily the case that  $\bar{X}^2 + \bar{X} = \frac{1}{n} \sum_{j=1}^n X_j^2$ .

This illustrates one of the disadvantages of MOM estimates—they may not be uniquely defined.

**Example 2.9.** Consider  $n$  systems with failure times  $X_1, X_2, \dots, X_n$  assumed to be i.i.d  $\text{Exp}(\lambda)$ ,  $\lambda > 0$ . Find the MOM estimators of  $\lambda$ .

**Solution:** It is not difficult to show that

$$\mathbb{E}(X_1) = \frac{1}{\lambda}, \quad \mathbb{E}(X_1^2) = \frac{2}{\lambda^2}.$$

Therefore

$$\lambda = \frac{1}{\mu_1} = \sqrt{\frac{2}{\mu_2}}.$$

The above equations lead to two different MOM estimators for  $\lambda$ ; the estimate based on the first moment is

$$\hat{\lambda} = \frac{1}{\hat{\mu}_1},$$

and the estimate based on the second moment is

$$\hat{\lambda} = \sqrt{\frac{2}{\hat{\mu}_2}}.$$

Once again, note the non-uniqueness of the estimates.

We finish up this section by some key observations about method of moments estimates.

- (i) The MOM principle generally leads to procedures that are easy to compute and which are therefore valuable as preliminary estimates.
- (ii) For large sample sizes, these estimates are likely to be close to the value being estimated (consistency).
- (iii) The prime disadvantage is that they do not provide a unique estimate and this has been illustrated before with examples.

### 3 Method of Maximum Likelihood

As before, we have i.i.d observations  $X_1, X_2, \dots, X_n$  with common probability density (or mass function)  $f(x, \theta)$ , where  $\theta \in \Omega \subseteq \mathbb{R}^k$  is a Euclidean parameter indexing the class of distributions being considered.

The goal is to estimate  $\theta$  or some  $\Psi(\theta)$  where  $\Psi$  is some known function of  $\theta$ .

**Definition 11** (Likelihood function). *The likelihood function for the sample  $\mathbf{X}_n = (X_1, X_2, \dots, X_n)$  is*

$$L_n(\theta) \equiv L_n(\theta, \mathbf{X}_n) := \prod_{i=1}^n f(X_i, \theta).$$

*This is simply the joint density (or mass function) but we now think of this as a function of  $\theta$  for a fixed  $\mathbf{X}_n$ ; namely the  $\mathbf{X}_n$  that is realized.*

**Intuition:** Suppose for the moment that  $X_i$ 's are discrete, so that  $f$  is actually a p.m.f. Then  $L_n(\theta)$  is exactly the probability that the observed data is realized or “happens”.

We now seek to obtain that  $\theta \in \Omega$  for which  $L_n(\theta)$  is maximized. Call this  $\hat{\theta}_n$  (assume that it exists). Thus  $\hat{\theta}_n$  is that value of the parameter that maximizes the likelihood function, or in other words, makes the observed data most likely.

It makes sense to pick  $\hat{\theta}_n$  as a guess for  $\theta$ .

When the  $X_i$ 's are continuous and  $f(x, \theta)$  is in fact a density we do the same thing – maximize the likelihood function as before and prescribe the maximizer as an estimate of  $\theta$ .

For obvious reasons,  $\hat{\theta}_n$  is called an **maximum likelihood estimate** (MLE).

**Remarks:**

- Note that  $\hat{\theta}_n$  is itself a deterministic function of  $\mathbf{X}_n = (X_1, X_2, \dots, X_n)$  and is therefore a random variable. Of course there is nothing that guarantees that  $\hat{\theta}_n$  is unique, even if it exists.
- Sometimes, in the case of multiple maximizers, we choose one which is more desirable according to some “sensible” criterion.

**Example 3.1.** Suppose that  $X_1, \dots, X_n$  are i.i.d Poisson( $\theta$ ),  $\theta > 0$ . Find the MLE of  $\theta$ .

**Solution:** In this case, it is easy to see that

$$L_n(\theta) = \prod_{i=1}^n \frac{e^{-\theta} \theta^{X_i}}{X_i!} = C(\mathbf{X}_n) e^{-n\theta} \theta^{\sum_{i=1}^n X_i}.$$



To maximize this expression, we set

$$\frac{\partial}{\partial \theta} \log L_n(\theta) = 0.$$

This yields that

$$\frac{\partial}{\partial \theta} \left[ -n\theta + \left( \sum_{i=1}^n X_i \right) \log \theta \right] = 0;$$

i.e.,

$$-n + \frac{\sum_{i=1}^n X_i}{\theta} = 0,$$

showing that

$$\hat{\theta}_n = \bar{X}.$$

It can be checked (by computing the second derivative at  $\hat{\theta}_n$ ) that the stationary point indeed gives (a unique) maximum (or by noting that the log-likelihood is a (strictly) concave function).

**Exercise 2:** Let  $X_1, X_2, \dots, X_n$  be i.i.d  $\text{Ber}(\theta)$  where  $0 \leq \theta \leq 1$ . What is the MLE of  $\theta$ ?

**Example 3.2.** Suppose  $X_1, X_2, \dots, X_n$  are i.i.d  $\text{Uniform}([0, \theta])$  random variables, where  $\theta > 0$ . We want to obtain the MLE of  $\theta$ .

**Solution:** The likelihood function is given by,

$$\begin{aligned} L_n(\theta) &= \prod_{i=1}^n \frac{1}{\theta} I_{[0, \theta]}(X_i) \\ &= \frac{1}{\theta^n} \prod_{i=1}^n I_{[X_i, \infty)}(\theta) \\ &= \frac{1}{\theta^n} I_{[\max_{i=1, \dots, n} X_i, \infty)}(\theta). \end{aligned}$$

It is then clear that  $L_n(\theta)$  is constant and equals  $1/\theta^n$  for  $\theta \geq \max_{i=1, \dots, n} X_i$  and is 0 otherwise. By plotting the graph of this function, you can see that

$$\hat{\theta}_n = \max_{i=1, \dots, n} X_i.$$

Here, differentiation will not help you to get the MLE because the likelihood function is not differentiable at the point where it hits the maximum.

**Example 3.3.** Suppose that  $X_1, X_2, \dots, X_n$  are i.i.d  $N(\mu, \sigma^2)$ . We want to find the MLEs of the mean  $\mu$  and the variance  $\sigma^2$ .

**Solution:** We write down the likelihood function first. This is,

$$L_n(\mu, \sigma^2) = \frac{1}{(2\pi)^{n/2}\sigma^n} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2\right).$$

It is easy to see that,

$$\begin{aligned} \log L_n(\mu, \sigma^2) &= -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 + \text{constant} \\ &= -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 - \frac{n}{2\sigma^2} (\bar{X}_n - \mu)^2. \end{aligned}$$

To maximize the above expression w.r.t  $\mu$  and  $\sigma^2$  we proceed as follows. For any  $(\mu, \sigma^2)$  we have,

$$\log L_n(\mu, \sigma^2) \leq \log L_n(\bar{X}_n, \sigma^2),$$

showing that we can choose  $\hat{\mu}_{MLE} = \bar{X}_n$ .

It then remains to maximize  $\log L_n(\bar{X}_n, \sigma^2)$  with respect to  $\sigma^2$  to find  $\hat{\sigma}_{MLE}^2$ .

Now,

$$\log L_n(\bar{X}_n, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Differentiating the left-side w.r.t  $\sigma^2$  gives,

$$-\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} n \hat{\sigma}^2 = 0,$$

where  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ . The above equation leads to,

$$\hat{\sigma}_{MLE}^2 = \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

The fact that this actually gives a global maximizer follows from the fact that the second derivative at  $\hat{\sigma}^2$  is negative.

Note that, once again, the MOM estimates coincide with the MLEs.

**Exercise 3:** We now tweak the above situation a bit. Suppose now that we restrict the parameter space, so that  $\mu$  has to be non-negative, i.e.,  $\mu \geq 0$ .

Thus we seek to maximize  $\log L_n(\mu, \sigma^2)$  but subject to the constraint that  $\mu \geq 0$  and  $\sigma^2 > 0$ . Find the MLEs in this scenario.

**Example 3.4** (MLEs might not be unique). Suppose that  $X_1, \dots, X_n$  form a random sample from the uniform distribution on the interval  $[\theta, \theta+1]$ , where  $\theta \in \mathbb{R}$  is unknown. Show that the MLE of  $\theta$  is not unique.

**Solution:** The likelihood has the form

$$L_n(\theta) = \prod_{i=1}^n I_{[\theta, \theta+1]}(X_i).$$

The condition that  $\theta \leq X_i$ , for all  $i = 1, \dots, n$ , is equivalent to the condition that  $\theta \leq \min\{X_1, \dots, X_n\} = X_{(1)}$ . Similarly, the condition that  $X_i \leq \theta + 1$ , for all  $i = 1, \dots, n$ , is equivalent to the condition that  $\theta \geq \max\{X_1, \dots, X_n\} - 1 = X_{(n)} - 1$ . Thus the likelihood can be written as

$$L_n(\theta) = I_{[X_{(n)}-1, X_{(1)}]}(\theta).$$

Hence it is possible to select as an MLE any value of  $\theta$  in the interval  $[X_{(n)} - 1, X_{(1)}]$ , and thus the MLE is not unique.

**Example 3.5** (MLEs might not exist). Consider a random variable  $X$  that can come with equal probability either from a  $N(0, 1)$  or from  $N(\mu, \sigma^2)$ , where both  $\mu$  and  $\sigma$  are unknown.

Thus, the p.d.f.  $f(\cdot, \mu, \sigma^2)$  of  $X$  is given by

$$f(x, \mu, \sigma^2) = \frac{1}{2} \left[ \frac{1}{\sqrt{2\pi}} e^{-x^2/2} + \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/(2\sigma^2)} \right].$$

Suppose now that  $X_1, \dots, X_n$  form an i.i.d sample from this distribution. As usual, the likelihood function

$$L_n(\mu, \sigma^2) = \prod_{i=1}^n f(X_i, \mu, \sigma^2).$$

We want to find the MLE of  $\theta = (\mu, \sigma^2)$ .

Let  $X_k$  denote one of the observed values. Note that

$$\max_{\mu \in \mathbb{R}, \sigma^2 > 0} L_n(\mu, \sigma^2) \geq L_n(X_k, \sigma^2) \geq \frac{1}{2^n} \left[ \frac{1}{\sqrt{2\pi}\sigma} \right] \prod_{i \neq k} \frac{1}{\sqrt{2\pi}} e^{-X_i^2/2}.$$

Thus, if we let  $\mu = X_k$  and let  $\sigma^2 \rightarrow 0$  then the factor  $f(X_k, \mu, \sigma^2)$  will grow large without bound, while each factor  $f(X_i, \mu, \sigma^2)$ , for  $i \neq k$ , will approach the value

$$\frac{1}{2\sqrt{2\pi}} e^{-X_i^2/2}.$$

Hence, when  $\mu = X_k$  and  $\sigma^2 \rightarrow 0$ , we find that  $L_n(\mu, \sigma^2) \rightarrow \infty$ .

Note that 0 is not a permissible estimate of  $\sigma^2$ , because we know in advance that  $\sigma > 0$ . Since the likelihood function can be made arbitrarily large by choosing  $\mu = X_k$  and choosing  $\sigma^2$  arbitrarily close to 0, it follows that the MLE *does not exist*.

### 3.1 Properties of MLEs

**Theorem 3.6** (Invariance property of MLEs). *If  $\hat{\theta}_n$  is the MLE of  $\theta$  and if  $\Psi$  is any function, then  $\Psi(\hat{\theta}_n)$  is the MLE of  $\Psi(\theta)$ .*

See Theorem 7.6.2 and Example 7.6.3 in the text book.

Thus if  $X_1, \dots, X_n$  be i.i.d  $N(\mu, \sigma^2)$ , then the MLE of  $\sigma$  is  $\sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$ .

Consider an estimation problem in which a random sample is to be taken from a distribution involving a parameter  $\theta$ . Then, under certain conditions, which are typically satisfied in practical problems, the sequence of MLEs is *consistent*, i.e.,

$$\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta, \quad \text{as } n \rightarrow \infty.$$

### 3.2 Computational methods for approximating MLEs

**Example:** Suppose that  $X_1, \dots, X_n$  are i.i.d from a Gamma distribution for which the p.d.f is as follows:

$$f(x, \alpha) = \frac{1}{\Gamma(\alpha)} x^{\alpha-1} e^{-x}, \quad \text{for } x > 0.$$

The likelihood function is

$$L_n(\alpha) = \frac{1}{\Gamma(\alpha)^n} \left( \prod_{i=1}^n X_i \right)^{\alpha-1} e^{-\sum_{i=1}^n X_i},$$

and thus the log-likelihood is

$$\ell_n(\alpha) \equiv \log L_n(\alpha) = -n \log \Gamma(\alpha) + (\alpha - 1) \sum_{i=1}^n \log(X_i) - \sum_{i=1}^n X_i,$$

The MLE of  $\alpha$  will be the value of  $\alpha$  that satisfies the equation

$$\begin{aligned} \frac{\partial}{\partial \alpha} \ell_n(\alpha) &= -n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_{i=1}^n \log(X_i) = 0 \\ \text{i.e., } \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} &= \frac{1}{n} \sum_{i=1}^n \log(X_i). \end{aligned}$$

Solving this equation for  $\alpha$  cannot be done analytically and requires computational methods to obtain an approximate solutions.

Read pages 428–430 and pages 434–439 of the text-book. I will cover this later, if time permits.

## 4 Principles of estimation

Setup: Our data  $X_1, X_2, \dots, X_n$  are i.i.d observations from the distribution  $P_\theta$  where  $\theta \in \Omega$ , the parameter space ( $\Omega$  is assumed to be the  $k$ -dimensional Euclidean space). We assume identifiability of the parameter, i.e.  $\theta_1 \neq \theta_2 \Rightarrow P_{\theta_1} \neq P_{\theta_2}$ .

**Estimation problem:** Consider now, the problem of estimating  $g(\theta)$  where  $g$  is some function of  $\theta$ .

In many cases  $g(\theta) = \theta$  itself.

Generally  $g(\theta)$  will describe some important aspect of the distribution  $P_\theta$ .

Our estimator of  $g(\theta)$  will be some function of our observed data  $\mathbf{X}_n = (X_1, X_2, \dots, X_n)$ .

In general there will be several different estimators of  $g(\theta)$  which may all seem reasonable from different perspectives — the question then becomes one of finding the most optimal one.

This requires an objective **measure of performance** of an estimator.

If  $T_n$  estimates  $g(\theta)$  a criterion that naturally suggests itself is the distance of  $T_n$  from  $g(\theta)$ . Good estimators are those for which  $|T_n - g(\theta)|$  is generally small.

Since  $T_n$  is a random variable no deterministic statement can be made about the *absolute deviation*; however what we can expect of a good estimator is a high chance of remaining close to  $g(\theta)$ .

Also as  $n$ , the sample size, increases we get hold of more information and hence expect to be able to do a better job of estimating  $g(\theta)$ .

These notions when coupled together give rise to the **consistency** requirement for a sequence of estimators  $T_n$ ; as  $n$  increases,  $T_n$  ought to converge in probability to  $g(\theta)$  (under the probability distribution  $P_\theta$ ). In other words, for any  $\epsilon > 0$ ,

$$\mathbb{P}_\theta (|T_n - g(\theta)| > \epsilon) \rightarrow 0.$$

The above is clearly a *large sample property*; what it says is that with probability increasing to 1 (as the sample size grows),  $T_n$  estimates  $g(\theta)$  to any pre-determined level of accuracy.

However, the consistency condition alone, does not tell us anything about how well we are performing for any particular sample size, or the rate at which the above probability is going to 0.

### 4.1 Mean squared error

**Question:** For a fixed sample size  $n$ , how do we measure the performance of an

estimator  $T_n$ ?

A way out of this difficulty is to obtain an average measure of the error, or in other words, average out  $|T_n - g(\theta)|$  over all possible realizations of  $T_n$ .

The resulting quantity is then still a function of  $\theta$  but no longer random. It is called the **mean absolute error** and can be written compactly (using acronym) as:

$$\text{MAD} := \mathbb{E}_\theta [|T_n - g(\theta)|] .$$

However, it is more common to avoid absolute deviations and work with the square of the deviation, integrated out as before over the distribution of  $T_n$ . This is called the **mean squared error** (MSE) and is defined as

$$\text{MSE}(T_n, g(\theta)) := \mathbb{E}_\theta [(T_n - g(\theta))^2] . \quad (1)$$

Of course, this is meaningful, only if the above quantity is finite for all  $\theta$ . Good estimators are those for which the MSE is generally not too high, whatever be the value of  $\theta$ .

There is a standard decomposition of the MSE that helps us understand its components.

**Theorem 4.1.** *For any estimator  $T_n$  of  $g(\theta)$ , we have*

$$\text{MSE}(T_n, g(\theta)) = \text{Var}_\theta(T_n) + b(T_n, g(\theta))^2 ,$$

where  $b(T_n, g(\theta)) = \mathbb{E}_\theta(T_n) - g(\theta)$  is the **bias** of  $T_n$  as an estimator of  $g(\theta)$ .

*Proof.* We have,

$$\begin{aligned} \text{MSE}(T_n, g(\theta)) &= \mathbb{E}_\theta [(T_n - g(\theta))^2] \\ &= \mathbb{E}_\theta [(T_n - \mathbb{E}_\theta(T_n) + \mathbb{E}_\theta(T_n) - g(\theta))^2] \\ &= \mathbb{E}_\theta [(T_n - \mathbb{E}_\theta(T_n))^2] + (\mathbb{E}_\theta(T_n) - g(\theta))^2 \\ &\quad + 2 \mathbb{E}_\theta[(T_n - \mathbb{E}_\theta(T_n))(\mathbb{E}_\theta(T_n) - g(\theta))] \\ &= \text{Var}_\theta(T_n) + b(T_n, g(\theta))^2 , \end{aligned}$$

where

$$b(T_n, g(\theta)) := \mathbb{E}_\theta(T_n) - g(\theta)$$

is the **bias** of  $T_n$  as an estimator of  $g(\theta)$ .

The cross product term in the above display vanishes since  $\mathbb{E}_\theta(T_n) - g(\theta)$  is a constant and  $\mathbb{E}_\theta(T_n - \mathbb{E}_\theta(T_n)) = 0$ .  $\square$

The bias measures, on an average, by how much  $T_n$  overestimates or underestimates  $g(\theta)$ . If we think of the expectation  $\mathbb{E}_\theta(T_n)$  as the center of the distribution of  $T_n$ , then the bias measures by how much the *center deviates from the target*.

The variance of  $T_n$ , of course, measures how closely  $T_n$  is clustered around its center. Ideally one would like to minimize both simultaneously, but unfortunately this is rarely possible.

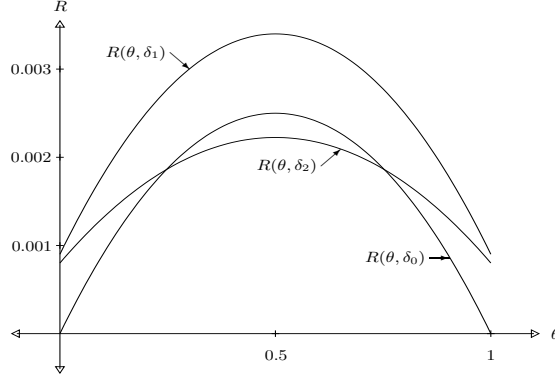


Figure 3: The plot shows the mean squared error for three estimators  $\delta_1$ ,  $\delta_2$  and  $\delta_0$ . Here  $R(\theta, \delta_i) = \mathbb{E}_\theta[(\delta_i(X) - \theta)^2]$  where  $i = 0, 1, 2$ .

## 4.2 Comparing estimators

Two estimators  $T_n$  and  $S_n$  can be compared on the basis of their MSEs. Under parameter value  $\theta$ ,  $T_n$  dominates  $S_n$  as an estimator if

$$\text{MSE}(T_n, \theta) \leq \text{MSE}(S_n, \theta) \quad \text{for all } \theta \in \Omega.$$

In this situation we say that  $S_n$  is *inadmissible* in the presence of  $T_n$ .

The use of the term “inadmissible” hardly needs explanation. If, for all possible values of the parameter, we incur less error using  $T_n$  instead of  $S_n$  as an estimate of  $g(\theta)$ , then clearly there is no point in considering  $S_n$  as an estimator at all.

Continuing along this line of thought, is there an estimate that improves all others? In other words, is there an estimator that makes every other estimator inadmissible? The answer is **no**, except in certain pathological situations.

**Example 4.2.** Suppose that  $X \sim \text{Binomial}(100, \theta)$ , where  $\theta \in [0, 1]$ . The goal is to estimate the unknown parameter  $\theta$ . A natural estimator of  $\theta$  in this problem is  $\delta_0(X) = X/100$  (which is also the MLE and the method of moments estimator). Then

$$R(\theta, \delta_0) := \text{MSE}(\delta_0(X), \theta) = \frac{\theta(1 - \theta)}{100}, \quad \text{for } \theta \in [0, 1].$$

The MSE of  $\delta_0(X)$  as a function of  $\theta$  is given in Figure 3.

We can also consider two other estimators in this problem:  $\delta_1(X) = (X + 3)/100$  and  $\delta_2(X) = (X + 3)/106$ . Figure 3 shows the MSEs of  $\delta_1$  and  $\delta_2$ , which can be shown to be (show this):

$$R(\theta, \delta_1) := \text{MSE}(\delta_1(X), \theta) = \frac{9 + 100\theta(1 - \theta)}{100^2}, \quad \text{for } \theta \in [0, 1],$$

and

$$R(\theta, \delta_2) := \text{MSE}(\delta_2(X), \theta) = \frac{(9 - 8\theta)(1 + 8\theta)}{106^2}, \quad \text{for } \theta \in [0, 1].$$

Looking at the plot,  $\delta_0$  and  $\delta_2$  are both better than  $\delta_1$ , but the comparison between  $\delta_0$  and  $\delta_2$  is ambiguous. When  $\theta$  is near  $1/2$ ,  $\delta_2$  is the preferable estimator, but if  $\theta$  is near 0 or 1,  $\delta_0$  is preferable. If  $\theta$  were known, we could choose between  $\delta_0$  and  $\delta_2$ . However, if  $\theta$  were known, there would be no need to estimate its value.

As we have noted before, it is generally not possible to find a universally best estimator. One way to try to construct optimal estimators is to restrict oneself to a subclass of estimators and try to find the best possible estimator in this subclass. One arrives at subclasses of estimators by constraining them to meet some desirable requirements. One such requirement is that of *unbiasedness*. Below, we provide a formal definition.

### 4.3 Unbiased estimators

An estimator  $T_n$  of  $g(\theta)$  is said to be *unbiased* if  $\mathbb{E}_\theta(T_n) = g(\theta)$  for all possible values of  $\theta$ ; i.e.,

$$b(T_n, g(\theta)) = 0 \quad \text{for all } \theta \in \Omega.$$

Thus, unbiased estimators, on an average, hit the target, for all parameter values. This seems to be a reasonable constraint to impose on an estimator and indeed produces meaningful estimates in a variety of situations.

**Lemma 4.3.** *Let  $X_1, \dots, X_n$  be iid with mean  $\mu$  and finite variance  $\sigma^2$ . Then*

$$\mathbb{E}[\bar{X}_n] = \mu, \quad \mathbb{E}[s_n^2] = \sigma^2.$$

*In other words, sample mean and sample variance are unbiased estimators for  $\mu$  and  $\sigma^2$ , respectively.*

*Proof.* Exercise! In particular, make sure you see in the calculation why we have to define  $s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  (with division by  $n-1$ ) to obtain an unbiased estimator.  $\square$

Note that for an unbiased estimator  $T_n$ , the MSE under  $\theta$  is simply the variance of  $T_n$  under  $\theta$ .

In a large class of models, it is possible to find an unbiased estimator of  $g(\theta)$  that has the smallest possible variance among all possible unbiased estimators. Such an estimate is called an **minimum variance unbiased estimator** (MVUE). Here is a formal definition.

**MVUE:** We call  $S_n$  an MVUE of  $g(\theta)$  if

$$(i) \quad \mathbb{E}_\theta(S_n) = g(\theta) \quad \text{for all } \theta \in \Omega$$

and (ii) if  $T_n$  is an unbiased estimate of  $g(\theta)$ , then  $\text{Var}_\theta(S_n) \leq \text{Var}_\theta(T_n)$ .

Here are a few examples to illustrate some of the various concepts discussed above.



- (a) Consider  $X_1, \dots, X_n$  i.i.d  $N(\mu, \sigma^2)$ .

A natural unbiased estimator of  $g_1(\theta) = \mu$  is  $\bar{X}_n$ , the sample mean. It is also consistent for  $\mu$  by the WLLN. It can be shown that this is also the MVUE of  $\mu$ .

In other words, *any* other unbiased estimate of  $\mu$  will have a larger variance than  $\bar{X}_n$ . Recall that the variance of  $\bar{X}_n$  is simply  $\sigma^2/n$ .

Consider now, the estimation of  $\sigma^2$ . Two estimates of this that we have considered in the past are

$$(i) \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \quad \text{and} \quad (ii) s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Out of these  $\hat{\sigma}_n^2$  is not unbiased for  $\sigma^2$  but  $s_n^2$  is. In fact  $s_n^2$  is the MVUE of  $\sigma^2$ .

- (b) Let  $X_1, X_2, \dots, X_n$  be i.i.d from some underlying density function or mass function  $f(x, \theta)$ . Let  $g(\theta) = \mathbb{E}_\theta(X_1)$ .

Then the sample mean  $\bar{X}_n$  is always an unbiased estimate of  $g(\theta)$  (see Lemma 4.3). Whether it is MVUE or not depends on the underlying structure of the model.

- (c) Suppose that  $X_1, X_2, \dots, X_n$  be i.i.d Ber( $\theta$ ). It can be shown that  $\bar{X}_n$  is the MVUE of  $\theta$ .

Now define  $g(\theta) = \theta/(1-\theta)$ . This is a quantity of interest because it is precisely the odds in favor of Heads. It can be shown that there is *no unbiased estimator* of  $g(\theta)$  in this model (**Why?**).

However an intuitively appealing estimate of  $g(\theta)$  is  $T_n \equiv \bar{X}_n/(1 - \bar{X}_n)$ . It is *not unbiased* for  $g(\theta)$ ; however it does converge in probability to  $g(\theta)$ .

This example illustrates an important point — unbiased estimators may not always exist. Hence imposing unbiasedness as a constraint may not be meaningful in all situations.

- (d) Unbiased estimators are not always better than biased estimators.

Remember, it is the MSE that gauges the performance of the estimator and a biased estimator may actually outperform an unbiased one owing to a significantly smaller variance.

**Example 4.4.** Consider  $X_1, X_2, \dots, X_n$  i.i.d Uniform( $[0, \theta]$ ) with  $\theta > 0$ . Here  $\Omega = (0, \infty)$ . The MLE for  $\theta$  is the maximum of the  $X_i$ 's, which we denote by  $X_{(n)}$ . Another estimate of  $\theta$  is obtained by observing that  $\bar{X}_n$  is an unbiased estimate of  $\theta/2$ , the common mean of the  $X_i$ 's; hence  $2\bar{X}_n$  is an unbiased estimate of  $\theta$ . Show that  $X_{(n)}$  in the sense of MSE outperforms  $2\bar{X}_n$  by an order of magnitude. The best unbiased estimator (MVUE) of  $\theta$  is  $(1 + n^{-1})X_{(n)}$ .

## 4.4 Sufficient Statistics

In some problems, there may not be any MLE, or there may be more than one. Even when an MLE is unique, it may not be a suitable estimator (as in the  $\text{Unif}(0, \theta)$  example, where the MLE always underestimates the value of  $\theta$ ).

In such problems, the search for a good estimator must be extended beyond the methods that have been introduced thus far.

In this section, we shall define the concept of a **sufficient statistic**, which can be used to simplify the search for a good estimator in many problems.

Suppose that in a specific estimation problem, two statisticians A and B must estimate the value of the parameter  $\theta$ .

Statistician A can observe the values of the observations  $X_1, X_2, \dots, X_n$  in a random sample, and statistician B cannot observe the individual values of  $X_1, X_2, \dots, X_n$  but can learn the value of a certain statistic  $T = \varphi(X_1, \dots, X_n)$ .

In this case, statistician A can choose any function of the observations  $X_1, X_2, \dots, X_n$  as an estimator of  $\theta$  (including a function of  $T$ ). But statistician B can use only a function of  $T$ . Hence, it follows that A will generally be able to find a better estimator than will B.

In some problems, however, B will be able to do just as well as A. In such a problem, *the single function  $T = \varphi(X_1, \dots, X_n)$  will in some sense summarize all the information contained in the random sample about  $\theta$* , and knowledge of the individual values of  $X_1, \dots, X_n$  will be irrelevant in the search for a good estimator of  $\theta$ . A statistic  $T$  having this property is called a **sufficient statistic**.

A statistic is **sufficient** with respect to a statistical model  $P_\theta$  and its associated unknown parameter  $\theta$  if it provides “all” the information on  $\theta$ ; e.g., if “no other statistic that can be calculated from the same sample provides any additional information as to the value of the parameter”. This intuition will be made rigorous at the end of this subsection.

**Definition 12** (Sufficient statistic). *Let  $X_1, X_2, \dots, X_n$  be a random sample from a distribution indexed by a parameter  $\theta \in \Omega$ . A statistic  $T$  is called a **sufficient statistic** for the parameter  $\theta$  if the following holds: For all possible values of  $t$  and no matter what the true value of  $\theta \in \Omega$  is, the conditional joint distribution of  $X_1, X_2, \dots, X_n$  given that  $T = t$  does not depend on  $\theta$ .*

So, if  $T$  is sufficient, and one observed only  $T$  instead of  $(X_1, \dots, X_n)$ , one could, at least in principle, simulate random variables  $(X'_1, \dots, X'_n)$  with the same joint distribution.

In this sense,  $T$  is sufficient for obtaining as much information about  $\theta$  as one could get from  $(X_1, \dots, X_n)$ .

**Example 4.5.** Suppose that  $X_1, \dots, X_n$  are i.i.d Poisson( $\theta$ ), where  $\theta > 0$ . Show that  $T = \sum_{i=1}^n X_i$  is sufficient. Let  $\mathbf{X} = (X_1, \dots, X_n)$ .

Note that

$$\mathbb{P}_\theta(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = t) = \frac{\mathbb{P}_\theta(\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = t)}{\mathbb{P}_\theta(T(\mathbf{X}) = t)}.$$

But,

$$\mathbb{P}_\theta(\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = t) = \begin{cases} 0 & T(\mathbf{x}) \neq t \\ \mathbb{P}_\theta(\mathbf{X} = \mathbf{x}) & T(\mathbf{x}) = t. \end{cases}$$

Because

$$\mathbb{P}(\mathbf{X} = \mathbf{x}) = \frac{e^{-n\theta} \theta^{T(\mathbf{x})}}{\prod_{i=1}^n x_i!}$$

and

$$\mathbb{P}_\theta(T(\mathbf{X}) = t) = \frac{e^{-n\theta} (n\theta)^t}{t!},$$

we have that

$$\frac{\mathbb{P}_\theta(\mathbf{X} = \mathbf{x})}{\mathbb{P}_\theta(T(\mathbf{X}) = t)} = \frac{t!}{\prod_{i=1}^n x_i! n^t},$$

which does not depend on  $\theta$ . So  $T = \sum_{i=1}^n X_i$  is a sufficient statistic for  $\theta$ .

Other sufficient statistics are:  $T = 3.7 \sum_{i=1}^n X_i$ ,  $T = (\sum_{i=1}^n X_i, X_4)$ , and  $T = (X_1, \dots, X_n)$ .

We shall now present a simple method for finding a sufficient statistic that can be applied in many problems.

**Theorem 4.6** (Factorization criterion). *Let  $X_1, X_2, \dots, X_n$  be iid from either a continuous distribution or a discrete distribution for which the p.d.f or the p.m.f is  $f(x, \theta)$ , where the value of  $\theta$  is unknown and belongs to a given parameter space  $\Omega$ . A statistic  $T = r(X_1, X_2, \dots, X_n)$  is a sufficient statistic for  $\theta$  if and only if the joint p.d.f or the joint p.m.f  $f_n(\mathbf{x}, \theta)$  of  $(X_1, X_2, \dots, X_n)$  can be factored as follows for all values of  $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$  and all values of  $\theta \in \Omega$ :*

$$f_n(\mathbf{x}, \theta) = u(\mathbf{x}) \nu(r(\mathbf{x}), \theta),$$

where

- $u$  and  $\nu$  are both non-negative,
- the function  $u$  may depend on  $\mathbf{x}$  but does not depend on  $\theta$ ,
- the function  $\nu$  will depend on  $\theta$  but depends on the observed value  $\mathbf{x}$  only through the value of the statistic  $r(\mathbf{x})$ .

**Example:** Suppose that  $X_1, \dots, X_n$  are i.i.d  $\text{Poi}(\theta)$ ,  $\theta > 0$ . Thus, for every non-negative integers  $x_1, \dots, x_n$ , the joint p.m.f  $f_n(\mathbf{x}, \theta)$  of  $(X_1, \dots, X_n)$  is

$$f_n(\mathbf{x}, \theta) = \prod_{i=1}^n \frac{e^{-\theta} \theta^{x_i}}{x_i!} = \frac{1}{\prod_{i=1}^n x_i!} e^{-n\theta} \theta^{\sum_{i=1}^n x_i}.$$

Thus, we can take  $u(\mathbf{x}) = 1/(\prod_{i=1}^n x_i!)$ ,  $r(\mathbf{x}) = \sum_{i=1}^n x_i$ ,  $\nu(t, \theta) = e^{-n\theta} \theta^t$ . It follows that  $T = \sum_{i=1}^n X_i$  is a sufficient statistic for  $\theta$ .

**Example:** Suppose that  $X_1, \dots, X_n$  are i.i.d  $\text{Gamma}(\alpha, \beta)$ ,  $\alpha, \beta > 0$ , where  $\alpha$  is known, and  $\beta$  is unknown. The joint p.d.f is

$$f_n(\mathbf{x}, \beta) = \left\{ [\Gamma(\alpha)]^n \left( \prod_{i=1}^n x_i \right)^{\alpha-1} \right\}^{-1} \times \left\{ \beta^{n\alpha} \exp(-\beta t) \right\}, \quad \text{where } t = \sum_{i=1}^n x_i.$$

$u(\mathbf{x})$   $\nu(t, \beta)$

The sufficient statistics is  $T_n = \sum_{i=1}^n X_i$ .

**Example:** Suppose that  $X_1, \dots, X_n$  are i.i.d  $\text{Gamma}(\alpha, \beta)$ ,  $\alpha, \beta > 0$ , where  $\alpha$  is unknown, and  $\beta$  is known.

The joint p.d.f in this exercise is the same as that given in the previous exercise. However, since the unknown parameter is now  $\alpha$  instead of  $\beta$ , the appropriate factorization is now

$$f_n(\mathbf{x}, \alpha) = \left\{ \exp \left( -\beta \sum_{i=1}^n x_i \right) \right\} \times \left\{ \frac{\beta^{n\alpha}}{[\Gamma(\alpha)]^n} t^{\alpha-1} \right\}, \quad \text{where } t = \sum_{i=1}^n x_i.$$

$u(\mathbf{x})$   $\nu(t, \alpha)$

The sufficient statistics is  $T_n = \sum_{i=1}^n X_i$ .

**Exercise:** Suppose that  $X_1, \dots, X_n$  are i.i.d  $\text{Unif}([0, \theta])$ ,  $\theta > 0$  is the unknown parameter. Show that  $T = \max\{X_1, \dots, X_n\}$  is a sufficient statistic.

Suppose that  $\mathbf{X} = (X_1, \dots, X_n)$  form a random sample from a distribution for which the p.d.f or p.m.f. is  $f(\cdot|\theta)$ , where the parameter  $\theta$  must belong to some parameter space  $\Omega$ . Let  $\mathbf{T}$  be a sufficient statistic for  $\theta$  in this problem.

We show how to improve upon an estimator that is not a function of a sufficient statistic by using an estimator that is a function of a sufficient statistic. Let  $\delta(\mathbf{X})$  be an estimator of  $g(\theta)$ . We define the estimator  $\delta_0(T)$  by the following conditional expectation:

$$\delta_0(\mathbf{T}) = \mathbb{E}_\theta[\delta(\mathbf{X})|\mathbf{T}].$$

Since  $\mathbf{T}$  is a sufficient statistic, the conditional expectation of the function  $\delta(\mathbf{X})$  will be the same for every value of  $\theta \in \Omega$ . It follows that the conditional expectation above will depend on the value of  $\mathbf{T}$  but will not actually depend on the value of  $\theta$ . In other words, the function  $\delta_0(\mathbf{T})$  is indeed an estimator of  $g(\theta)$  because it depends only on the observations  $\mathbf{X}$  and does not depend on the unknown value of  $\theta$ .

We can now state the following theorem, which was established independently by D. Blackwell and C. R. Rao in the late 1940s.

**Theorem 4.7** (Rao–Blackwell theorem). *For every value of  $\theta \in \Omega$ ,*

$$\text{MSE}(\delta_0(\mathbf{T}), g(\theta)) \leq \text{MSE}(\delta(\mathbf{X}), g(\theta)).$$

The above result is proved in Theorem 7.9.1 of the textbook.

## 5 The sampling distribution of a statistic

A **statistic** is a function of the data, and hence is itself a random variable with a distribution. This distribution is called its **sampling distribution**. It tells us what values the statistic is likely to assume and how likely is it to take these values. Formally, suppose that  $X_1, \dots, X_n$  are i.i.d with p.d.f/p.m.f  $f_\theta(\cdot)$ , where  $\theta \in \Omega \subset \mathbb{R}^k$ . Let  $T$  be a statistic, i.e., suppose that  $T = \varphi(X_1, \dots, X_n)$ . The distribution of  $T$  (with  $\theta$  fixed) is called the **sampling distribution** of  $T$ .

**Example:** Suppose that  $X_1, \dots, X_n$  are i.i.d  $N(\mu, \sigma^2)$ . Then we know that

$$\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right).$$

### 5.1 The gamma and the $\chi^2$ distributions

#### 5.1.1 The gamma distribution

The gamma function is a real-valued non-negative function defined on  $(0, \infty)$  in the following manner

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx, \quad \alpha > 0.$$

The Gamma function enjoys some nice properties. Two of these are listed below:

$$(a) \Gamma(\alpha + 1) = \alpha \Gamma(\alpha), \quad (b) \Gamma(n) = (n-1)! \quad (n \text{ integer}).$$

Property (b) is an easy consequence of Property (a). Start off with  $\Gamma(n)$  and use Property (a) recursively along with the fact that  $\Gamma(1) = 1$ . Another important fact is that  $\Gamma(1/2) = \sqrt{\pi}$ .

**Definition 13.** The **gamma distribution** with parameters  $\alpha > 0, \lambda > 0$  (denoted by  $\text{Gamma}(\alpha, \lambda)$ ) is defined through the following density function:

$$f(x; \alpha, \lambda) = \frac{\lambda^\alpha}{\Gamma(\alpha)} e^{-\lambda x} x^{\alpha-1} I_{(0, \infty)}(x).$$

The first parameter  $\alpha$  is called the **shape** parameter and the second parameter  $\lambda$  is called the **scale** parameter.

By definition, the  $\text{Gamma}(1, \lambda)$  is nothing else but the exponential distribution with rate  $\lambda$ . For fixed  $\lambda$  the shape parameter regulates the shape of the gamma density. Here is a simple exercise that justifies the term “scale parameter” for  $\lambda$ .

The following properties can be proved using techniques from probability theory (see Chapter 5.7 of textbook):

- **Scaling property:** Let  $X$  be a random variable following  $\text{Gamma}(\alpha, \lambda)$ . Then

$$cX \sim \text{Gamma}(\alpha, \frac{\lambda}{c}).$$

- **Reproductive property:** Let  $X_1, X_2, \dots, X_n$  be independent random variables with  $X_i \sim \text{Gamma}(\alpha_i, \lambda)$ , for  $i = 1, \dots, n$ . Then,

$$\sum_{i=1}^n X_i \sim \text{Gamma}\left(\sum_{i=1}^n \alpha_i, \lambda\right).$$

- **Moments:** If  $X$  follows the  $\text{Gamma}(\alpha, \lambda)$  distribution, then

$$\mathbb{E}(X) = \frac{\alpha}{\lambda} \quad \text{and} \quad \text{Var}(X) = \frac{\alpha}{\lambda^2}.$$

In general, if  $k$  is a positive integer,

$$\mathbb{E}(X^k) = \frac{\prod_{i=1}^k (\alpha + i - 1)}{\lambda^k}.$$

**Exercise:** Here is an exercise that should follow from the discussion above. Let  $S_n \sim \text{Gamma}(n, \lambda)$ , where  $\lambda > 0$ . Show that for large  $n$ , the distribution of  $S_n$  is well approximated by a normal distribution (with parameters that you need to identify).

### 5.1.2 The Chi-squared distribution

We now introduce an important family of distributions, called the chi-squared family. To do so, we first define the **chi-squared distribution** with 1 degree of freedom (for brevity, we call it “chi-squared one” and write it as  $\chi_1^2$ ).

**The  $\chi_1^2$  distribution:** Let  $Z \sim N(0, 1)$ . Then the distribution of  $W := Z^2$  is called the  $\chi_1^2$  distribution, and  $W$  itself is called a  $\chi_1^2$  random variable.

**Exercise:** Show that  $W$  follows a  $\text{Gamma}(1/2, 1/2)$  distribution. You can do this by relating the c.d.f. of  $W$  to that of  $Z$  and differentiation.

For any integer  $d > 0$  we can now define the  $\chi_d^2$  distribution (chi-squared  $d$  distribution, or equivalently, the chi-squared distribution with  $d$  degrees of freedom).

**The  $\chi_d^2$  distribution:** Let  $Z_1, Z_2, \dots, Z_d$  be i.i.d  $N(0, 1)$  random variables. Then the distribution of

$$W_d := Z_1^2 + Z_2^2 + \dots + Z_d^2$$

is called the  $\chi_d^2$  distribution and  $W_d$  itself is called a  $\chi_d^2$  random variable.

**Exercise:** Using the reproductive property of the Gamma distribution, show that  $W_d \sim \text{Gamma}(d/2, 1/2)$ .

**Exercise:** Let  $Z_1, Z_2, Z_3$  be i.i.d  $N(0, 1)$  random variables. Consider the vector  $(Z_1, Z_2, Z_3)$  as a random point in 3-dimensional space. Let  $R$  be the length of the radius vector connecting this point to the origin. Find the density functions of (a)  $R$  and (b)  $R^2$ .

**Theorem 5.1.** If  $X \sim \chi_m^2$ , then  $\mathbb{E}(X) = m$  and  $\text{Var}(X) = 2m$ .

**Theorem 5.2.** Suppose that  $X_1, \dots, X_k$  are independent and  $X_i \sim \chi_{m_i}^2$  then the sum

$$X_1 + \dots + X_k \sim \chi_{\sum_{i=1}^k m_i}^2.$$

In particular, the sum of  $k$  i.i.d  $\chi_1^2$  random variables is a  $\chi_k^2$  random variable.

## 5.2 Sampling from a normal population

Let  $X_1, X_2, \dots, X_n$  be i.i.d  $N(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$ ,  $\sigma > 0$  are unknown. We have seen that the MLE and MOM estimators of the mean and the variance are given by

$$\hat{\mu}_n = \bar{X}_n \quad \text{and} \quad \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Because  $\hat{\sigma}_n^2$  is biased, we will use a slightly different estimator of  $\sigma^2$  than the one proposed above. We will use the sample variance

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

Besides being an unbiased estimator (see Lemma 4.3), it turns out that  $s_n^2$  has a nice interpretation as the multiple of a  $\chi^2$  random variable. Here is an interesting (and fairly profound) proposition.

**Proposition 5.3.** Let  $X_1, X_2, \dots, X_n$  be an i.i.d sample from some distribution  $F$  with mean  $\mu$  and variance  $\sigma^2$ . Then  $F$  is the  $N(\mu, \sigma^2)$  distribution if and only if for all  $n$ ,  $\bar{X}_n$  and  $s_n^2$  are independent random variables. Moreover, when  $F$  is  $N(\mu, \sigma^2)$ , then

$$\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right), \quad \text{and} \quad s_n^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2.$$

The “if” part is the profound part. It says that the independence of the natural estimates of the mean and the variance for any sample size forces the underlying distribution to be normal. We will sketch a proof of the “only if” part, i.e., we will assume that  $F$  is  $N(\mu, \sigma^2)$  and show that  $\bar{X}_n$  and  $s_n^2$  are independent.



*Proof.* Suppose that  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ . Define the *standardized versions* of the  $X_i$ 's as

$$Y_i = \frac{X_i - \mu}{\sigma}.$$

These are i.i.d.  $N(0, 1)$  random variables. Now, note that:

$$\bar{X} = \bar{Y} \sigma + \mu \quad \text{and} \quad s^2 = \frac{\sigma^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}{n-1}.$$

From the above display, we see that it suffices to show the independence of  $\bar{Y}$  and  $\sum_{i=1}^n (Y_i - \bar{Y})^2$ .

The way this proceeds is outlined below: Let  $\mathbf{Y}$  denote the  $n \times 1$  column vector  $(Y_1, Y_2, \dots, Y_n)^T$  and let  $A$  be an  $n \times n$  orthogonal matrix with the first row of  $A$  (which has length  $n$ ) being  $(1/\sqrt{n}, 1/\sqrt{n}, \dots, 1/\sqrt{n})$ .

Recall that an orthogonal matrix satisfies

$$A^\top A = AA^\top = I$$

where  $I$  is the identity matrix. Using standard linear algebra techniques it can be shown that such a  $A$  can always be constructed. For instance, in the case  $n = 2$ , we have

$$A = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}.$$

**Exercise:** Show that this  $A$  is orthogonal.

Now define a new random vector

$$\mathbf{W} = A\mathbf{Y}$$

and use the following result:

**Theorem 5.4.** *If  $Y_1, \dots, Y_n$  are i.i.d  $N(0, 1)$  and  $A$  is an orthogonal matrix and*

$$\mathbf{W} = A\mathbf{Y},$$

*then the random variables  $W_1, \dots, W_n$  are i.i.d  $N(0, 1)$ .*

*Proof of Theorem 5.4.* The joint p.d.f of  $\mathbf{Y} = (Y_1, \dots, Y_n)$  is

$$f_n(\mathbf{y}) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n y_i^2\right), \quad \text{for } \mathbf{y} \in \mathbb{R}^n.$$

Note that as  $\mathbf{Y} \mapsto A\mathbf{Y}$  is a linear transformation. The joint p.d.f of  $\mathbf{W} = A\mathbf{Y}$  is

$$g_n(\mathbf{w}) = \frac{1}{|\det A|} f_n(A^{-1}\mathbf{w}), \quad \text{for } \mathbf{w} \in \mathbb{R}^n.$$

Let  $\mathbf{y} = A^{-1}\mathbf{w}$ . Since  $A$  is orthogonal,  $|\det A| = 1$  and  $\mathbf{w}^\top \mathbf{w} = \sum_{i=1}^n w_i^2 = \mathbf{y}^\top \mathbf{y} = \sum_{i=1}^n y_i^2$ . So,

$$g_n(\mathbf{w}) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n w_i^2\right), \quad \text{for } \mathbf{w} \in \mathbb{R}^n.$$

Thus,  $\mathbf{W}$  has the same joint p.d.f as  $\mathbf{Y}$ . □

**Exercise:** Compute  $W$  if  $n = 2$  and

$$A = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}.$$

Returning to the proof of Proposition 5.3, note that

$$\mathbf{W}^\top \mathbf{W} = (A\mathbf{Y})^\top A\mathbf{Y} = \mathbf{Y}^\top A^\top A\mathbf{Y} = \mathbf{Y}^\top \mathbf{Y}$$

by the orthogonality of  $A$ —in other words,  $\sum_{i=1}^n W_i^2 = \sum_{i=1}^n Y_i^2$ . Also,

$$W_1 = Y_1/\sqrt{n} + Y_2/\sqrt{n} + \cdots + Y_n/\sqrt{n} = \sqrt{n} \bar{Y}_n.$$

Note that  $W_1$  is independent of  $W_2^2 + W_3^2 + \cdots + W_n^2$ . But

$$\sum_{i=2}^n W_i^2 = \sum_{i=1}^n W_i^2 - W_1^2 = \sum_{i=1}^n Y_i^2 - n \bar{Y}_n^2 = \sum_{i=1}^n (Y_i - \bar{Y}_n)^2.$$

It therefore follows that  $\sqrt{n} \bar{Y}_n$  and  $\sum_{i=1}^n (Y_i - \bar{Y}_n)^2$  are independent – which implies that  $\bar{Y}_n$  and  $\sum_{i=1}^n (Y_i - \bar{Y}_n)^2$  are independent.

Finally, we prove the distributional properties of  $\bar{X}_n$  and  $s_n^2$ . Note that  $\bar{Y}_n \sim N(0, 1/n)$ . Deduce that  $\bar{X}_n$  follows  $N(\mu, \sigma^2/n)$ . Since  $\sum_{i=1}^n (Y_i - \bar{Y}_n)^2 = W_2^2 + W_3^2 + \cdots + W_n^2$ , it follows that

$$\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \sum_{i=1}^n (Y_i - \bar{Y}_n)^2 \sim \chi_{n-1}^2.$$

Thus,

$$s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2. \quad (2)$$

□

### 5.3 The $t$ -distribution

**Definition 14.** Let  $Z \sim N(0, 1)$  and let  $V \sim \chi_n^2$  be independent of each other. Then,

$$T = \frac{Z}{\sqrt{V/n}}$$

is said to follow the  **$t$ -distribution** on  $n$  degrees of freedom. We write  $T \sim t_n$ .

The density of the  $t$ -distribution is derived in the text book (see Chapter 8.4). With a little bit of patience, you can also work it out, using the change of variable theorem appropriately (I won't go into the computational details here).

Here are some important facts about the  $t$ -distribution. Let  $T \sim t_n$ .

- (a)  $T$  and  $-T$  have the same distribution. Thus, the distribution of  $T$  is symmetric about 0 and it has an even density function.

Indeed, by definition,

$$-T = \frac{-Z}{\sqrt{V/n}} = \frac{\tilde{Z}}{\sqrt{V/n}},$$

where  $\tilde{Z} \equiv -Z$  follows  $N(0, 1)$ , and is independent of  $V$  where  $V$  follows  $\chi_n^2$ . Thus, by definition,  $-T$  also follows the  $t$ -distribution on  $n$  degrees of freedom.

- (b) As  $n \rightarrow \infty$ , the  $t_n$  distribution converges to the  $N(0, 1)$  distribution.

This follows from the law of large numbers. Consider the term  $V/n$  in the denominator of  $T$  for large  $n$ . As  $V$  follows  $\chi_n^2$  it has the same distribution as  $K_1 + K_2 + \cdots + K_n$  where  $K_i$ 's are i.i.d  $\chi_1^2$  random variables. But by the WLLN we know that

$$\frac{K_1 + K_2 + \cdots + K_n}{n} \xrightarrow{\mathbb{P}} \mathbb{E}(K_1) = 1 \quad (\text{check!}).$$

Thus  $V/n$  converges in probability to 1; hence the denominator in  $T$  converges in probability to 1 and  $T$  converges in distribution to  $Z$ , where  $Z$  is  $N(0, 1)$ .

**Theorem 5.5.** Suppose that  $X_1, \dots, X_n$  form a random sample from the normal distribution with mean  $\mu$  and variance  $\sigma^2$ . Then

$$\frac{\bar{X}_n - \mu}{\sqrt{s_n^2/n}} \sim t_{n-1}.$$

## 6 Confidence intervals

Confidence intervals (CIs) provide a method of quantifying uncertainty to an estimator  $\hat{\theta}$  when we wish to estimate an unknown parameter  $\theta$ . We want to find an interval  $(A, B)$  that we think has high probability of containing  $\theta$ .

**Definition:** Suppose that  $\mathbf{X}_n = (X_1, \dots, X_n)$  is a random sample from a distribution  $P_\theta$ ,  $\theta \in \Omega$ . Suppose that we want to estimate  $g(\theta)$ , a real-valued function of  $\theta$ . If  $A \leq B$  are two statistics with the property that for all values of  $\theta$ ,

$$\mathbb{P}_\theta(A \leq g(\theta) \leq B) \geq 1 - \alpha,$$

where  $\alpha \in (0, 1)$ , then the random interval  $(A, B)$  is called a **confidence interval** for  $g(\theta)$  with **(confidence) level**  $1 - \alpha$ . If the inequality “ $\geq 1 - \alpha$ ” is an equality for all  $\theta$ , the CI is called **exact**.

### 6.1 Construction of confidence interval using a pivot

How do we construct confidence intervals? One approach is to use a so-called pivot.

**Definition:** A random variable  $\Psi(X_1, X_2, \dots, X_n, g(\theta))$  is called a **pivot** for  $g(\theta)$  if its distribution is independent of  $\theta$ .

**Example 1:** Find a level  $(1 - \alpha)$  CI for  $\mu$  from data  $X_1, X_2, \dots, X_n$  which are i.i.d.  $N(\mu, \sigma^2)$  where  $\sigma$  is **known**. Here  $\theta = \mu$  and  $g(\theta) = \mu$ .

The most intuitive estimator of  $\mu$  here is the sample mean  $\bar{X}_n$ . We know that

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

The standardized version of the sample mean follows  $N(0, 1)$  and can therefore act as a pivot. In other words, construct,

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

for every value of  $\mu$ .

With  $z_\beta$  denoting the **upper  $\beta$ -quantile** of  $N(0, 1)$  (i.e.,  $\mathbb{P}(Z > z_\beta) = \beta$  where  $Z$  follows  $N(0, 1)$ ) we can write:

$$\mathbb{P}_\mu \left( -z_{\alpha/2} \leq \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \leq z_{\alpha/2} \right) = 1 - \alpha.$$

From the above display we can find limits for  $\mu$  such that the above inequalities are simultaneously satisfied. On doing the algebra, we get:

$$\mathbb{P}_\mu \left( \bar{X}_n - \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \leq \mu \leq \bar{X}_n + \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \right) = 1 - \alpha.$$

Thus

$$\left( \bar{X} - \frac{\sigma}{\sqrt{n}} z_{\alpha/2}, \bar{X} + \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \right)$$

is an exact CI for  $\mu$  with level  $1 - \alpha$ .

The general **method of pivots** works as follows:

- (1) Construct a function  $\Psi$  of the data  $\mathbf{X}_n$  and  $g(\theta)$ , say  $\Psi(\mathbf{X}_n, g(\theta))$ , such that the distribution of this random variable under parameter value  $\theta$  *does not depend on*  $\theta$  and is known. Such a  $\Psi$  gives a pivot for  $g(\theta)$ .
- (2) Let  $G$  denote the distribution function of this pivot. The idea now is to get a range of plausible values of the pivot. The level of confidence  $1 - \alpha$  is to be used to get the appropriate range.

This can be done in a variety of ways but the following is standard. Denote by  $q(G; \beta)$  the  $\beta$ -**quantile** of  $G$ , i.e.,

$$\mathbb{P}_\theta[\Psi(\mathbf{X}_n, g(\theta)) \leq q(G; \beta)] = \beta.$$

- (3) Choose  $0 \leq \beta_1, \beta_2 \leq \alpha$  such that  $\beta_1 + \beta_2 = \alpha$ . Then,

$$\mathbb{P}_\theta[q(G; \beta_1) \leq \Psi(\mathbf{X}_n, g(\theta)) \leq q(G; 1 - \beta_2)] = 1 - \beta_2 - \beta_1 = 1 - \alpha.$$

- (4) Solve the inequalities  $q(G; \beta_1) \leq \Psi(\mathbf{X}_n, g(\theta)) \leq q(G; 1 - \beta_2)$  for  $\theta$  to obtain a confidence interval for  $g(\theta)$ .

**Example 2:** The data are the same as in Example 1 but now  $\sigma^2$  is no longer known. Thus, the parameter of unknowns  $\theta = (\mu, \sigma^2)$  and we are interested in finding a CI for  $g(\theta) = \mu$ .

Clearly, setting

$$\Psi(\mathbf{X}_n, \mu) = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$

will not work smoothly here. This certainly has a known  $(N(0, 1))$  distribution but involves the *nuisance parameter*  $\sigma$  (an unknown parameter that we are not primarily interested in).

However, one can replace  $\sigma$  by  $s_n$ , where  $s_n^2$  is the sample variance. So, set:

$$\Psi(\mathbf{X}_n, \mu) = \frac{\bar{X}_n - \mu}{\sqrt{s_n^2/n}}.$$

This only depends on the data and  $g(\theta) = \mu$ .

This is indeed a pivot: By Theorem 5.5, the new  $\Psi(\mathbf{X}_n, \mu)$  has a  $t_{n-1}$  distribution (which is independent of  $\mu$ ). Thus,  $G$  here is the  $t_{n-1}$  distribution and we can choose

the quantiles to be  $q(t_{n-1}; \alpha/2)$  and  $q(t_{n-1}; 1 - \alpha/2)$ . By symmetry of the  $t_{n-1}$  distribution about 0, we have,  $q(t_{n-1}; \alpha/2) = -q(t_{n-1}; 1 - \alpha/2)$ . It follows that,

$$\mathbb{P}_{\mu, \sigma^2} \left[ -q(t_{n-1}; 1 - \alpha/2) \leq \frac{\sqrt{n}(\bar{X}_n - \mu)}{s_n} \leq q(t_{n-1}; 1 - \alpha/2) \right] = 1 - \alpha.$$

As with Example 1, direct algebraic manipulations show that this is the same as the statement:

$$\mathbb{P}_{\mu, \sigma^2} \left[ \bar{X}_n - \frac{s_n}{\sqrt{n}} q(t_{n-1}; 1 - \alpha/2) \leq \mu \leq \bar{X}_n + \frac{s_n}{\sqrt{n}} q(t_{n-1}; 1 - \alpha/2) \right] = 1 - \alpha.$$

This gives a level  $1 - \alpha$  confidence set for  $\mu$ .

**Remark:** Both  $z_\alpha$  and  $q(t_n, \alpha)$  can be found from the table on p. 860 of the textbook. For example,

$$z_{0.05} = q(N(0, 1), 1 - 0.05) = q(t_\infty, 0.95) = 1.645, \quad q(t_{16}, 0.99) = 2.583.$$

**Food for thought:** In each of the above examples there are innumerable ways of decomposing  $\alpha$  as  $\beta_1 + \beta_2$ . It turns out that when  $\alpha$  is split equally the level  $1 - \alpha$  CIs obtained in Examples 1 and 2 are the shortest.

What are desirable properties of confidence sets? On one hand, we require high levels of confidence; in other words, we would like  $\alpha$  to be as small as possible. On the other hand we would like our CIs to be shortest possible. Unfortunately, we cannot simultaneously make the confidence levels of our CIs go up and the lengths of our CIs go down.

In Example 1, the length of the level  $(1 - \alpha)$  CI is

$$2\sigma \frac{z_{\alpha/2}}{\sqrt{n}}.$$

As we reduce  $\alpha$  (for higher confidence),  $z_{\alpha/2}$  increases, making the CI wider.

However, we can reduce the length of our CI for a fixed  $\alpha$  by increasing the sample size. If my sample size is 4 times yours, I will end up with a CI which has the same level as yours but has half the length of your CI.

Can we hope to get absolute confidence, i.e.  $\alpha = 0$ ? That is too much of an ask. When  $\alpha = 0$ ,  $z_{\alpha/2} = \infty$  and the CIs for  $\mu$  are infinitely large. The same can be verified for Example 2.

## 6.2 Asymptotic confidence intervals

The CLT allows us to construct an *approximate pivot* for large sample sizes for estimating the population mean  $\mu$  for any underlying distribution  $F$ .

Let  $X_1, X_2, \dots, X_n$  be i.i.d observations from some common distribution  $F$  and let

$$\mathbb{E}(X_1) = \mu \quad \text{and} \quad \text{Var}(X_1) = \sigma^2.$$

We are interested in constructing an approximate level  $(1 - \alpha)$  CI for  $\mu$ , *assuming that  $\sigma$  is known*.

By the CLT we have  $\bar{X}_n \sim_{\text{approx}} N(\mu, \sigma^2/n)$  for large  $n$ ; in other words,

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \sim_{\text{approx}} N(0, 1).$$

If  $\sigma$  is known the above quantity is an approximate pivot and following Example 1, we can therefore write,

$$\mathbb{P}_\mu \left( -z_{\alpha/2} \leq \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \leq z_{\alpha/2} \right) \approx 1 - \alpha.$$

As before, this translates to

$$\mathbb{P}_\mu \left( \bar{X}_n - \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \leq \mu \leq \bar{X}_n + \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \right) \approx 1 - \alpha.$$

This gives an approximate level  $(1 - \alpha)$  CI for  $\mu$  when  $\sigma$  is known. The approximation will improve as the sample size  $n$  increases. Note that the true coverage of the above CI may be different from  $1 - \alpha$  and can depend heavily on the nature of  $F$  and the sample size  $n$ .

Realistically, however,  $\sigma$  is unknown and is replaced by  $s_n$ . Since we are dealing with large sample sizes,  $s_n$  is with very high probability close to  $\sigma$  and the interval

$$\left( \bar{X}_n - \frac{s_n}{\sqrt{n}} z_{\alpha/2}, \bar{X}_n + \frac{s_n}{\sqrt{n}} z_{\alpha/2} \right),$$

still remains an approximate level  $(1 - \alpha)$  CI.

If  $\sigma^2$  is unknown but can be expressed in terms of the unknown parameter  $\mu$ , there is a better approach than just using  $s_n$ .

**Exercise:** Suppose  $X_1, X_2, \dots, X_n$  are i.i.d Bernoulli( $\theta$ ). The sample size  $n$  is large.

Thus

$$\mathbb{E}(X_1) = \theta \quad \text{and} \quad \text{Var}(X_1) = \theta(1 - \theta).$$

We want to find an approximate CI for  $\theta$  at level  $1 - \alpha$ . Note that both mean and variance are unknown but  $\sigma^2 = \theta(1 - \theta)$  is a function of  $\theta$ .

Show that if  $\hat{\theta}$  is natural estimate of  $\theta$  obtained by computing the sample proportion of 1's, then

$$\left[ \hat{\theta} - \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n-1}} z_{\alpha/2}, \hat{\theta} + \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n-1}} z_{\alpha/2} \right]$$

is an approximate level  $(1 - \alpha)$  CI for  $\theta$ .

**Interpretation of confidence intervals:** Let  $(A, B)$  be a coefficient  $\gamma$  confidence interval for a parameter  $\theta$ . Let  $(a, b)$  be the observed value of the interval.

It is NOT correct to say that “ $\theta$  lies in the interval  $(a, b)$  with *probability*  $\gamma$ ”.

It is true that “ $\theta$  will lie in the random intervals having endpoints  $A(X_1, \dots, X_n)$  and  $B(X_1, \dots, X_n)$  with probability  $\gamma$ ”.

After observing the specific values  $A(X_1, \dots, X_n) = a$  and  $B(X_1, \dots, X_n) = b$ , it is not possible to assign a probability to the event that  $\theta$  lies in the specific interval  $(a, b)$  without regarding  $\theta$  as a random variable.

We usually say that there is *confidence*  $\gamma$  that  $\theta$  lies in the interval  $(a, b)$ .

$$\frac{\sqrt{n}(\hat{\theta} - \theta)}{\sqrt{\theta(1-\theta)}} \sim_{\text{appx}} N(0, 1),$$

so that

$$P_{\theta} \left[ -z_{\alpha/2} \leq \frac{\sqrt{n}(\hat{\theta} - \theta)}{\sqrt{\theta(1-\theta)}} \leq z_{\alpha/2} \right] =_{\text{appx}} 1 - \alpha.$$

An approximate level  $(1 - \alpha)$  CI can now be obtained by solving for all  $\theta$  for which the above inequalities are both satisfied. This amounts to solving a quadratic and will yield a different C.I. than the one proposed in the Exercise. You should try to work out what this gives you.



## 7 The Cramér–Rao Information Inequality

We saw in the last lecture that for a variety of different models one could differentiate the log-likelihood function with respect to the parameter  $\theta$  and set this equal to 0 to obtain the MLE of  $\theta$ .

In these examples, the log-likelihood as a function of  $\theta$  is strictly concave (looks like an inverted bowl) and hence solving for the stationary point gives us the unique maximizer of the log-likelihood.

We start this section by introducing some notation. Let  $X$  be a random variable with p.d.f  $f(\cdot, \theta)$ , where  $\theta \in \Omega$ , and

$$\ell(x, \theta) = \log f(x, \theta) \quad \text{and} \quad \dot{\ell}(x, \theta) = \frac{\partial}{\partial \theta} \ell(x, \theta).$$

As before,  $\mathbf{X}_n$  denotes the vector  $(X_1, X_2, \dots, X_n)$  and  $\mathbf{x}$  denotes a particular value  $(x_1, x_2, \dots, x_n)$  assumed by the random vector  $\mathbf{X}_n$ .

We denote by  $f_n(\mathbf{x}, \theta)$  the value of the density of  $\mathbf{X}_n$  at the point  $\mathbf{x}$ . Then,

$$f_n(\mathbf{x}, \theta) = \prod_{i=1}^n f(x_i, \theta).$$

Thus,

$$L_n(\theta, \mathbf{X}_n) = \prod_{i=1}^n f(X_i, \theta) = f_n(\mathbf{X}_n, \theta)$$

and

$$\ell_n(\mathbf{X}_n, \theta) = \log L_n(\theta, \mathbf{X}_n) = \sum_{i=1}^n \ell(X_i, \theta).$$

Differentiating with respect to  $\theta$  yields

$$\dot{\ell}_n(\mathbf{X}_n, \theta) = \frac{\partial}{\partial \theta} \log f_n(\mathbf{X}_n, \theta) = \sum_{i=1}^n \dot{\ell}(X_i, \theta).$$

We call  $\dot{\ell}(x, \theta)$  the **score function** and

$$\dot{\ell}_n(\mathbf{X}_n, \theta) = 0$$

the **score equation**. If differentiation is permissible for the purpose of obtaining the MLE, then  $\hat{\theta}_n$ , the MLE, solves the equation

$$\dot{\ell}_n(\mathbf{X}_n, \theta) \equiv \sum_{i=1}^n \dot{\ell}(X_i, \theta) = 0.$$

In this section, our first goal is to find a (nontrivial) **lower bound** on the **variance of unbiased estimators** of  $g(\theta)$  where  $g : \Omega \rightarrow \mathbb{R}$  is some differentiable function.

If we can indeed find such a bound (albeit under some regularity conditions) and there is an unbiased estimator of  $g(\theta)$  that attains this lower bound, we can conclude that it is the MVUE of  $g(\theta)$ .

We now impose the following restrictions (regularity conditions) on the model.

(A.1) The set  $A_\theta = \{x : f(x, \theta) > 0\}$  actually does NOT depend on  $\theta$  and is subsequently denoted by  $A$ .

(A.2) If  $W(\mathbf{X}_n)$  is a statistic such that  $\mathbb{E}_\theta(|W(\mathbf{X}_n)|) < \infty$  for all  $\theta$ , then,

$$\frac{\partial}{\partial \theta} \mathbb{E}_\theta[W(\mathbf{X}_n)] = \frac{\partial}{\partial \theta} \int_{A^n} W(\mathbf{x}) f_n(\mathbf{x}, \theta) d\mathbf{x} = \int_{A^n} W(\mathbf{x}) \frac{\partial}{\partial \theta} f_n(\mathbf{x}, \theta) d\mathbf{x}.$$

(A.3) The quantity  $\frac{\partial}{\partial \theta} \log f(x, \theta)$  exists for all  $x \in A$  and all  $\theta \in \Omega$  as a well-defined finite quantity.

The first condition says that the set of possible values of the data vector on which the distribution of  $\mathbf{X}_n$  is supported does not vary with  $\theta$ ; this therefore rules out families of distribution like the uniform.

The second assumption is a “smoothness assumption” on the family of densities and is generally happily satisfied for most parametric models we encounter in statistics.

There are various types of simple sufficient conditions that one can impose on  $f(x, \theta)$  to make the interchange of integration and differentiation possible — we shall however not bother about these for the moment.

## 7.1 Information

For most of the sequel, for notational simplicity, we will assume that the parameter space  $\Omega \subset \mathbb{R}$ . We define the **Fisher information** about the parameter  $\theta$  in the model, namely  $I(\theta)$ , by

$$I(\theta) := \mathbb{E}_\theta[\dot{\ell}^2(X, \theta)],$$

provided it exists as a finite quantity for every  $\theta \in \Omega$ .

We then have the following theorem.

**Theorem 7.1** (Cramér–Rao inequality). *All notation being as above, if  $T(\mathbf{X}_n)$  is an unbiased estimator of  $g(\theta)$ , then*

$$\text{Var}_\theta(T(\mathbf{X}_n)) \geq \frac{[g'(\theta)]^2}{nI(\theta)},$$

*provided assumptions A.1, A.2 and A.3 hold, and  $I(\theta)$  exists and is finite for all  $\theta$ .*

The above inequality is the celebrated **Cramér–Rao inequality** (or the information inequality) and is one of the most well-known inequalities in statistics and has important ramifications in even more advanced forms of inference.

Notice that if we take  $g(\theta) = \theta$  then  $n^{-1}I(\theta)^{-1}$  gives us a lower bound on the variance of unbiased estimators of  $\theta$  in the model.

If  $I(\theta)$  is small, the lower bound is large, so unbiased estimators are doing a poor job in general—in other words, the data is not that informative about  $\theta$  (within the context of unbiased estimation).

On the other hand, if  $I(\theta)$  is big, the lower bound is small, and so if we have a best unbiased estimator of  $\theta$  that actually attains this lower bound, we are doing a good job. That is why  $I(\theta)$  is referred to as the information about  $\theta$ .

**Proof of Theorem 7.1:** By the Cauchy–Schwarz inequality (see Theorem 4.6.3 in the textbook),

$$\text{Cov}_\theta^2\left(T(\mathbf{X}_n), \dot{\ell}_n(\mathbf{X}_n, \theta)\right) \leq \text{Var}_\theta(T(\mathbf{X}_n))\text{Var}_\theta(\dot{\ell}_n(\mathbf{X}_n, \theta)). \quad (3)$$

As

$$1 = \int f_n(\mathbf{x}, \theta) d\mathbf{x}, \quad \text{for all } \theta \in \Omega,$$

on differentiating both sides of the above identity with respect to  $\theta$  and using (A.2) with  $W(\mathbf{x}) \equiv 1$  we obtain,

$$\begin{aligned} 0 &= \int \frac{\partial}{\partial \theta} f_n(\mathbf{x}, \theta) d\mathbf{x} = \int \left( \frac{\partial}{\partial \theta} f_n(\mathbf{x}, \theta) \right) \frac{1}{f_n(\mathbf{x}, \theta)} f_n(\mathbf{x}, \theta) d\mathbf{x} \\ &= \int \left( \frac{\partial}{\partial \theta} \log f_n(\mathbf{x}, \theta) \right) f_n(\mathbf{x}, \theta) d\mathbf{x}. \end{aligned}$$

The last expression in the above display is precisely  $\mathbb{E}_\theta[\dot{\ell}_n(\mathbf{X}_n, \theta)]$  which therefore is equal to 0. Note that

$$\mathbb{E}_\theta[\dot{\ell}_n(\mathbf{X}_n, \theta)] = \mathbb{E}_\theta \left[ \sum_{i=1}^n \dot{\ell}(X_i, \theta) \right] = n\mathbb{E}_\theta[\dot{\ell}(X, \theta)],$$

since the  $\dot{\ell}(X_i, \theta)$ 's are i.i.d. Thus, we have  $\mathbb{E}_\theta[\dot{\ell}(X_1, \theta)] = 0$ . This implies that

$$I(\theta) = \text{Var}_\theta(\dot{\ell}(X, \theta)).$$

Further, let  $I_n(\theta) := \mathbb{E}_\theta[\dot{\ell}_n^2(\mathbf{X}_n, \theta)]$ . Then

$$\begin{aligned} I_n(\theta) &= \text{Var}_\theta(\dot{\ell}_n(\mathbf{X}_n, \theta)) = \text{Var}_\theta \left( \sum_{i=1}^n \dot{\ell}(X_i, \theta) \right) \\ &= \sum_{i=1}^n \text{Var}_\theta(\dot{\ell}(X_i, \theta)) = nI(\theta). \end{aligned}$$

We will refer to  $I_n(\theta)$  as the **Fisher information in the sample  $\mathbf{X}_n$** . Since  $\mathbb{E}_\theta[\dot{\ell}_n(\mathbf{X}_n, \theta)] = 0$ , it follows that

$$\begin{aligned} \text{Cov}_\theta \left( T(\mathbf{X}_n), \dot{\ell}_n(\mathbf{X}_n, \theta) \right) &= \int T(\mathbf{x}) \dot{\ell}_n(\mathbf{x}, \theta) f_n(\mathbf{x}, \theta) d\mathbf{x} \\ &= \int T(\mathbf{x}) \left( \frac{\partial}{\partial \theta} f_n(\mathbf{x}, \theta) \right) d\mathbf{x} \\ &= \frac{\partial}{\partial \theta} \int T(\mathbf{x}) f_n(\mathbf{x}, \theta) d\mathbf{x} \quad (\text{by (A.2)}) \\ &= \frac{\partial}{\partial \theta} g(\theta) = g'(\theta). \end{aligned}$$

Using the above in conjunction in (3) we get,

$$[g'(\theta)]^2 \leq \text{Var}_\theta(T(\mathbf{X}_n)) I_n(\theta)$$

which is equivalent to what we set out to prove.  $\square$

There is an alternative expression for the Fisher information  $I(\theta)$  in terms of the second derivative of the log-likelihood with respect to  $\theta$ . If

$$\ddot{\ell}(x, \theta) := \frac{\partial^2}{\partial \theta^2} \log f(x, \theta)$$

exists for all  $x \in A$  and for all  $\theta \in \Theta$ , then we have the following identity:

$$I(\theta) = \mathbb{E}_\theta[\dot{\ell}(X, \theta)^2] = -\mathbb{E}_\theta[\ddot{\ell}(X, \theta)],$$

provided we can differentiate twice under the integral sign; more concretely, if

$$\int \frac{\partial^2}{\partial \theta^2} f(x, \theta) dx = \frac{\partial^2}{\partial \theta^2} \int f(x, \theta) dx = 0 \quad (\star).$$

To prove the above identity, first note that,

$$\dot{\ell}(x, \theta) = \frac{1}{f(x, \theta)} \left[ \frac{\partial}{\partial \theta} f(x, \theta) \right].$$

Now,

$$\begin{aligned} \ddot{\ell}(x, \theta) &= \frac{\partial}{\partial \theta} \left( \dot{\ell}(x, \theta) \right) = \frac{\partial}{\partial \theta} \left( \frac{1}{f(x, \theta)} \frac{\partial}{\partial \theta} f(x, \theta) \right) \\ &= \frac{\partial^2}{\partial \theta^2} f(x, \theta) \frac{1}{f(x, \theta)} - \frac{1}{f^2(x, \theta)} \left( \frac{\partial}{\partial \theta} f(x, \theta) \right)^2 \\ &= \frac{\partial^2}{\partial \theta^2} f(x, \theta) \frac{1}{f(x, \theta)} - \dot{\ell}(x, \theta)^2. \end{aligned}$$

Thus,

$$\begin{aligned}\mathbb{E}_\theta[\ddot{\ell}(X, \theta)] &= \int \ddot{\ell}(x, \theta) f(x, \theta) dx \\ &= \int \frac{\partial^2}{\partial \theta^2} f(x, \theta) dx - \mathbb{E}_\theta[\dot{\ell}^2(X, \theta)] \\ &= 0 - \mathbb{E}_\theta[\dot{\ell}^2(X, \theta)],\end{aligned}$$

where the first term on the right side vanishes by virtue of  $(\star)$ . This establishes the desired equality. It follows that,

$$I_n(\theta) = \mathbb{E}_\theta[-\ddot{\ell}_n(\mathbf{X}_n, \theta)],$$

where  $\ddot{\ell}_n(\mathbf{X}_n, \theta)$  is the second partial derivative of  $\ell_n(\mathbf{X}_n, \theta)$  with respect to  $\theta$ . To see this, note that,

$$\ddot{\ell}_n(\mathbf{X}_n, \theta) = \frac{\partial^2}{\partial \theta^2} \left( \sum_{i=1}^n \ell(X_i, \theta) \right) = \sum_{i=1}^n \ddot{\ell}(X_i, \theta),$$

so that

$$\mathbb{E}_\theta[\ddot{\ell}_n(\mathbf{X}_n, \theta)] = \sum_{i=1}^n \mathbb{E}_\theta[\ddot{\ell}(X_i, \theta)] = n \mathbb{E}_\theta[\ddot{\ell}(X, \theta)] = -n I(\theta).$$

We have just established the following result.

**Theorem 7.2.** *If (A.1)–(A.3) as well as  $(\star)$  hold, then*

$$I(\theta) = -\mathbb{E}_\theta[\ddot{\ell}(X, \theta)] = \text{Var}_\theta(\dot{\ell}(X, \theta)).$$

## 7.2 Examples

We now look at some applications of the Cramér–Rao inequality.

**Example 1:** Let  $X_1, X_2, \dots, X_n$  be i.i.d  $\text{Pois}(\theta)$ ,  $\theta > 0$ . Then

$$\mathbb{E}_\theta(X_1) = \theta \quad \text{and} \quad \text{Var}_\theta(X_1) = \theta.$$

Let us first write down the likelihood of the data. We have,

$$f_n(\mathbf{x}, \theta) = \prod_{i=1}^n \frac{e^{-\theta} \theta^{x_i}}{x_i!} = e^{-n\theta} \theta^{\sum_{i=1}^n x_i} \left( \prod_{i=1}^n x_i! \right)^{-1}.$$

Thus,

$$\begin{aligned}\ell_n(\mathbf{x}, \theta) &= -n\theta + \log \theta \left( \sum_{i=1}^n x_i \right) - \log \prod_{i=1}^n x_i! \\ \dot{\ell}_n(\mathbf{x}, \theta) &= -n + \frac{1}{\theta} \sum_{i=1}^n x_i.\end{aligned}$$

Thus the Fisher information about  $\theta$  in the sample  $X_1, \dots, X_n$  is given by

$$I_n(\theta) = \text{Var}_\theta \left( -n + \frac{1}{\theta} \sum_{i=1}^n X_i \right) = \frac{1}{\theta^2} \text{Var}_\theta \left( \sum_{i=1}^n X_i \right) = \frac{n\theta}{\theta^2} = \frac{n}{\theta}.$$

The assumptions needed for the Cram r–Rao inequality to hold are all satisfied for this model, and it follows that for any unbiased estimator  $T(\mathbf{X}_n)$  of  $g(\theta) = \theta$  we have,

$$\text{Var}_\theta(T(\mathbf{X}_n)) \geq \frac{1}{I_n(\theta)} = \frac{\theta}{n}.$$

Since  $\bar{X}_n$  is unbiased for  $\theta$  and has variance  $\theta/n$  we conclude that  $\bar{X}_n$  is the MVUE of  $\theta$ .

**Example 2:** Let  $X_1, X_2, \dots, X_n$  be i.i.d  $N(0, V)$ . Consider once again, the joint density of the  $n$  observations:

$$f_n(\mathbf{x}, V) = \frac{1}{(2\pi V)^{n/2}} \exp \left( -\frac{1}{2V} \sum_{i=1}^n x_i^2 \right).$$

Now,

$$\begin{aligned} \dot{\ell}_n(\mathbf{x}, V) &= \frac{\partial}{\partial V} \left( -\frac{n}{2} \log 2\pi - \frac{n}{2} \log V - \frac{1}{2V} \sum_{i=1}^n x_i^2 \right) \\ &= -\frac{n}{2V} + \frac{1}{2V^2} \sum_{i=1}^n x_i^2. \end{aligned}$$

Differentiating yet again we obtain,

$$\ddot{\ell}_n(\mathbf{x}, V) = \frac{n}{2V^2} - \frac{1}{V^3} \sum_{i=1}^n x_i^2.$$

Then, the Fisher information for  $V$  based on  $X_1, \dots, X_n$  is

$$I_n(V) = -\mathbb{E}_V \left( \frac{n}{2V^2} - \frac{1}{V^3} \sum_{i=1}^n X_i^2 \right) = \frac{n}{2V^2} + \frac{1}{V^3} nV = \frac{n}{2V^2}.$$

Now consider the problem of estimating  $g(V) = V$ . For any unbiased estimator  $S(\mathbf{X}_n)$  of  $V$ , the Cram r–Rao inequality tells us that

$$\text{Var}_V(S(\mathbf{X}_n)) \geq I_n(V)^{-1} = \frac{2V^2}{n}.$$

Consider,  $\sum_{i=1}^n X_i^2/n$  as an estimator of  $V$ . This is clearly unbiased for  $V$  and the variance is given by,

$$\text{Var}_V \left( \frac{1}{n} \sum_{i=1}^n X_i^2 \right) = \frac{1}{n} \text{Var}_V(X_1^2) = \frac{V^2}{n} \text{Var}_V \left( \frac{X_1^2}{V} \right) = \frac{2V^2}{n},$$

since  $X_1^2/V \sim \chi_1^2$  which has variance 2. It follows that  $\sum X_i^2/n$  is the MVUE of  $V$ .

### 7.3 Large sample properties of the MLE

In this subsection we study some of the large sample properties of the MLE in standard parametric models and how these can be used to construct confidence sets for  $\theta$  or a function of  $\theta$ . We will see in this section that in the long run MLEs are the best possible estimators in a variety of different models.

We will stick to models satisfying the restrictions (A.1)–(A.3) imposed in the last section. Hence our results will not apply to the uniform distribution (or ones similar to the uniform).

Let us throw our minds back to the Cramér–Rao inequality. When does an unbiased estimator  $T(\mathbf{X}_n)$  of  $g(\theta)$  attain the bound given by this inequality? This requires:

$$\text{Var}_\theta(T(\mathbf{X}_n)) = \frac{(g'(\theta))^2}{n I(\theta)}.$$

But this is equivalent to the assertion that the correlation between  $T(\mathbf{X}_n)$  and  $\dot{\ell}_n(\mathbf{X}_n, \theta)$  is equal to 1 or -1.

This means that  $\dot{\ell}_n(\mathbf{X}_n, \theta)$  can be expressed as a *linear function* of  $T(\mathbf{X}_n)$ .

In fact, this is a necessary and sufficient condition for the information bound to be attained by the variance of  $T(\mathbf{X}_n)$ .

It turns out that this is generally difficult to achieve. Thus, there will be many different functions of  $\theta$ , for which best unbiased estimators will exist but whose variance will not hit the information bound. The example below will illustrate this point.

**Example:** Let  $X_1, X_2, \dots, X_n$  be i.i.d Ber( $\theta$ ). We have,

$$f(x, \theta) = \theta^x (1 - \theta)^{1-x} \quad \text{for } x = 0, 1.$$

Thus,

$$\ell(x, \theta) = x \log \theta + (1 - x) \log(1 - \theta),$$

$$\dot{\ell}(x, \theta) = \frac{x}{\theta} - \frac{1 - x}{1 - \theta}$$

and

$$\ddot{\ell}(x, \theta) = -\frac{x}{\theta^2} - \frac{1 - x}{(1 - \theta)^2}.$$

Thus,

$$\dot{\ell}_n(\mathbf{X}_n, \theta) = \sum_{i=1}^n \dot{\ell}(X_i, \theta) = \frac{\sum_{i=1}^n X_i}{\theta} - \frac{n - \sum_{i=1}^n X_i}{1 - \theta}.$$

Recall that the MLE solves  $\dot{\ell}_n(\mathbf{X}_n, \theta) = 0$ .

Check that in this situation, this gives you precisely  $\bar{X}_n$  as your MLE.

Let us compute the Fisher information  $I(\theta)$ . We have,

$$I(\theta) = -\mathbb{E}_\theta[\ddot{\ell}(X_1, \theta)] = \mathbb{E}_\theta \left( \frac{X_1}{\theta^2} + \frac{1 - X_1}{(1 - \theta)^2} \right) = \frac{1}{\theta} + \frac{1}{1 - \theta} = \frac{1}{\theta(1 - \theta)}.$$

Thus,

$$I_n(\theta) = nI(\theta) = \frac{n}{\theta(1 - \theta)}.$$

Consider unbiased estimation of  $\Psi(\theta) = \theta$  based on  $\mathbf{X}_n$ . Let  $T(\mathbf{X}_n)$  be an unbiased estimator of  $\theta$ . Then, by the information inequality,

$$\text{Var}_\theta(T(\mathbf{X}_n)) \geq \frac{\theta(1 - \theta)}{n}.$$

Note that the variance of  $\bar{X}_n$  is precisely  $\theta(1 - \theta)/n$ , so that it is the MVUE of  $\theta$ . Note that

$$\dot{\ell}_n(\mathbf{X}_n, \theta) = \frac{n\bar{X}}{\theta} - \frac{n(1 - \bar{X})}{1 - \theta} = \left( \frac{n}{\theta} + \frac{n}{1 - \theta} \right) \bar{X} - \frac{n}{1 - \theta}.$$

Thus,  $\bar{X}_n$  is indeed linear in  $\dot{\ell}_n(\mathbf{X}_n, \theta)$ .

Consider now estimating a different function of  $\theta$ , say  $g(\theta) = \theta^2$ .

This is the probability of getting two consecutive heads. Suppose we try to find an unbiased estimator of this parameter.

Then  $S(\mathbf{X}_n) = X_1X_2$  is an unbiased estimator ( $\mathbb{E}_\theta(X_1X_2) = \mathbb{E}_\theta(X_1)\mathbb{E}_\theta(X_2) = \theta^2$ ), but then so is  $X_iX_j$  for any  $i \neq j$ .

We can find the MVUE of  $\theta^2$  in this model by using techniques beyond the scope of this course—it can be shown that any estimator  $T(\mathbf{X}_n)$  that can be written as a function of  $\bar{X}_n$  and is unbiased for  $\theta^2$  is an MVUE (and indeed there is one such).

Verify that,

$$T^*(\mathbf{X}_n) = \frac{n\bar{X}_n^2 - \bar{X}_n}{n - 1}$$

is unbiased for  $\theta^2$  and is therefore an (in fact *the*) MVUE.

However, the variance of  $T^*(\mathbf{X}_n)$  does not attain the information bound for estimating  $g(\theta)$  which is  $4\theta^3(1 - \theta)/n$  (Exercise). This can be checked by direct (somewhat tedious) computation or by noting that  $T^*(\mathbf{X}_n)$  is not a linear function of  $\dot{\ell}_n(\mathbf{X}_n, \theta)$ .

The question then is whether we can propose an estimator of  $\theta^2$  that does achieve the bound, at least approximately, in the long run.

It turns out that this is actually possible. Since the MLE of  $\theta$  is  $\bar{X}$ , the MLE of  $g(\theta)$  is proposed as the plug-in value  $g(\bar{X}) = \bar{X}^2$ .



This is *not an unbiased estimator of  $g(\theta)$*  in finite samples, but has excellent behavior in the long run. In fact,

$$\sqrt{n}(g(\bar{X}_n) - g(\theta)) \rightarrow_d N(0, 4\theta^3(1 - \theta)).$$

**Exercise:** Prove the last statement.

Thus for large values of  $n$ ,  $g(\bar{X})$  behaves approximately like a normal random variable with mean  $g(\theta)$  and variance  $4\theta^3(1 - \theta)/n$ .

In this sense,  $g(\bar{X}_n)$  is *asymptotically (in the long run) unbiased and asymptotically efficient* (in the sense that it has minimum variance).

Here is an important proposition that establishes the limiting behavior of the MLE.

**Theorem 7.3.** *If  $\hat{\theta}_n$  is the MLE of  $\theta$  obtained by solving*

$$\sum_{i=1}^n \dot{\ell}(X_i, \theta) = 0,$$

*then the following representation for the MLE is valid:*

$$\sqrt{n}(\hat{\theta}_n - \theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n I(\theta)^{-1} \dot{\ell}(X_i, \theta) + r_n,$$

*where  $r_n$  converges to 0 in probability. It follows by a direct application of the CLT that,*

$$\sqrt{n}(\hat{\theta}_n - \theta) \rightarrow_d N(0, I(\theta)^{-1}).$$

The above result shows MLE  $\hat{\theta}$  is (asymptotically) the best possible estimator: Not only does its long term distribution center around  $\theta$ , the quantity of interest, its distribution is also less spread out than that of any “reasonable” estimator of  $\theta$ . If  $S_n$  is a “reasonable” estimator of  $\theta$ , with

$$\sqrt{n}(S_n - \theta) \rightarrow_d N(0, \xi^2(\theta)),$$

then  $\xi^2(\theta) \geq I(\theta)^{-1}$ .

Recall the delta method.

**Proposition 7.4** (Delta method). *Suppose  $T_n$  is an estimator of  $\theta$  (based on i.i.d observations,  $X_1, X_2, \dots, X_n$  from  $P_\theta$ ) that satisfies:*

$$\sqrt{n}(T_n - \theta) \rightarrow_d N(0, \sigma^2(\theta)).$$

*Here  $\sigma^2(\theta)$  is the limiting variance and depends on the underlying parameter  $\theta$ . Then, for a continuously differentiable function  $h$  such that  $h'(g(\theta)) \neq 0$ , we have:*

$$\sqrt{n}(g(T_n) - g(\theta)) \rightarrow_d N(0, (g'(\theta))^2 \sigma^2(\theta)).$$

We can now deduce the limiting behavior of the MLE of  $g(\theta)$  given by  $g(\hat{\theta}_n)$  for any smooth function  $g$  such that  $g'(\theta) \neq 0$ .

Combining Proposition 7.3 with Proposition 7.4 yields (take  $T_n = \hat{\theta}_n$ )

$$\sqrt{n}(g(\hat{\theta}_n) - g(\theta)) \rightarrow_d N(0, g'(\theta)^2 I(\theta)^{-1}).$$

Thus, for large  $n$ ,

$$g(\hat{\theta}_n) \sim_{\text{approx}} N(g(\theta), g'(\theta)^2 (n I(\theta))^{-1}).$$

Thus  $g(\hat{\theta}_n)$  is asymptotically unbiased for  $g(\theta)$  (unbiased in the long run) and its variance is approximately the information bound for unbiased estimators of  $g(\theta)$ .

**Constructing confidence sets for  $\theta$ :** Suppose that, for simplicity,  $\theta$  takes values in a subset of  $\mathbb{R}$ . Since,

$$\sqrt{n}(\hat{\theta}_n - \theta) \rightarrow_d N(0, I(\theta)^{-1}),$$

it follows that

$$\sqrt{n I(\theta)}(\hat{\theta}_n - \theta) \rightarrow_d N(0, 1).$$

Thus, the left side acts as an *approximate pivot* for  $\theta$ . We have,

$$\mathbb{P}_\theta \left( -z_{\alpha/2} \leq \sqrt{n I(\theta)}(\hat{\theta}_n - \theta) \leq z_{\alpha/2} \right) \approx 1 - \alpha.$$

An approximate level  $1 - \alpha$  confidence set for  $\theta$  is obtained as

$$\left\{ \theta : -z_{\alpha/2} \leq \sqrt{n I(\theta)}(\hat{\theta}_n - \theta) \leq z_{\alpha/2} \right\}.$$

To find the above confidence set, one needs to solve for all values of  $\theta$  satisfying the inequalities in the above display; this can however be a potentially complicated exercise depending on the functional form for  $I(\theta)$ .

However, if the sample size  $n$  is large,  $I(\hat{\theta}_n)$  can be expected to be close to  $I(\theta)$  with high probability and hence the following is also valid:

$$P_\theta \left[ -z_{\alpha/2} \leq \sqrt{n I(\hat{\theta}_n)}(\hat{\theta}_n - \theta) \leq z_{\alpha/2} \right] \approx 1 - \alpha. \quad (\star\star)$$

This immediately gives an approximate level  $1 - \alpha$  CI for  $\theta$  as:

$$\left[ \hat{\theta}_n - \frac{1}{\sqrt{n I(\hat{\theta}_n)}} z_{\alpha/2}, \hat{\theta}_n + \frac{1}{\sqrt{n I(\hat{\theta}_n)}} z_{\alpha/2} \right].$$

Let's see what this implies for the Bernoulli example discussed above. Recall that  $I(\theta) = (\theta(1 - \theta))^{-1}$  and  $\hat{\theta} = \bar{X}$ . The approximate  $(1 - \alpha)$ -CI is then given by,

$$\left[ \bar{X}_n - \sqrt{\frac{\bar{X}_n(1 - \bar{X}_n)}{n}} z_{\alpha/2}, \bar{X}_n + \sqrt{\frac{\bar{X}_n(1 - \bar{X}_n)}{n}} z_{\alpha/2} \right].$$

**Exercise:** Find explicitly

$$\left\{ \theta : -z_{\alpha/2} \leq \sqrt{n I(\theta)} (\hat{\theta}_n - \theta) \leq z_{\alpha/2} \right\}$$

in the following cases (a)  $X_1, X_2, \dots, X_n$  are i.i.d Bernoulli( $\theta$ ). (b)  $X_1, X_2, \dots, X_n$  are i.i.d Pois( $\theta$ ).

You will see that this involves solving for the roots of a quadratic equation. As in the Bernoulli example, one can also get an approximate CI for  $\theta$  in the Poisson setting on using (\*\*). Verify that this yields the following level  $1 - \alpha$  CI for  $\theta$ :

$$\left[ \bar{X}_n - \sqrt{\frac{\bar{X}_n}{n}} z_{\alpha/2}, \bar{X}_n + \sqrt{\frac{\bar{X}_n}{n}} z_{\alpha/2} \right] .$$

The recipe (\*\*) is somewhat unsatisfactory because it involves one more level of approximation in that  $I(\theta)$  is replaced by  $I(\hat{\theta})$  (note that there is already one level of approximation in that the pivots being considered are only approximately  $N(0, 1)$  by the CLT).

## 8 Bayesian paradigm

### Frequentist versus Bayesian statistics:

Frequentist:

- Data are a repeatable random sample — there is a frequency.
- *Parameters are fixed.*
- Underlying parameters remain constant during this repeatable process.

Bayesian:

- Parameters are unknown and described probabilistically.
- *Analysis is done conditioning on the observed data; i.e., data is treated as fixed.*

### 8.1 Prior distribution

**Definition 15** (Prior distribution). *Suppose that one has a statistical model with parameter  $\theta$ . If one treats  $\theta$  as random, then the distribution that one assigns to  $\theta$  before observing the data is called its **prior distribution**. Its pdf/pmf is called the **prior pdf/pmf** of  $\theta$ .*

Thus, now  $\theta$  is random and will be denoted by  $\Theta$  (note the change of notation).

**Example:** Let  $\Theta$  denote the probability of obtaining a head when a certain coin is tossed.

- Case 1: Suppose that it is known that the coin either is fair or has a head on each side. Then  $\Theta$  only takes two values, namely  $1/2$  and  $1$ . If the prior probability that the coin is fair is  $0.8$ , then the prior p.m.f of  $\Theta$  is  $\xi(1/2) = 0.8$  and  $\xi(1) = 0.2$ .
- Case 2: Suppose that  $\Theta$  can take any value between  $(0, 1)$  with a prior distribution given by a Beta distribution with parameters  $(1, 1)$ .

Suppose that the observable data  $X_1, X_2, \dots, X_n$  are modeled as random sample from a distribution indexed by  $\theta$ . Suppose  $f(\cdot|\theta)$  denote the p.m.f/p.d.f of a single random variable under the distribution indexed by  $\theta$ .

When we treat the unknown parameter  $\Theta$  as random, then the joint distribution of the observable random variables (i.e., data) indexed by  $\theta$  is understood as the **conditional distribution** of the data given  $\Theta = \theta$ .

Thus, in general we will have  $X_1, \dots, X_n | \Theta = \theta$  are i.i.d with p.d.f/p.m.f  $f(\cdot | \theta)$ , and that  $\Theta \sim \xi$ , i.e.,

$$f_n(\mathbf{x} | \theta) = f(x_1 | \theta) \cdots f(x_n | \theta),$$

where  $f_n$  is the joint conditional distribution of  $\mathbf{X} = (X_1, \dots, X_n)$  given  $\Theta = \theta$ .

## 8.2 Posterior distribution

**Definition 16** (Posterior distribution). *Consider a statistical inference problem with parameter  $\theta$  and random variables  $X_1, \dots, X_n$  to be observed. The conditional distribution of  $\Theta$  given  $X_1, \dots, X_n$  is called the **posterior distribution** of  $\theta$ . The conditional p.m.f/p.d.f of  $\Theta$  given  $X_1 = x_1, \dots, X_n = x_n$  is called the **posterior p.m.f/p.d.f** of  $\theta$  and is usually denoted by  $\xi(\cdot | x_1, \dots, x_n)$ .*

**Theorem 8.1.** *Suppose that the  $n$  random variables  $X_1, \dots, X_n$  form a random sample from a distribution for which the p.d.f/p.m.f is  $f(\cdot | \theta)$ . Suppose also that the value of the parameter  $\theta$  is unknown and the prior p.d.f/p.m.f of  $\theta$  is  $\xi(\cdot)$ . Then the posterior p.d.f/p.m.f of  $\theta$  is*

$$\xi(\theta | \mathbf{x}) = \frac{f(x_1 | \theta) \cdots f(x_n | \theta) \xi(\theta)}{g_n(\mathbf{x})}, \quad \text{for } \theta \in \Omega,$$

where  $g_n$  is the marginal joint p.d.f/p.m.f of  $X_1, \dots, X_n$ .

**Example 8.2** (Sampling from a Bernoulli distribution). Suppose that  $X_1, \dots, X_n$  form a random sample from the Bernoulli distribution with mean  $\theta > 0$ , where  $0 < \theta < 1$  is unknown. Suppose that the prior distribution of  $\Theta$  is Beta( $\alpha, \beta$ ), where  $\alpha, \beta > 0$ .

Then the posterior distribution of  $\Theta$  given  $X_i = x_i$ , for  $i = 1, \dots, n$ , is Beta( $\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i$ ).

*Proof.* The joint p.m.f of the data is

$$f_n(\mathbf{x} | \theta) = f(x_1 | \theta) \cdots f(x_n | \theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i}.$$

Therefore the posterior density of  $\Theta | X_1 = x_1, \dots, X_n = x_n$  is given by

$$\begin{aligned} \xi(\theta | \mathbf{x}) &\propto \theta^{\alpha-1} (1 - \theta)^{\beta-1} \cdot \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i} \\ &= \theta^{\sum_{i=1}^n x_i + \alpha - 1} (1 - \theta)^{\beta + n - \sum_{i=1}^n x_i - 1}, \end{aligned}$$

for  $\theta \in (0, 1)$ . Thus,  $\Theta | X_1 = x_1, \dots, X_n = x_n \sim \text{Beta}(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i)$ .  $\square$

**Example 8.3** (Sampling from a Poisson distribution). Suppose that  $X_1, \dots, X_n$  form a random sample from the Poisson distribution with mean  $\theta > 0$ , where  $\theta$  is unknown. Suppose that the prior distribution of  $\Theta$  is  $\text{Gamma}(\alpha, \beta)$ , where  $\alpha, \beta > 0$ . Show that the posterior distribution of  $\Theta$  given  $X_i = x_i$ , for  $i = 1, \dots, n$ , is  $\text{Gamma}(\alpha + \sum_{i=1}^n x_i, \beta + n)$ .

**Definition:** Let  $X_1, X_2, \dots$ , be conditionally i.i.d given  $\Theta = \theta$  with p.m.f/p.d.f  $f(\cdot|\theta)$ , where  $\theta \in \Omega$ . Let  $\Psi$  be a family of possible distributions over the parameter space  $\Omega$ . Suppose that no matter which prior distribution  $\xi$  we choose from  $\Psi$ , no matter how many observations  $\mathbf{X} = (X_1, \dots, X_n)$  we observe, and no matter what their observed values  $\mathbf{x} = (x_1, \dots, x_n)$  are, the posterior distribution  $\xi(\cdot|\mathbf{x})$  is a member of  $\Psi$ . Then  $\Psi$  is called a *conjugate family of prior distributions* for samples from the distributions  $f(\cdot|\theta)$ .

**Example 8.4** (Sampling from an Exponential distribution). Suppose that the distribution of the lifetime of fluorescent tubes of a certain type is the exponential distribution with parameter  $\theta$ . Suppose that  $X_1, \dots, X_n$  is a random sample of lamps of this type. Also suppose that  $\Theta \sim \text{Gamma}(\alpha, \beta)$ . Then

$$f_n(\mathbf{x}|\theta) = \prod_{i=1}^n \theta e^{-\theta x_i} = \theta^n e^{-\theta \sum_{i=1}^n x_i}.$$

Then the posterior distribution of  $\Theta$  given the data is

$$\xi(\theta|\mathbf{x}) \propto \theta^n e^{-\theta \sum_{i=1}^n x_i} \cdot \theta^{\alpha-1} e^{-\beta\theta} = \theta^{n+\alpha-1} e^{-(\beta + \sum_{i=1}^n x_i)\theta}.$$

Therefore,  $\Theta|\mathbf{X}_n = \mathbf{x} \sim \text{Gamma}(\alpha + n, \beta + \sum_{i=1}^n x_i)$ .

### 8.3 Bayes Estimators

**Definition:** A *loss function* is a real-valued function of two variables,  $L(\theta, a)$ , where  $\theta \in \Omega$  and  $a \in \mathbb{R}$ .

The interpretation is that the statistician loses  $L(\theta, a)$  if the parameter equals  $\theta$  and the estimate equals  $a$ .

**Example:** (Squared error loss)  $L(\theta, a) = (\theta - a)^2$ .

(Absolute error loss)  $L(\theta, a) = |\theta - a|$ .

Suppose that  $\Theta \sim \xi(\cdot)$  is a p.d.f/p.m.f. Consider the problem of estimating  $\Theta$  without being able to observe the data. If the statistician chooses a particular estimate  $a$ , then their expected loss will be

$$\mathbb{E}[L(\Theta, a)] = \int L(\theta, a) \xi(\theta) d\theta.$$

It is sensible that the statistician wishes to choose an estimate  $a$  for which the expected loss is *minimal*.

**Definition:** Suppose now that the statistician can observe the value  $\mathbf{x}$  of data  $\mathbf{X}_n$ , and let  $\xi(\cdot|\mathbf{x})$  denote the posterior p.d.f/p.m.f of  $\theta \in \Omega$ . For each estimate  $a$  that the statistician might use, their expected loss in this case will be

$$\mathbb{E}[L(\theta, a)|\mathbf{x}] = \int_{\Omega} L(\theta, a)\xi(\theta|\mathbf{x})d\theta. \quad (4)$$

For each possible value  $\mathbf{x}$  of  $\mathbf{X}_n$ , let  $\delta^*(\mathbf{x})$  denote a value of the estimate  $a$  for which the expected loss (4) is minimum. Then  $\delta^*(\mathbf{x})$  is called the **Bayes estimate** of  $\theta$ . Plugging in  $\mathbf{X}_n$  instead of  $\mathbf{x}$ , we obtain  $\delta^*(\mathbf{X}_n)$ , which is called the **Bayes estimator** of  $\theta$ .

Thus, a Bayes estimator is an estimator that is chosen to minimize the *posterior mean* of some measure of how far the estimator is from the parameter.

**Corollary 8.5.** *Let  $\theta \in \Omega \subset \mathbb{R}$ . Suppose that the squared error loss function is used and the posterior mean of  $\Theta$ , i.e.,  $\mathbb{E}(\Theta|\mathbf{X}_n)$ , is finite. Then the Bayes estimator of  $\theta$  is*

$$\delta^*(\mathbf{X}_n) = \mathbb{E}(\Theta|\mathbf{X}_n).$$

**Example 8.6** (Bernoulli distribution with Beta prior). Suppose that  $X_1, \dots, X_n$  form a random sample from the Bernoulli distribution with mean  $\theta > 0$ , where  $0 < \theta < 1$  is unknown. Suppose that the prior distribution of  $\Theta$  is Beta( $\alpha, \beta$ ), where  $\alpha, \beta > 0$ . Recall that  $\Theta|X_1 = x_1, \dots, X_n = x_n \sim \text{Beta}(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i)$ . Thus,

$$\delta^*(\mathbf{X}_n) = \frac{\alpha + \sum_{i=1}^n X_i}{\alpha + \beta + n}.$$

## 8.4 Sampling from a normal distribution

**Theorem 8.7.** *Suppose that  $X_1, \dots, X_n$  form a random sample from  $N(\theta, \sigma^2)$ , where  $\theta$  is unknown and the value of the variance  $\sigma^2 > 0$  is known. Suppose that  $\Theta \sim N(\mu_0, v_0^2)$ . Then*

$$\Theta|X_1 = x_1, \dots, X_n = x_n \sim N(\mu_1, v_1^2),$$

where

$$\mu_1 = \frac{\sigma^2 \mu_0 + n v_0^2 \bar{x}_n}{\sigma^2 + n v_0^2} \quad \text{and} \quad v_1^2 = \frac{\sigma^2 v_0^2}{\sigma^2 + n v_0^2}.$$

*Proof.* The joint density has the form

$$f_n(\mathbf{x}|\theta) \propto \exp \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2 \right].$$

The method of completing the squares tells us that

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

Thus, by omitting the factor that involves  $x_1, \dots, x_n$  but does not depend on  $\theta$ , we may rewrite  $f_n(\mathbf{x}|\theta)$  as

$$f_n(\mathbf{x}|\theta) \propto \exp \left[ -\frac{n}{2\sigma^2}(\theta - \bar{x}_n)^2 \right].$$

Since the prior density has the form

$$\xi(\theta) \propto \exp \left[ -\frac{1}{2v_0^2}(\theta - \mu_0)^2 \right],$$

it follows that the posterior p.d.f  $\xi(\theta|\mathbf{x})$  satisfies

$$\xi(\theta|\mathbf{x}) \propto \exp \left[ -\frac{n}{2\sigma^2}(\theta - \bar{x}_n)^2 - \frac{1}{2v_0^2}(\theta - \mu_0)^2 \right].$$

Completing the squares again establishes the following identity:

$$\frac{n}{\sigma^2}(\theta - \bar{x}_n)^2 + \frac{1}{v_0^2}(\theta - \mu_0)^2 = \frac{1}{v_1^2}(\theta - \mu_1)^2 + \frac{n}{\sigma^2 + nv_0^2}(\bar{x}_n - \mu_0)^2.$$

The last term on the right side does not involve on  $\theta$ . Thus,

$$\xi(\theta|\mathbf{x}) \propto \exp \left[ -\frac{1}{2v_1^2}(\theta - \mu_1)^2 \right].$$

□

Thus, the Bayes estimator (under the squared error loss) in this problem is

$$\delta^*(\mathbf{X}_n) = \frac{\sigma^2\mu_0 + nv_0^2\bar{X}_n}{\sigma^2 + nv_0^2}.$$



## 9 Hypothesis Testing

### 9.1 Principles of Hypothesis Testing

We are given data (say  $X_1, \dots, X_n$  i.i.d  $P_\theta$ ) from a model that is parametrized by  $\theta$ . We consider a statistical problem involving  $\theta$  whose value is unknown but must lie in a certain space  $\Omega$ . We consider the testing problem

$$H_0 : \theta \in \Omega_0 \quad \text{versus} \quad H_1 : \theta \in \Omega_1, \quad (5)$$

where  $\Omega_0 \cap \Omega_1 = \emptyset$  and  $\Omega_0 \cup \Omega_1 = \Omega$ .

Here the hypothesis  $H_0$  is called the **null hypothesis** and  $H_1$  is called the **alternative hypothesis**.

**Question:** Is there enough evidence in the data against the null hypothesis (in which case we reject it) or should we continue to stick to it?

Such questions arise very naturally in many different fields of application.

---

**Definition 17** (One-sided and two-sided hypotheses). *Let  $\theta$  be a one-dimensional parameter.*

- *one-sided hypotheses*
  - $H_0 : \theta \leq \theta_0$ , and  $H_1 : \theta > \theta_0$ , or
  - $H_0 : \theta \geq \theta_0$ , and  $H_1 : \theta < \theta_0$
- *two-sided hypotheses*  $H_0 : \theta = \theta_0$ , and  $H_1 : \theta \neq \theta_0$ .

$H_0$  is *simple* if  $\Omega_0$  is a set with only one point; otherwise,  $H_0$  is *composite*.

---

**Testing for a normal mean:** Suppose that  $X_1, X_2, \dots, X_n$  is a sample from a  $N(\mu, \sigma^2)$  distribution and let, initially,  $\sigma^2$  be known.

We want to test the *null hypothesis*  $H_0 : \mu = \mu_0$  against the alternative  $H_1 : \mu \neq \mu_0$ .

**Example:** For concreteness,  $X_1, X_2, \dots, X_n$  could be the heights of  $n$  individuals in some tribal population. The distribution of heights in a (homogeneous) population is usually normal, so that a  $N(\mu, \sigma^2)$  model is appropriate. If we have some a-priori reason to believe that the average height in this population is around 60 inches, we could postulate a null hypothesis of the form  $H_0 : \mu = \mu_0 \equiv 60$ ; the alternative hypothesis is  $H_1 : \mu \neq 60$ .

---

## 9.2 Critical regions and test statistics

Consider a problem in which we wish to test the following hypotheses:

$$H_0 : \theta \in \Omega_0, \quad \text{and} \quad H_1 : \theta \in \Omega_1. \quad (6)$$

**Question:** How do we do the test?

The statistician must decide, after observing data, which of the hypothesis  $H_0$  or  $H_1$  appears to be true.

A procedure for deciding which hypothesis to choose is called a **test procedure** of simply a **test**. We will denote a test by  $\delta$ .

Suppose we can observe a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  drawn from a distribution that involves the unknown parameter  $\theta$ , e.g., suppose that  $X_1, \dots, X_n$  are i.i.d  $P_\theta$ ,  $\theta \in \Omega$ .

Let  $S$  denote the set of all possible values of the  $n$ -dimensional random vector  $\mathbf{X}$ .

We specify a test procedure by partitioning  $S$  into two subsets:  $S = S_0 \cup S_1$ , where

- the **critical region**  $S_1$  contains the values of  $\mathbf{X}$  for which we will reject  $H_0$ , and
- the other subset  $S_0$  (usually called the **acceptance region**) contains the values of  $\mathbf{X}$  for which we will not reject  $H_0$ .

A test procedure is determined by specifying the critical region  $S_1$  of the test.

In most hypothesis-testing problems, the critical region is defined in terms of a statistic,  $T = \varphi(\mathbf{X})$ .

---

**Definition 18** (Test statistic/rejection region). *Let  $\mathbf{X}$  be a random sample from a distribution that depends on a parameter  $\theta$ . Let  $T = \varphi(\mathbf{X})$  be a statistic, and let  $R$  be a subset of the real line. Suppose that a test procedure is of the form:*

$$\text{reject } H_0 \quad \text{if} \quad T \in R.$$

*Then we call  $T$  a **test statistic**, and we call  $R$  the **rejection region** of the test, and the critical region reduces to*

$$S_1 = \{\mathbf{x} : \varphi(\mathbf{x}) \in R\}.$$

Typically, the rejection region for a test based on a test statistic  $T$  will be some fixed interval or the complement of some fixed interval.

If the test rejects  $H_0$  when  $T \geq c$ , the rejection region is the interval  $[c, \infty)$ . Indeed, most of the tests can be written in the form:

$$\text{reject } H_0 \quad \text{if} \quad T \geq c.$$

---

**Example:** Suppose that  $X_1, \dots, X_n$  are i.i.d  $N(\mu, \sigma^2)$  where  $\mu \in \mathbb{R}$  is unknown, and  $\sigma > 0$  is assumed *known*.

Suppose that we want to test  $H_0 : \mu = \mu_0$  versus  $H_1 : \mu \neq \mu_0$ .

Some of these procedures can be justified using formal paradigms. Under the null hypothesis the  $X_i$ 's are i.i.d  $N(\mu_0, \sigma^2)$  and the sample mean  $\bar{X}$  follows  $N(\mu_0, \sigma^2/n)$ .

Thus, it is reasonable to take  $T = \varphi(\mathbf{X}) = |\bar{X} - \mu_0|$ .

Large deviations of the observed value of  $\bar{X}$  from  $\mu_0$  would lead us to suspect that the null hypothesis might not be true.

Thus, a reasonable test can be to reject  $H_0$  if  $T = |\bar{X} - \mu_0| > c$ , for some “large” constant  $c$ .

**But how large is large?** We will discuss this soon...

---

Associated with the test procedure  $\delta$  are two different kinds of error that we can commit. These are called *Type 1 error* and *Type 2 error* (Draw the  $2 \times 2$  table!).

Decision	Fail to reject $H_0$	Reject $H_0$
State		
$H_0$ True	Correct	Type 1 error
$H_1$ True	Type 2 error	Correct

Table 1: Hypothesis test.

**Type 1 error** occurs if we reject the null hypothesis when actually  $H_0$  is true.

**Type 2 error** occurs if we do not reject the null hypothesis when actually  $H_0$  is false.

---

## 9.3 Power function and types of error

Let  $\delta$  be a test procedure. If  $S_1$  denotes the critical region of  $\delta$ , then the **power function** of the test  $\delta$ ,  $\pi(\theta|\delta)$ , is defined by the relation

$$\pi(\theta|\delta) = \mathbb{P}_\theta(\mathbf{X} \in S_1) \quad \text{for} \quad \theta \in \Omega.$$

Thus, the power function  $\pi(\theta|\delta)$  specifies for each possible value of  $\theta$ , the *probability that  $\delta$  will reject  $H_0$* . If  $\delta$  is described in terms of a test statistic  $T$  and rejection region  $R$ , the power function is

$$\pi(\theta|\delta) = \mathbb{P}_\theta(T \in R) \quad \text{for } \theta \in \Omega.$$

**Example:** Suppose that  $X_1, \dots, X_n$  are i.i.d Uniform( $0, \theta$ ), where  $\theta > 0$  is unknown.

Suppose that we are interested in the following hypotheses:

$$H_0 : 2.9 \leq \theta \leq 4, \quad \text{versus} \quad H_1 : \theta < 2.9, \text{ or } \theta > 4.$$

We know that the MLE of  $\theta$  is  $X_{(n)} = \max\{X_1, \dots, X_n\}$ .

Note that  $X_{(n)} < \theta$ .

Suppose that we use a test  $\delta$  given by the critical region

$$S_1 = \{\mathbf{x} \in \mathbb{R}^n : x_{(n)} \leq 2.9 \text{ or } x_{(n)} \geq 4\}.$$

**Question:** Find the power function  $\pi(\theta|\delta)$ ?

**Solution:** The power function of  $\delta$  is

$$\pi(\theta|\delta) = \mathbb{P}_\theta(X_{(n)} \leq 2.9 \text{ or } X_{(n)} > 4) = \mathbb{P}_\theta(X_{(n)} \leq 2.9) + \mathbb{P}_\theta(X_{(n)} \geq 4).$$

Case (i): Suppose that  $\theta \leq 2.9$ . Then

$$\pi(\theta|\delta) = \mathbb{P}_\theta(X_{(n)} \leq 2.9) = 1.$$

Case (ii): Suppose that  $2.9 < \theta < 4$ . Then

$$\pi(\theta|\delta) = \mathbb{P}_\theta(X_{(n)} \leq 2.9) = \left(\frac{2.9}{\theta}\right)^n.$$

Case (iii): Suppose that  $\theta > 4$ . Then

$$\pi(\theta|\delta) = \left(\frac{2.9}{\theta}\right)^n + \left[1 - \left(\frac{4}{\theta}\right)^n\right].$$

The ideal power function would be one for which

- $\pi(\theta|\delta) = 0$  for every value of  $\theta \in \Omega_0$ , and

- $\pi(\theta|\delta) = 1$  for every value of  $\theta \in \Omega_1$ .

If the power function of a test  $\delta$  actually had these values, then regardless of the actual value of  $\theta$ ,  $\delta$  would lead to the correct decision with probability 1.

In a practical problem, however, there would seldom exist any test procedure having this ideal power function.

- Type-I error: rejecting  $H_0$  given that  $\theta \in \Omega_0$ . It occurs with probability  $\pi(\theta|\delta)$ .
- Type-II error: not rejecting  $H_0$  given that  $\theta \in \Omega_1$ . It occurs with probability  $1 - \pi(\theta|\delta)$ .

Ideal goals: we would like the power function  $\pi(\theta|\delta)$  to be **low** for values of  $\theta \in \Omega_0$ , and **high** for  $\theta \in \Omega_1$ .

Generally, these two goals work against each other. That is, if we choose  $\delta$  to make  $\pi(\theta|\delta)$  small for  $\theta \in \Omega_0$ , we will usually find that  $\pi(\theta|\delta)$  is small for  $\theta \in \Omega_1$  as well.

Examples:

- The test procedure  $\delta_0$  that never rejects  $H_0$ , regardless of what data are observed, will have  $\pi(\theta|\delta_0) = 0$  for all  $\theta \in \Omega_0$ . However, for this procedure  $\pi(\theta|\delta_0) = 0$  for all  $\theta \in \Omega_1$  as well.
- Similarly, the test  $\delta_1$  that always rejects  $H_0$  will have  $\pi(\theta|\delta_1) = 1$  for all  $\theta \in \Omega_1$ , but it will also have  $\pi(\theta|\delta_1) = 1$  for all  $\theta \in \Omega_0$ .

Hence, there is a need to strike an appropriate balance between the two goals of

*low power in  $\Omega_0$  and high power in  $\Omega_1$ .*

1. The most popular method for striking a balance between the two goals is to choose a number  $\alpha_0 \in (0, 1)$  and require that

$$\pi(\theta|\delta) \leq \alpha_0, \quad \text{for all } \theta \in \Omega_0. \quad (7)$$

This  $\alpha_0$  will usually be a small positive fraction (historically .05 or .01) and will be called the **level of significance** or simply *level*.

Then, among all tests that satisfy (7), the statistician seeks a test whose power function is as high as can be obtained for  $\theta \in \Omega_1$ .

2. Another method of balancing the probabilities of type I and type II errors is to minimize a linear combination of the different probabilities of error.

## 9.4 Significance level

**Definition 19** (level/size). *(of the test)*

- A test that satisfies (7) is called a level  $\alpha_0$  test, and we say that the test has level of significance  $\alpha_0$ .
- The size  $\alpha(\delta)$  of a test  $\delta$  is defined as follows:

$$\alpha(\delta) = \sup_{\theta \in \Omega_0} \pi(\theta|\delta).$$

It follows from Definition 19 that:

- A test  $\delta$  is a level  $\alpha_0$  test iff  $\alpha(\delta) \leq \alpha_0$ .
- If the null hypothesis is simple (that is,  $H_0 : \theta = \theta_0$ ), then  $\alpha(\delta) = \pi(\theta_0|\delta)$ .

### Making a test have a specific significance level

Suppose that we wish to test the hypotheses

$$H_0 : \theta \in \Omega_0, \quad \text{versus} \quad H_1 : \theta \in \Omega_1.$$

Let  $T$  be a test statistic, and suppose that our test will reject the null hypothesis if  $T \geq c$ , for some constant  $c$ . Suppose also that we desire our test to have the level of significance  $\alpha_0$ . The power function of our test is  $\pi(\theta|\delta) = \mathbb{P}_\theta(T \geq c)$ , and we want that

$$\sup_{\theta \in \Omega_0} \mathbb{P}_\theta(T \geq c) \leq \alpha_0. \tag{8}$$

---

#### Remarks:

1. It is clear that the power function, and hence the left side of (8), are non-increasing functions of  $c$ .  
Hence, (8) will be satisfied for large values of  $c$ , but not for small values.  
If  $T$  has a continuous distribution, then it is usually simple to find an appropriate  $c$ .
2. Whenever we choose a test procedure, we should also examine the power function. If one has made a good choice, then the power function should generally be larger for  $\theta \in \Omega_1$  than for  $\theta \in \Omega_0$ .

---

**Example:** Suppose that  $X_1, \dots, X_n$  are i.i.d  $N(\mu, \sigma^2)$  where  $\mu \in \mathbb{R}$  is unknown, and  $\sigma > 0$  is assumed *known*. We want to test  $H_0 : \mu = \mu_0$  versus  $H_1 : \mu \neq \mu_0$ .

Suppose that the null hypothesis  $H_0$  is true.

If the variance of the sample mean is, say, 100, a deviation of  $\bar{X}$  from  $\mu_0$  by 15 is not really unusual.

On the other hand if the variance is 10, then a deviation of the sample mean from  $\mu_0$  by 15 is really sensational.

Thus the quantity  $|\bar{X} - \mu_0|$  in itself is not sufficient to formulate a decision regarding rejection of the null hypothesis.

We need to adjust for the underlying variance. This is done by computing the so-called  $z$ -statistic,

$$Z := \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \equiv \frac{\sqrt{n}(\bar{X} - \mu_0)}{\sigma}$$

and rejecting the null hypothesis for large absolute values of this statistic.

Under the null hypothesis  $Z$  follows  $N(0, 1)$ ; thus an absolute  $Z$ -value of 3.5 is quite unlikely. Therefore if we observe an absolute  $Z$ -value of 3.5 we might rule in favor of the alternative hypothesis.

You can see now that we need a threshold value, or in other words a critical point such that if the  $Z$ -value exceeds that point we reject. Our test procedure  $\delta$  then looks like,

$$\text{reject } H_0 \quad \text{if} \quad \left| \frac{\sqrt{n}(\bar{X} - \mu_0)}{\sigma} \right| > c_{n, \alpha_0}$$

where  $c_{n, \alpha_0}$  is the *critical value* and will depend on  $\alpha_0$  which is the tolerance for the Type 1 error, i.e., the level that we set beforehand.

The quantity  $c_{n, \alpha_0}$  is determined using the relation

$$\mathbb{P}_{\mu_0} \left( \left| \frac{\sqrt{n}(\bar{X} - \mu_0)}{\sigma} \right| > c_{n, \alpha_0} \right) = \alpha_0.$$

Straightforward algebra then yields that

$$P_{\mu_0} \left( -c_{n, \alpha_0} \frac{\sigma}{\sqrt{n}} \leq \bar{X} - \mu_0 \leq c_{n, \alpha_0} \frac{\sigma}{\sqrt{n}} \right) = 1 - \alpha_0,$$

whence we can choose  $c_{n, \alpha_0} = z_{\alpha_0/2}$ , the  $\frac{\alpha_0}{2}$ -th quantile of the  $N(0, 1)$  distribution.

The acceptance region  $\mathcal{A}$  (or  $S_0$ ) for the null hypothesis is therefore

$$\mathcal{A} = \left\{ \mathbf{x} = (x_1, x_2, \dots, x_n) : \mu_0 - \frac{\sigma}{\sqrt{n}} z_{\alpha_0/2} \leq \bar{x} \leq \mu_0 + \frac{\sigma}{\sqrt{n}} z_{\alpha_0/2} \right\}.$$

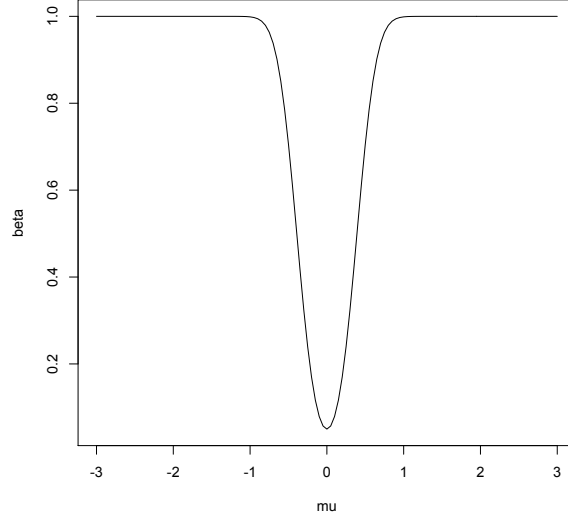


Figure 4: The power function  $\pi(\mu|\delta)$  for  $\mu_0 = 0$ ,  $\sigma = 1$  and  $n = 25$ .

So we accept whenever  $\bar{X}$  lies in a certain window of  $\mu_0$ , the postulated value under the null, and reject otherwise which is in accordance with intuition.

The length of the window is determined by the tolerance level  $\alpha_0$ , the underlying variance  $\sigma^2$  and of course the sample size  $n$ .

## 9.5 *P*-value

The ***p*-value** is the smallest level  $\alpha_0$  such that we would reject  $H_0$  at level  $\alpha_0$  with the observed data.

For this reason, the *p*-value is also called the *observed level of significance*.

Example: If the observed value of  $Z$  was 2.78, and that the corresponding *p*-value = 0.0054. It is then said that the observed value of  $Z$  is just significant at the level of significance 0.0054.

### **Advantages:**

1. No need to select beforehand an arbitrary level of significance  $\alpha_0$  at which to carry out the test.
2. When we learn that the observed value of  $Z$  was just significant at the level of significance 0.0054, we immediately know that  $H_0$  would be rejected for every larger value of  $\alpha_0$  and would not be rejected for any smaller value.



## 9.6 Testing simple hypotheses: optimal tests

Let the random vector  $\mathbf{X} = (X_1, \dots, X_n)$  come from a distribution for which the joint p.m.f/p.d.f is either  $f_0(\mathbf{x})$  or  $f_1(\mathbf{x})$ . Let  $\Omega = \{\theta_0, \theta_1\}$ . Then,

- $\theta = \theta_0$  stands for the case in which the data have p.m.f/p.d.f  $f_0(\mathbf{x})$ ,
- $\theta = \theta_1$  stands for the case in which the data have p.m.f/p.d.f  $f_1(\mathbf{x})$ .

We are then interested in testing the following simple hypotheses:

$$H_0 : \theta = \theta_0 \quad \text{versus} \quad H_1 : \theta = \theta_1.$$

In this case, we have special notation for the probabilities of type I and type II errors:

$$\begin{aligned} \alpha(\delta) &= \mathbb{P}_{\theta_0}(\text{Rejecting } H_0), \\ \beta(\delta) &= \mathbb{P}_{\theta_1}(\text{Not rejecting } H_0). \end{aligned}$$

### 9.6.1 Minimizing the $\mathbb{P}(\text{Type-II error})$

Suppose that the probability  $\alpha(\delta)$  of an error of type I is not permitted to be greater than a specified level of significance, and it is desired to find a procedure  $\delta$  for which  $\beta(\delta)$  will be a minimum.

**Theorem 9.1** (Neyman-Pearson lemma). *Suppose that  $\delta'$  is a test procedure that has the following form for some constant  $k > 0$ :*

- $H_0$  is not rejected if  $f_1(\mathbf{x}) < kf_0(\mathbf{x})$ ,
- $H_0$  is rejected if  $f_1(\mathbf{x}) > kf_0(\mathbf{x})$ , and
- $H_0$  can be either rejected or not if  $f_1(\mathbf{x}) = kf_0(\mathbf{x})$ .

Let  $\delta$  be another test procedure. Then,

$$\begin{aligned} \text{if } \alpha(\delta) &\leq \alpha(\delta'), \quad \text{then it follows that } \beta(\delta) \geq \beta(\delta') \\ \text{if } \alpha(\delta) &< \alpha(\delta'), \quad \text{then it follows that } \beta(\delta) > \beta(\delta'). \end{aligned}$$

**Example:** Suppose that  $\mathbf{X} = (X_1, \dots, X_n)$  is a random sample from the normal distribution with unknown mean  $\theta$  and known variance 1. We are interested in testing:

$$H_0 : \theta = 0 \quad \text{versus} \quad H_1 : \theta = 1.$$

We want to find a test procedure for which  $\beta(\delta)$  will be a minimum among all test procedures for which  $\alpha(\delta) \leq 0.05$ .

We have,

$$f_0(\mathbf{x}) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n x_i^2\right) \quad \text{and} \quad f_1(\mathbf{x}) = \frac{1}{(2\pi)^{n/2}} \exp\left[-\frac{1}{2} \sum_{i=1}^n (x_i - 1)^2\right].$$

After some algebra, the likelihood ratio  $f_1(\mathbf{x})/f_0(\mathbf{x})$  can be written in the form

$$\frac{f_1(\mathbf{x})}{f_0(\mathbf{x})} = \exp\left[n\left(\bar{x} - \frac{1}{2}\right)\right].$$

Thus, rejecting  $H_0$  when the likelihood ratio is greater than a specified positive constant  $k$  is equivalent to rejecting  $H_0$  when the sample mean  $\bar{X}$  is greater than  $k' := 1/2 + \log k/n$ , another constant. Thus, we want to find,  $k'$  such that

$$\mathbb{P}_0(\bar{X} > k') = 0.05.$$

Now,

$$\begin{aligned} \mathbb{P}_0(\bar{X} > k') &= \mathbb{P}_0(\sqrt{n}\bar{X} > \sqrt{nk'}) = \mathbb{P}_0(Z > \sqrt{nk'}) = 0.05 \\ \Rightarrow \sqrt{nk'} &= 1.645. \end{aligned}$$

## 9.7 Uniformly most powerful (UMP) tests

We suppose that  $\Omega_0$  and  $\Omega_1$  are disjoint subsets of  $\Omega$ , and the hypotheses to be tested are

$$H_0 : \theta \in \Omega_0 \quad \text{versus} \quad H_1 : \theta \in \Omega_1. \quad (9)$$

- The subset  $\Omega_1$  contains at least two distinct values of  $\theta$ , in which case the alternative hypothesis  $H_1$  is composite.
- The null hypothesis  $H_0$  may be either simple or composite.

We consider *only* procedures in which

$$\mathbb{P}_\theta(\text{Rejecting } H_0) \leq \alpha_0 \quad \forall \theta \in \Omega_0.$$

that is

$$\pi(\theta|\delta) \leq \alpha_0 \quad \forall \theta \in \Omega_0$$

or

$$\alpha(\delta) \leq \alpha_0. \quad (10)$$

---

Finally, among all test procedures that satisfy the requirement (10), we want to find one such that

- the probability of type II error is as small as possible for every  $\theta \in \Omega_1$ , or
- the value of  $\pi(\theta|\delta)$  is as large as possible for every value of  $\theta \in \Omega_1$ .

There might be no single test procedure  $\delta$  that maximizes the power function  $\pi(\theta|\delta)$  simultaneously for every value of  $\theta \in \Omega_1$ .

In some problems, however, there will exist a test procedure that satisfies this criterion. Such a procedure, when it exists, is called a UMP test.

**Definition 20** (Uniformly most powerful (UMP) test). *A test procedure  $\delta^*$  is a uniformly most powerful (UMP) test of the hypotheses (9) at the level of significance  $\alpha_0$  if*

$$\alpha(\delta^*) \leq \alpha_0$$

*and, for every other test procedure  $\delta$  such that  $\alpha(\delta) \leq \alpha_0$ , it is true that*

$$\pi(\theta|\delta) \leq \pi(\theta|\delta^*)$$

*for every value of  $\theta \in \Omega_1$ .*

Usually no test will uniformly most powerful against ALL alternatives, except in the special case of “monotone likelihood ratio” (MLR).

*Example:* Suppose that  $X_1, \dots, X_n$  form a random sample from a normal distribution for which the mean  $\mu$  (unknown) and the variance  $\sigma^2$  (known). Consider testing  $H_0 : \mu = \mu_0$  versus  $H_1 : \mu \neq \mu_0$ . Even in this simple example, there is no UMP test.

## 9.8 The $t$ -test

### 9.8.1 Testing hypotheses about the mean with unknown variance

- Problem: testing hypotheses about the **mean** of a normal distribution when both the mean and the variance are unknown.
- The random variables  $X_1, \dots, X_n$  form a random sample from a normal distribution for which the mean  $\mu$  and the variance  $\sigma^2$  are unknown.
- The parameter space  $\Omega$  in this problem comprises every two-dimensional vector  $(\mu, \sigma^2)$ , where  $-\infty < \mu < \infty$  and  $\sigma^2 > 0$ .
- $H_0 : \mu = \mu_0$  versus  $H_1 : \mu \neq \mu_0$
- Define

$$U_n = \frac{\bar{X}_n - \mu_0}{s_n/\sqrt{n}}, \tag{11}$$

$$\text{where } s_n = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2}.$$

- We reject  $H_0$  if

$$|U_n| \geq T_{n-1}^{-1} \left( 1 - \frac{\alpha_0}{2} \right),$$

the  $(1 - \alpha_0/2)$ -quantile of the  $t$ -distribution with  $n - 1$  degrees of freedom and  $U_n$  is defined in (11).

- $p$ -values for  $t$ -tests: The  $p$ -value from the observed data and a specific test is the smallest  $\alpha_0$  such that we would reject the null hypothesis at level of significance  $\alpha_0$ .

Let  $u$  be the observed value of the statistic  $U_n$ . Thus the  $p$ -value of the test is

$$\mathbb{P}(|U_n| > |u|),$$

where  $U_n \sim T_{n-1}$ , under  $H_0$ .

- The  $p$ -value is  $2[1 - T_{n-1}(|u|)]$ , where  $u$  be the observed value of the statistic  $U_n$ .

**Example 9.2** (Example 9.5.9, p.582). We take  $\alpha_0 = 0.05$  and test

$$H_0 : \mu = 140 \quad \text{vs.} \quad H_1 : \mu \neq 140.$$

Furthermore we assume that we have iid.  $N(\mu, \sigma^2)$  samples, where both  $\mu$  and  $\sigma^2$  are unknown. We also have  $n = 30$ ,  $\bar{X}_{30} = 131.37$ ,  $s_{30} = 5.129$ . We compute

$$U_{30} = \sqrt{30} \left( \frac{131.37 - 140}{5.129} \right) = -9.219$$

and compare it to

$$T_{29}^{-1}(0.975) = 2.045.$$

In conclusion we reject  $H_0$  at 0.05-level.

## The Complete power function

Before we study the case when  $\sigma > 0$  is unknown, let us go back to the case when  $\sigma$  is known.

Our test  $\delta$  is “reject  $H_0$  if  $\left| \frac{\sqrt{n}(\bar{X} - \mu_0)}{\sigma} \right| > z_{\alpha/2}$ ”.

Thus we have,

$$\pi(\mu|\delta) = \mathbb{P}_\mu \left( \left| \frac{\sqrt{n}(\bar{X} - \mu_0)}{\sigma} \right| > z_{\alpha/2} \right),$$

which is just,

$$\mathbb{P}_\mu \left( \left| \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} + \frac{\sqrt{n}(\mu - \mu_0)}{\sigma} \right| > z_{\alpha/2} \right).$$

But when  $\mu$  is the population mean,  $\sqrt{n}(\bar{X} - \mu)/\sigma$  is  $N(0, 1)$ . If  $Z$  denotes a  $N(0, 1)$  variable then,

$$\begin{aligned} \pi(\mu|\delta) &= \mathbb{P}_\mu \left( \left| Z + \frac{\sqrt{n}(\mu - \mu_0)}{\sigma} \right| > z_{\alpha/2} \right) \\ &= \mathbb{P}_\mu \left( Z + \frac{\sqrt{n}(\mu - \mu_0)}{\sigma} > z_{\alpha/2} \right) + \mathbb{P} \left( Z + \frac{\sqrt{n}(\mu - \mu_0)}{\sigma} < -z_{\alpha/2} \right) \\ &= 1 - \Phi \left( z_{\alpha/2} - \frac{\sqrt{n}(\mu - \mu_0)}{\sigma} \right) + \Phi \left( -z_{\alpha/2} - \frac{\sqrt{n}(\mu - \mu_0)}{\sigma} \right) \\ &= \Phi \left( -z_{\alpha/2} + \frac{\sqrt{n}(\mu - \mu_0)}{\sigma} \right) + \Phi \left( -z_{\alpha/2} - \frac{\sqrt{n}(\mu - \mu_0)}{\sigma} \right). \end{aligned}$$

Check from the above calculations that  $\pi(\mu_0|\delta) = \alpha$ , the level of the test  $\delta$ .

Notice that the test function  $\delta$  depends on the value  $\mu_0$  under the null but it does not depend on any value in the alternative.

It is easy to check that  $\pi(\mu|\delta)$  diverges to 1 as  $\mu$  diverges to  $\infty$  or  $-\infty$ .

Moreover the power function is symmetric around  $\mu_0$ . In other words,  $\pi(\mu_0 + \Delta|\delta) = \pi(\mu_0 - \Delta|\delta)$  where  $\Delta > 0$ .

To see this, note that

$$\pi(\mu_0 + \Delta|\delta) = \Phi \left( -z_{\alpha/2} + \frac{\sqrt{n}\Delta}{\sigma} \right) + \Phi \left( -z_{\alpha/2} - \frac{\sqrt{n}\Delta}{\sigma} \right).$$

Check that you get the same expression for  $\pi(\mu_0 - \Delta|\delta)$ .

**Question:** What happens when  $\sigma > 0$  is unknown?

We can rewrite  $U_n$  as

$$U_n = \frac{\sqrt{n}(\bar{X}_n - \mu_0)/\sigma}{s_n/\sigma},$$

- The numerator has the normal distribution with mean  $\sqrt{n}(\mu - \mu_0)/\sigma$  and variance 1.
- The denominator is the square-root of a  $\chi^2$ -random variable divided by its degrees of freedom,  $n - 1$ .

- When the mean of the numerator is not 0,  $U_n$  has a *non-central  $t$ -distribution*.

**Definition 21** (Noncentral  $t$ -distributions). *Let  $W$  and  $Y_m$  be independent random variables  $W \sim \mathcal{N}(\psi, 1)$  and  $Y \sim \chi_m^2$ . Then the distribution of*

$$X := \frac{W}{\sqrt{Y_m/m}}$$

*is called the **non-central  $t$ -distribution** with  $m$  degrees of freedom and non-centrality parameter  $\psi$ . We define*

$$T_m(t|\psi) = \mathbb{P}(X \leq t)$$

*as the c.d.f of this distribution.*

- The non-central  $t$ -distribution with  $m$  degrees of freedom and non-centrality parameter  $\psi = 0$  is also the  $t$ -distribution with  $m$  degrees of freedom.
- The distribution of the statistic  $U_n$  in (11) is the non-central  $t$ -distribution with  $n - 1$  degrees of freedom and non-centrality parameter

$$\psi := \sqrt{n} \frac{(\mu - \mu_0)}{\sigma}.$$

- The power function of  $\delta$  (see Figure 9.14) is

$$\pi(\mu, \sigma^2|\delta) = T_{n-1}(-c|\psi) + 1 - T_{n-1}(c|\psi),$$

where  $c := T_{n-1}^{-1}(1 - \alpha_0/2)$ .

### 9.8.2 One-sided alternatives

We consider testing the following hypotheses:

$$H_0 : \mu \leq \mu_0, \quad \text{versus} \quad H_1 : \mu > \mu_0. \quad (12)$$

- When  $\mu = \mu_0$ ,  $U_n \sim t_{n-1}$ , regardless of the value of  $\sigma^2$ .
- The test rejects  $H_0$  if

$$U_n \geq c,$$

where  $c := T_{n-1}^{-1}(1 - \alpha_0)$  (the  $(1 - \alpha_0)$ -quantile) of the  $t$ -distribution with  $n - 1$  degrees of freedom.

- $\pi(\mu, \sigma^2|\delta) = 1 - T_{n-1}(c|\psi)$ .

### Power function of the $t$ -test

Let  $\delta$  be the test that rejects  $H_0$  in (12) if  $U_n \geq c$ .

The  $p$ -value for the hypotheses in (12) is  $1 - T_{n-1}(u)$ , where  $u$  is the observed value of the statistic  $U_n$ .

The power function  $\pi(\mu, \sigma^2|\delta)$  has the following properties:

1.  $\pi(\mu, \sigma^2|\delta) = \alpha_0$  when  $\mu = \mu_0$ ,
2.  $\pi(\mu, \sigma^2|\delta) < \alpha_0$  when  $\mu < \mu_0$ ,
3.  $\pi(\mu, \sigma^2|\delta) > \alpha_0$  when  $\mu > \mu_0$ ,
4.  $\pi(\mu, \sigma^2|\delta) \rightarrow 0$  as  $\mu \rightarrow -\infty$ ,
5.  $\pi(\mu, \sigma^2|\delta) \rightarrow 1$  as  $\mu \rightarrow \infty$ ,
6.  $\sup_{\theta \in \Omega_0} \pi(\theta|\delta) = \alpha_0$ .

---

When we want to test

$$H_0 : \mu \geq \mu_0 \quad \text{versus} \quad H_1 : \mu < \mu_0. \quad (13)$$

the test rejects  $H_0$  if  $U_n \leq c$ , where  $c = T_{n-1}^{-1}(\alpha_0)$  (the  $\alpha_0$ -quantile) of the  $t$ -distribution with  $n - 1$  degrees of freedom.

### Power function of the $t$ test

Let  $\delta$  be the test that rejects  $H_0$  in (13) if  $U_n \leq c$ .

The  $p$ -value for the hypotheses in (13) is  $T_{n-1}(u)$ . Observe that  $\pi(\mu, \sigma^2|\delta) = T_{n-1}(c|\psi)$ .

The power function  $\pi(\mu, \sigma^2|\delta)$  has the following properties:

1.  $\pi(\mu, \sigma^2|\delta) = \alpha_0$  when  $\mu = \mu_0$ ,
2.  $\pi(\mu, \sigma^2|\delta) > \alpha_0$  when  $\mu < \mu_0$ ,
3.  $\pi(\mu, \sigma^2|\delta) < \alpha_0$  when  $\mu > \mu_0$ ,
4.  $\pi(\mu, \sigma^2|\delta) \rightarrow 1$  as  $\mu \rightarrow -\infty$ ,
5.  $\pi(\mu, \sigma^2|\delta) \rightarrow 0$  as  $\mu \rightarrow \infty$ ,
6.  $\sup_{\theta \in \Omega_0} \pi(\theta|\delta) = \alpha_0$ .

## 9.9 Comparing the means of two normal distributions (two-sample $t$ test)

### 9.9.1 One-sided alternatives

Random samples are available from **two** normal distributions with common unknown variance  $\sigma^2$ , and it is desired to determine which distribution has the larger mean. Specifically,

- $\mathbf{X} = (X_1, \dots, X_m)$  random sample of  $m$  observations from a normal distribution for which both the mean  $\mu_1$  and the variance  $\sigma^2$  are unknown, and
- $\mathbf{Y} = (Y_1, \dots, Y_n)$  form an independent random sample of  $n$  observations from another normal distribution for which both the mean  $\mu_2$  and the variance  $\sigma^2$  are unknown.
- We shall assume that the variance  $\sigma^2$  is the same for both distributions, even though the value of  $\sigma^2$  is unknown.

If we are interested in testing hypotheses such as

$$H_0 : \mu_1 \leq \mu_2 \quad \text{versus} \quad H_1 : \mu_1 > \mu_2, \quad (14)$$

We reject  $H_0$  in (14) if the difference between the sample means is large. For all values of  $\theta = (\mu_1, \mu_2, \sigma^2)$  such that  $\mu_1 = \mu_2$ , the test statistics

$$U_{m,n} = \frac{\sqrt{m+n-2}(\bar{X}_m - \bar{Y}_n)}{\sqrt{(\frac{1}{m} + \frac{1}{n})(S_X^2 + S_Y^2)}}$$

follows the  $t$ -distribution with  $m+n-2$  degrees of freedom, where

$$S_X^2 = \sum_{i=1}^m (X_i - \bar{X}_m)^2, \quad \text{and} \quad S_Y^2 = \sum_{j=1}^n (Y_j - \bar{Y}_n)^2.$$

We reject  $H_0$  if

$$U_{m,n} \geq T_{m+n-2}^{-1}(1 - \alpha_0).$$

The  $p$ -value for the hypotheses in (14) is  $1 - T_{m+n-2}(u)$ , where  $u$  is the observed value of the statistic  $U_{m,n}$ .

If we are interested in testing hypotheses such as

$$H_0 : \mu_1 \geq \mu_2 \quad \text{versus} \quad H_1 : \mu_1 < \mu_2, \quad (15)$$

we reject  $H_0$  if

$$U_{m,n} \leq -T_{m+n-2}^{-1}(1 - \alpha_0) = T_{m+n-2}^{-1}(\alpha_0).$$



The  $p$ -value for the hypotheses in (15) is  $T_{m+n-2}(u)$ , where  $u$  is the observed value of the statistic  $U_{m,n}$ .

**Example 9.3** (Example 9.6.2, p.590). We test

$$H_0 : \mu_1 \leq \mu_2 \quad \text{vs.} \quad H_1 : \mu_1 > \mu_2.$$

Take  $\alpha_0 = 0.01$ ,  $n = m = 26$ ,  $\bar{X}_{26} = 5.13$ ,  $\bar{Y}_{26} = 3.99$ ,  $S_X^2 = 63.96$ ,  $S_Y^2 = 67.39$ . Then

$$U_{26,26} = \frac{(26 + 26 - 2)(5.13 - 3.99)}{\left(\frac{1}{26} + \frac{1}{26}\right)^{1/2} (63.96 + 67.39)^{1/2}} = 2.544$$

and  $T_{50}^{-1}(0.99) = 2.403$ , so we reject  $H_0$  at level  $\alpha_0 = 0.01$ .

### 9.9.2 Two-sided alternatives

If we are interested in testing hypotheses such as

$$H_0 : \mu_1 = \mu_2 \quad \text{versus} \quad H_1 : \mu_1 \neq \mu_2, \quad (16)$$

we reject  $H_0$  if

$$|U_{m,n}| \geq T_{m+n-2}^{-1}(1 - \frac{\alpha_0}{2}).$$

The  $p$ -value for the hypotheses in (16) is  $2[1 - T_{m+n-2}(|u|)]$ , where  $u$  is the observed value of the statistic  $U_{m,n}$ .

The power function of the two-sided two-sample  $t$  test is based on the non-central  $t$ -distribution in the same way as was the power function of the one-sample two-sided  $t$ -test. The test  $\delta$  that rejects  $H_0$  when  $|U_{m,n}| \geq c$  has power function

$$\pi(\mu_1, \mu_2, \sigma^2 | \delta) = T_{m+n-2}(-c|\psi) + 1 - T_{m+n-2}(c|\psi),$$

where  $T_{m+n-2}(\cdot | \psi)$  is the c.d.f of the non-central  $t$ -distribution with  $m + n - 2$  degrees of freedom and non-centrality parameter  $\psi$  given by

$$\psi = \frac{\mu_1 - \mu_2}{\sqrt{\sigma^2 \left(\frac{1}{m} + \frac{1}{n}\right)}}.$$

## 9.10 Comparing the variances of two normal distributions ( $F$ -test)

- $\mathbf{X} = (X_1, \dots, X_m)$  random sample of  $m$  observations from a normal distribution for which both the mean  $\mu_1$  and the variance  $\sigma_1^2$  are unknown, and

- $\mathbf{Y} = (Y_1, \dots, Y_n)$  form an independent random sample of  $n$  observations from another normal distribution for which both the mean  $\mu_2$  and the variance  $\sigma_2^2$  are unknown.

Suppose that we want to test the hypothesis of equality of the population variances, i.e.,  $H_0 : \sigma_1^2 = \sigma_2^2$ .

---

**Definition 22** (*F-distribution*). *Let  $Y$  and  $W$  be independent random variables such that  $Y \sim \chi_m^2$  and  $W \sim \chi_n^2$ . Then the distribution of*

$$X = \frac{Y/m}{W/n}$$

*is called the F-distribution with  $m$  and  $n$  degrees of freedom.*

---

The test statistic

$$V_{m,n}^* = \frac{\frac{S_X^2}{\sigma_1^2}/(m-1)}{\frac{S_Y^2}{\sigma_2^2}/(n-1)} = \frac{\sigma_2^2 S_X^2/(m-1)}{\sigma_1^2 S_Y^2/(n-1)}$$

follows the *F*-distribution with  $m-1$  and  $n-1$  degrees of freedom. In particular, if  $\sigma_1^2 = \sigma_2^2$ , then the distribution of

$$V_{m,n} = \frac{S_X^2/(m-1)}{S_Y^2/(n-1)}$$

is the *F*-distribution with  $m-1$  and  $n-1$  degrees of freedom.

Let  $\nu$  be the observed value of the statistic  $V_{m,n}$  below, and let  $G_{m-1,n-1}(\cdot)$  be the c.d.f of the *F*-distribution with  $m-1$  and  $n-1$  degrees of freedom.

### 9.10.1 One-sided alternatives

If we are interested in testing hypotheses such as

$$H_0 : \sigma_1^2 \leq \sigma_2^2 \quad \text{versus} \quad H_1 : \sigma_1^2 > \sigma_2^2, \quad (17)$$

we reject  $H_0$  if

$$V_{m,n} \geq G_{m-1,n-1}^{-1}(1 - \alpha_0).$$

The *p*-value for the hypotheses in (17) when  $V_{m,n} = \nu$  is observed equals  $1 - G_{m-1,n-1}(\nu)$ .

### 9.10.2 Two-sided alternatives

If we are interested in testing hypotheses such as

$$H_0 : \sigma_1^2 = \sigma_2^2, \quad \text{versus} \quad H_1 : \sigma_1^2 \neq \sigma_2^2, \quad (18)$$

we reject  $H_0$  if either  $V_{m,n} \leq c_1$  or  $V_{m,n} \geq c_2$ , where  $c_1$  and  $c_2$  are constants such that

$$\mathbb{P}(V_{m,n} \leq c_1) + \mathbb{P}(V_{m,n} \geq c_2) = \alpha_0$$

when  $\sigma_1^2 = \sigma_2^2$ . The most convenient choice of  $c_1$  and  $c_2$  is the one that makes

$$\mathbb{P}(V_{m,n} \leq c_1) = \mathbb{P}(V_{m,n} \geq c_2) = \frac{\alpha_0}{2},$$

that is,

$$c_1 = G_{m-1, n-1}^{-1}(\alpha_0/2) \quad \text{and} \quad c_2 = G_{m-1, n-1}^{-1}(1 - \alpha_0/2).$$

**Example 9.4** (Example 9.7.4, p.601 and Example 9.6.2). We test

$$H_0 : \sigma_1^2 = \sigma_2^2, \quad \text{versus} \quad H_1 : \sigma_1^2 \neq \sigma_2^2, \quad (19)$$

Here  $m = n = 26$  and thus

$$V_{26,26} = \frac{63.96}{67.39} = 0.9491.$$

We compare this to

$$F_{25,25}^{-1}(0.025) = 0.4484 \quad \text{and} \quad F_{25,25}^{-1}(0.975) = 2.2303.$$

Thus we fail to reject  $H_0$ .

## 9.11 Likelihood ratio test

A very popular form of hypothesis test is the **likelihood ratio test**.

Suppose that we want to test

$$H_0 : \theta \in \Omega_0, \quad \text{and} \quad H_1 : \theta \in \Omega_1. \quad (20)$$

In order to compare these two hypotheses, we might wish to see whether the likelihood function is higher on  $\Omega_0$  or on  $\Omega_1$ .

The **likelihood ratio statistic** (LRS) is defined as

$$\Lambda(\mathbf{X}) = \frac{\sup_{\theta \in \Omega_0} L_n(\theta, \mathbf{X})}{\sup_{\theta \in \Omega} L_n(\theta, \mathbf{X})}, \quad (21)$$

where  $\Omega = \Omega_0 \cup \Omega_1$ .

A likelihood ratio test of the hypotheses (20) rejects  $H_0$  when

$$\Lambda(\mathbf{x}) \leq k,$$

for some constant  $k$ .

Interpretation: we reject  $H_0$  if the likelihood function on  $\Omega_0$  is sufficiently small compared to the likelihood function on all of  $\Omega$ .

Generally,  $k$  is to be chosen so that the test has a desired level  $\alpha_0$ .

**Example 9.5.** Suppose that  $\mathbf{X}_n = (X_1, \dots, X_n)$  is a random sample from a normal distribution with unknown mean  $\mu$  and known variance  $\sigma^2$ . We wish to test the hypotheses

$$H_0 : \mu = \mu_0 \quad \text{versus} \quad H_1 : \mu \neq \mu_0$$

at the level  $\alpha_0$ . We note that the MLE for  $\mu$  is given by  $\bar{x}$  and we compute

$$L_n(\mu_0, \mathbf{X}) = \left( \frac{1}{2\pi\sigma^2} \right)^{n/2} \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu_0)^2 \right),$$

$$L_n(\bar{x}, \mathbf{X}) = \left( \frac{1}{2\pi\sigma^2} \right)^{n/2} \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2 \right).$$

Thus

$$\Theta(x) = \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu_0)^2 - (x_i - \bar{x})^2 \right) = \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (\mu_0 - \bar{x})^2 \right)$$

and the LRT is equivalent to the  $z$ -test.

**Theorem 9.6.** Let  $\Omega$  be a open set of a  $k$ -dimensional space, and suppose that  $H_0$  specifies that  $p$  coordinates ( $p \leq k$ ) of  $\theta$  are equal to  $p$  specific values. Assume that  $H_0$  is true and that the likelihood function satisfies the conditions needed to prove that the MLE is asymptotically normal and asymptotically efficient. Then, as  $n \rightarrow \infty$ ,

$$-2 \log \Lambda(\mathbf{X}) \xrightarrow{d} \chi_p^2.$$

## 9.12 Equivalence of hypothesis tests and confidence sets

**Example:** Suppose that  $X_1, \dots, X_n$  are i.i.d  $N(\mu, \sigma^2)$  where  $\mu$  is unknown and  $\sigma^2$  is known.

We now illustrate how the testing procedure ties up naturally with the CI construction problem.

Consider testing  $H_0 : \mu = \mu_0$  versus  $H_1 : \mu \neq \mu_0$ .

First note that the acceptance region of the derived test  $\delta$  can be written as:

$$S_0 = \mathcal{A}_{\mu_0} = \left\{ \mathbf{x} = (x_1, x_2, \dots, x_n) : \bar{x} - \frac{\sigma}{\sqrt{n}} z_{\alpha_0/2} \leq \mu_0 \leq \bar{x} + \frac{\sigma}{\sqrt{n}} z_{\alpha_0/2} \right\}.$$

Now, consider a fixed data set  $(X_1, X_2, \dots, X_n)$  and based on this consider testing a family of null hypotheses:

$$\{H_{0,\tilde{\mu}} : \mu = \tilde{\mu} : \tilde{\mu} \in \mathbb{R}\}.$$

We can now ask the following question: Based on the observed data and the above testing procedure, *what values of  $\tilde{\mu}$  would fail to be rejected by the level  $\alpha_0$  test?* This means that  $\tilde{\mu}$  would have to fall in the interval

$$\bar{X} - \frac{\sigma}{\sqrt{n}} z_{\alpha_0/2} \leq \tilde{\mu} \leq \bar{X} + \frac{\sigma}{\sqrt{n}} z_{\alpha_0/2}.$$

Thus, the set of  $\tilde{\mu}$ 's for which the null hypothesis would fail to be rejected by the level  $\alpha_0$  test is the set:

$$\left[ \bar{X} - \frac{\sigma}{\sqrt{n}} z_{\alpha_0/2}, \bar{X} + \frac{\sigma}{\sqrt{n}} z_{\alpha_0/2} \right].$$

But this is precisely the level  $1 - \alpha_0$  CI that we obtained before!

Thus, we obtain a level  $1 - \alpha_0$  CI for  $\mu$ , the population mean, by *compiling all possible  $\tilde{\mu}$ 's for which the null hypothesis  $H_{0,\tilde{\mu}} : \mu = \tilde{\mu}$  fails to be rejected by the level  $\alpha_0$  test.*

**From hypothesis testing to CIs:** Let  $X_1, X_2, \dots, X_n$  be i.i.d observations from some underlying distribution  $F_\theta$ ; here  $\theta$  is a “parameter” indexing a family of distributions. The goal is to construct a CI for  $\theta$  using hypothesis testing.

For each  $\tilde{\theta}$  consider testing the null hypothesis  $H_{0,\tilde{\theta}} : \theta = \tilde{\theta}$ . Suppose, there exists a level  $\alpha_0$  test  $\delta_{\tilde{\theta}}(\mathbf{X})$  for this problem with

$$\mathcal{A}_{\tilde{\theta}} = \{\mathbf{x} : T_{\tilde{\theta}}(\mathbf{x}) \leq c_{\alpha_0}\}$$

being the acceptance region of  $\delta_{\tilde{\theta}}$  and  $\mathbb{P}_{\tilde{\theta}}(\mathbf{X} \in \mathcal{A}_{\tilde{\theta}}) \geq 1 - \alpha_0$ . Then a level  $1 - \alpha$  confidence set for  $\theta$  is:

$$\mathcal{S}(\mathbf{X}) = \{\tilde{\theta} : \mathbf{X} \in \mathcal{A}_{\tilde{\theta}}\}.$$

We need to verify that for any  $\theta$ ,  $\mathbb{P}_\theta[\theta \in \mathcal{S}(\mathbf{X})] \geq 1 - \alpha$ . But

$$\mathbb{P}_\theta(\theta \in \mathcal{S}(\mathbf{X})) = \mathbb{P}_\theta(\mathbf{X} \in \mathcal{A}_\theta) \geq 1 - \alpha_0.$$

**Theorem 9.7.** For each  $\theta_0 \in \Omega$ , let  $\mathcal{A}(\theta_0)$  be the acceptance region of a level  $\alpha$  test of  $H_0 : \theta = \theta_0$ . For each  $\mathbf{x} \in \mathcal{X}$  ( $\mathcal{X}$  is the space of all data values), define a set  $\mathcal{S}(\mathbf{x})$  in the parameter space by

$$\mathcal{S}(\mathbf{x}) = \{\theta_0 : \mathbf{x} \in \mathcal{A}(\theta_0)\}.$$

Then the random set  $\mathcal{S}(\mathbf{X})$  is a  $1 - \alpha$  confidence set. Conversely, let  $\mathcal{S}(\mathbf{X})$  be a  $1 - \alpha$  confidence set. For any  $\theta_0 \in \Omega$ , define

$$\mathcal{A}(\theta_0) = \{\mathbf{x} : \theta_0 \in \mathcal{S}(\mathbf{x})\}.$$

Then  $\mathcal{A}(\theta_0)$  is the acceptance region of a level  $\alpha$  test of  $H_0 : \theta = \theta_0$ .

*Proof.* The first part is essentially done above!

For the second part, the type I error probability for the test of  $H_0 : \theta = \theta_0$  with acceptance region  $\mathcal{A}(\theta_0)$  is  $\mathbb{P}_{\theta_0}(\mathbf{X} \notin \mathcal{A}_{\theta_0}) = \mathbb{P}_{\theta_0}[\theta_0 \notin \mathcal{S}(\mathbf{X})] \leq \alpha$ .  $\square$

**Remark:** The more useful part of the theorem is the first part, i.e., given a level  $\alpha$  test (which is usually easy to construct) we can get a confidence set by inverting the family of tests.

**Example 9.8.** Suppose that  $X_1, \dots, X_n$  are i.i.d  $\text{Exp}(\lambda)$ . We want to test  $H_0 : \lambda = \lambda_0$  versus  $H_1 : \lambda \neq \lambda_0$ . Find the LRT.

The acceptance region is given by

$$\mathcal{A}(\lambda_0) = \left\{ \mathbf{x} : \left( \lambda_0 \sum x_i \right)^n e^{-\lambda_0 \sum x_i} \geq c^* \right\},$$

where  $c^*$  is a constant chosen to satisfy

$$\mathbb{P}_{\lambda_0}(\mathbf{X} \in \mathcal{A}(\lambda_0)) = 1 - \alpha.$$

Inverting this acceptance region gives the  $1 - \alpha$  confidence set

$$\mathcal{S}(\mathbf{x}) = \left\{ \lambda : \left( \lambda \sum x_i \right)^n e^{-\lambda \sum x_i} \geq c^* \right\}.$$

This can be shown to be an interval in the parameter space.

**Example 9.9.** Suppose that  $X_1, \dots, X_n$  are i.i.d with p.d.f/p.m.f  $f(\cdot|\theta)$  where the unknown parameter is  $\theta \equiv (\theta_1, \dots, \theta_k) \in \Omega \subset \mathbb{R}^k$ . We want to construct an (approximate)  $(1 - \alpha)$ -CI for  $\theta_1 \subset \mathbb{R}$ . Note that, in this generality, it is not obvious how to construct an approximate CI for  $\theta_1$ . Here we explain how the *duality* between hypothesis testing and CIs can be used to solve this problem.

In that setting, consider testing the statistical hypothesis

$$H_0 : \theta_1 = \theta_1^0 \quad \text{versus} \quad H_1 : \theta_1 \neq \theta_1^0,$$

where  $\theta_1^0 \in \mathbb{R}$  is a fixed number (e.g.,  $\theta_1^0 = 0$ ). Note that this is (possibly) a composite null hypothesis as the values of  $\theta_2, \dots, \theta_k$  are left unspecified. We can use Theorem 9.6 to find an approximate  $1 - \alpha$  acceptance region for the above test, i.e.,

$$\mathcal{A}(\theta_1^0) := \left\{ \mathbf{x} : -2 \log(\Lambda_{\theta_1^0}(\mathbf{x})) \leq z_{\alpha/2}^2 \right\} \quad (22)$$

where  $\Lambda_{\theta_1^0}(\mathbf{X})$  is the LRS for testing  $H_0 : \theta_1 = \theta_1^0$ . Here the critical value of the above test is  $z_{\alpha/2}^2$  as  $\chi_1^2 \equiv Z^2$  where  $Z \sim N(0, 1)$ .

Inverting this acceptance region (for every value of  $\theta_1^0$ ) gives an approximate  $1 - \alpha$  confidence set for  $\theta_1$ :

$$\mathcal{S}(\mathbf{X}) := \{\theta_1^0 \in \mathbb{R} : \mathbf{X} \in \mathcal{A}(\theta_1^0)\} = \left\{ \theta_1^0 \in \mathbb{R} : -2 \log(\Lambda_{\theta_1^0}(\mathbf{X})) \leq z_{\alpha/2}^2 \right\},$$

where  $\mathbf{X} = (X_1, \dots, X_n)$  is the observed data.

This idea can be easily implemented in a computer: Consider a fine grid of value  $\theta_1^1 < \theta_1^2 < \dots < \theta_1^M$  on the real line (for some  $M \geq 1$ ). We can use (22) to test the hypothesis  $H_0 : \theta_1 = \theta_1^j$  for every  $j = 1, \dots, M$ . Then an approximation of  $\mathcal{S}(\mathbf{X})$  can be the smallest interval containing the set

$$\{\theta_1^j \in \mathbb{R} : \mathbf{X} \in \mathcal{A}(\theta_1^j), j = 1, \dots, M\},$$

i.e., we just consider the values  $\theta_1^j$  for which the corresponding null hypothesis is accepted.

## 10 Linear regression

- We are often interested in understanding the *relationship* between two or more variables.
- Want to model a functional relationship between a “predictor” (input, independent variable) and a “response” variable (output, dependent variable, etc.).
- But the real world is noisy, no  $f = ma$  (Force = mass  $\times$  acceleration). We have observation noise, weak relationship, etc.

### Examples:

- How is the *sales price* of a house related to its size, number of rooms and property tax?
- How does the probability of *surviving* a particular surgery change as a function of the patient’s age and general health condition?
- How does the *weight* of an individual depend on his/her height?

### Notation:

Suppose that we have  $n$  data points  $(x_1, Y_1), \dots, (x_n, Y_n)$ . We want to predict  $Y$  given a value of  $x$ .

- $Y_i$  is the value of the **response** variable for the  $i$ -th observation.
  - $x_i$  is the value of the **predictor** (covariate/explanatory variable) for the  $i$ -th observation.
  - **Scatter plot:** Plot the data and try to visualize the relationship.
- 

### 10.1 Simple linear regression

- Suppose that we think that  $Y$  is a **linear** function (actually here a more appropriate term is “affine”) of  $x$ , i.e.,

$$Y_i \approx \beta_0 + \beta_1 x_i,$$

and we want to find the “best” such linear function.

- The model for **simple linear regression** can be stated as follows:

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, \dots, n.$$



- Observations:  $\{(x_i, Y_i) : i = 1, \dots, n\}$ .
- $\beta_0$ ,  $\beta_1$  and  $\sigma^2$  are *unknown* parameters.
- $\epsilon_i$  is a (unobserved) **random error** term whose distribution is unspecified:

$$\mathbb{E}(\epsilon_i) = 0, \quad \text{Var}(\epsilon_i) = \sigma^2, \quad \text{Cov}(\epsilon_i, \epsilon_j) = 0 \quad \text{for } i \neq j.$$

- $x_i$ 's will be treated as known *constants*. Even if the  $x_i$ 's are random, we condition on the predictors and want to understand the **conditional distribution** of  $Y$  given  $X$ .
- **Regression function: Conditional mean** on  $Y$  given  $x$ , i.e.,

$$m(x) := \mathbb{E}[Y] = \beta_0 + \beta_1 x.$$

To emphasize the dependence on  $x$  we sometimes also write  $m(x) = \mathbb{E}[Y \mid x]$ .

- The regression function shows how the mean of  $Y$  changes as a *function* of  $x$ .
- $\mathbb{E}[Y_i] = \mathbb{E}(\beta_0 + \beta_1 x_i + \epsilon_i) = \beta_0 + \beta_1 x_i$
- $\text{Var}(Y_i) = \text{Var}(\beta_0 + \beta_1 x_i + \epsilon_i) = \text{Var}(\epsilon_i) = \sigma^2$ .

### 10.1.1 Interpretation

- The slope  $\beta_1$  has units “y-units per x-units”.
  - For every 1 inch increase in height, the model predicts a  $\beta_1$  *pounds increase* in the mean weight.
- The intercept term  $\beta_0$  is not always meaningful (depending on the application).
- The model is *only valid* for values of the explanatory variable in the domain of the data.

## 10.2 Method of least squares

- After formulating the model we use the observed data to *estimate* the *unknown* parameters.
- Three unknown parameters:  $\beta_0, \beta_1$  and  $\sigma^2$ .
- We are interested in finding the estimates of these parameters that *best fit* the data.

- Intuition: For the correct parameter values  $\beta_0$  and  $\beta_1$ , the *deviation* of the observed values to its expected value, i.e.,

$$Y_i - \beta_0 - \beta_1 x_i,$$

should be *small*.

- We try to *minimize* the sum of the  $n$  squared deviations, i.e., we can try to minimize

$$Q(b_0, b_1) = \sum_{i=1}^n (Y_i - b_0 - b_1 x_i)^2$$

as a function of  $b_0$  and  $b_1$ . In other words, we want to minimize the sum of the squares of the vertical deviations of all the points from the line. Minimisers of  $Q$  are called the **least squares** estimators of  $\beta_0$  and  $\beta_1$ .

- The least squares estimators can be found by differentiating  $Q$  with respect to  $b_0$  and  $b_1$  and setting the partial derivatives equal to 0.
- Find  $b_0$  and  $b_1$  that solve:

$$\begin{aligned} \frac{\partial Q}{\partial b_0} &= -2 \sum_{i=1}^n (Y_i - b_0 - b_1 x_i) = 0 \\ \frac{\partial Q}{\partial b_1} &= -2 \sum_{i=1}^n x_i (Y_i - b_0 - b_1 x_i) = 0. \end{aligned}$$

### 10.2.1 Normal equations

- The values of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  that minimize  $Q$  are given by the solution to the *normal equations*:

$$\sum_{i=1}^n Y_i = n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i \quad (23)$$

$$\sum_{i=1}^n x_i Y_i = \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2. \quad (24)$$

- Solving the normal equations gives us the following point estimates:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}) Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad (25)$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}, \quad (26)$$

where  $\bar{x} = \sum_{i=1}^n x_i / n$  and  $\bar{Y} = \sum_{i=1}^n Y_i / n$ . For this we have to subtract  $\bar{x}$  times (23) from (24) and use that  $\bar{x} \sum x_i = n\bar{x}^2$  and  $\sum x_i^2 - n\bar{x}^2 = \sum (x_i - \bar{x})^2$ .

In general, if we can parametrize the form of the functional dependence between  $Y$  and  $x$  in a linear fashion (linear in the parameters), then the method of least squares can be used to estimate the function. For example,

$$Y_i \approx \beta_0 + \beta_1 x_i + \beta_2 x_i^2$$

is still linear in the parameters.

### 10.2.2 Estimated regression function

- We estimate the regression function:

$$\mathbb{E}[Y] = \beta_0 + \beta_1 x$$

using

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x.$$

- The term

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, \quad i = 1, \dots, n,$$

is called the **fitted** or *predicted* value for the  $i$ -th observation, while  $Y_i$  is the observed value.

- The *residual*, denoted  $e_i$ , is the difference between the observed and the predicted value of  $Y_i$ , i.e.,

$$e_i = Y_i - \hat{Y}_i.$$

- The residuals show how far the individual data points fall from the regression function.

### 10.2.3 Properties

- As  $\hat{\beta}_0, \hat{\beta}_1$  satisfy (23) and (24) we conclude

$$\begin{aligned} \sum_{i=1}^n e_i &= \sum_{i=1}^n (Y_i - \hat{Y}_i) = \sum_{i=1}^n (Y_i - \hat{\beta}_0 + \hat{\beta}_1 x_i) = 0, \\ \sum_{i=1}^n x_i e_i &= \sum_{i=1}^n x_i (Y_i - \hat{Y}_i) = \sum_{i=1}^n x_i (Y_i - \hat{\beta}_0 + \hat{\beta}_1 x_i) = 0. \end{aligned}$$

- In particular, the regression line always passes through the point  $(\bar{x}, \bar{Y})$ .

### 10.2.4 Estimation of $\sigma^2$

- Recall:  $\sigma^2 = \text{Var}(\epsilon_i)$ .
- Naive idea: use  $\hat{\sigma}^2 = \frac{\sum_{i=1}^n (\epsilon_i - \bar{\epsilon})^2}{n-1}$ . But  $\epsilon_i$ 's are not *observed*!
- Better idea: use  $e_i$ 's, i.e.,  $s^2 := \frac{\sum_{i=1}^n (e_i - \bar{e})^2}{n-2} = \frac{\sum_{i=1}^n e_i^2}{n-2}$ .
- The divisor  $n - 2$  in  $s^2$  is the number of **degrees of freedom** associated with the estimate.
- To obtain  $s^2$ , the two parameters  $\beta_0$  and  $\beta_1$  must first be estimated, which results in a loss of *two* degrees of freedom.
- Using  $n - 2$  makes  $s^2$  an *unbiased* estimator of  $\sigma^2$ , i.e.,  $\mathbb{E}(s^2) = \sigma^2$ .

### 10.2.5 Gauss-Markov theorem

The least squares estimators  $\hat{\beta}_0, \hat{\beta}_1$  are **unbiased**, i.e.,

$$\begin{aligned}\mathbb{E}[\hat{\beta}_1] &= \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \mathbb{E}[Y_i] \\ &= \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) (\beta_0 + \beta_1 x_i) \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \beta_1 = \beta_1,\end{aligned}$$

where we recall that  $\sum_{i=1}^n x_i (x_i - \bar{x}) = \sum_{i=1}^n (x_i - \bar{x})^2$ , and

$$\mathbb{E}[\hat{\beta}_0] = \mathbb{E}[\bar{Y}] - \mathbb{E}[\hat{\beta}_1] \bar{x} = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0.$$

**Definition 23.** A linear estimator of  $\beta_j$  ( $j = 0, 1$ ) is an estimator of the form

$$\tilde{\beta}_j = \sum_{i=1}^n c_i Y_i,$$

where the coefficients  $c_1, \dots, c_n$  are only allowed to depend on  $x_i$ .

Note that  $\hat{\beta}_0, \hat{\beta}_1$  are linear estimators (this will be an exercise).

**Result:** No matter what the distribution of the error terms  $\epsilon_i$ , the least squares method provides *unbiased* point estimates that have **minimum** variance among all **unbiased linear estimators**.

The Gauss-Markov theorem states that in a linear regression model in which the errors have **expectation zero** and are **uncorrelated** and have **equal variances**, the *best linear unbiased estimator* (BLUE) of the coefficients is given by the **ordinary least squares estimators**.

## 10.3 Normal simple linear regression

To perform *inference* we need to make assumptions regarding the distribution of  $\epsilon_i$ .

We often assume that  $\epsilon_i$ 's are *normally* distributed.

The *normal error* version of the model for simple linear regression can be written:

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, \dots, n.$$

Here  $\epsilon_i$ 's are independent  $N(0, \sigma^2)$ ,  $\sigma^2$  unknown.

Hence,  $Y_i$ 's are independent normal random variables with mean  $\beta_0 + \beta_1 x_i$  and variance  $\sigma^2$ .

### 10.3.1 Maximum likelihood estimation

When the probability distribution of  $Y_i$  is *specified*, the estimates can be obtained using the method of *maximum likelihood*.

This method chooses as estimates those values of the parameter that are most *consistent* with the observed data.

The *likelihood* is the *joint density* of the  $Y_i$ 's viewed as a function of the unknown parameters, which we denote  $L(\beta_0, \beta_1, \sigma^2)$ .

Since the  $Y_i$ 's are *independent* this is simply the *product* of the density of individual  $Y_i$ 's.

We seek the values of  $\beta_0, \beta_1$  and  $\sigma^2$  that maximize  $L(\beta_0, \beta_1, \sigma^2)$  for the given  $x$  and  $Y$  values in the sample.

According to our model:

$$Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2), \quad \text{for } i = 1, 2, \dots, n.$$

The likelihood function for the  $n$  independent observations  $Y_1, \dots, Y_n$  is given by

$$\begin{aligned} L(\beta_0, \beta_1, \sigma^2) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (Y_i - \beta_0 - \beta_1 x_i)^2 \right\} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2 \right\}. \end{aligned} \tag{27}$$

The value of  $(\beta_0, \beta_1, \sigma^2)$  that maximizes the likelihood function are called *maximum likelihood estimates* (MLEs).

The MLE of  $\beta_0$  and  $\beta_1$  are *identical* to the ones obtained using the method of *least squares*, i.e.,

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}) Y_i}{S_x^2},$$

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

The MLE of  $\sigma^2$ :  $\hat{\sigma}^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n}$ .

### 10.3.2 Inference

Our model describes the *linear* relationship between the two variables  $x$  and  $Y$ .

Different samples from the same population will produce different point estimates of  $\beta_0$  and  $\beta_1$ .

Hence,  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are random variables with sampling distributions that describe *what values* they can take and *how often* they take them.

Hypothesis tests about  $\beta_0$  and  $\beta_1$  can be constructed using these distributions.

The next step is to perform *inference*, including:

- Tests and confidence intervals for the *slope* and *intercept*.
- Confidence intervals for the *mean response*.
- *Prediction intervals* for new observations.

**Theorem 10.1.** *Under the assumptions of the normal linear model,*

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} \sim N_2 \left( \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}, \sigma^2 \begin{pmatrix} \frac{1}{n} + \frac{\bar{x}^2}{S_x^2} & -\frac{\bar{x}}{S_x^2} \\ -\frac{\bar{x}}{S_x^2} & \frac{1}{S_x^2} \end{pmatrix} \right)$$

where  $S_x^2 = \sum_{i=1}^n (x_i - \bar{x})^2$ . Also, if  $n \geq 3$ ,  $\hat{\sigma}^2$  is independent of  $(\hat{\beta}_0, \hat{\beta}_1)$  and  $n\hat{\sigma}^2/\sigma^2$  has a  $\chi^2$ -distribution with  $n - 2$  degrees of freedom.

Note that if the  $x_i$ 's are random, the above theorem is still valid if we condition on the values of the predictor  $x_i$ 's.

### 10.3.3 Inference about $\beta_1$

We often want to perform tests about the *slope*:

$$H_0 : \beta_1 = 0 \quad \text{versus} \quad H_1 : \beta_1 \neq 0.$$

Under the null hypothesis there is *no linear relationship* between  $Y$  and  $x$  – the *means* of probability distributions of  $Y$  are equal at all levels of  $x$ , i.e.,  $\mathbb{E}(Y|x) = \beta_0$ , for all  $x$ .

**Lemma 10.2.** *The sampling distribution of  $\hat{\beta}_1$  under  $H_0$  is given by*

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{S_x^2}\right).$$

*Proof.* Need to show that:  $\hat{\beta}_1$  is normally distributed,

$$\mathbb{E}(\hat{\beta}_1) = \beta_1, \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{S_x^2}.$$

**Result:** When  $Z_1, \dots, Z_k$  are *independent* normal random variables, the linear combination

$$a_1 Z_1 + \dots + a_k Z_k$$

is also *normally* distributed.

Since  $\hat{\beta}_1$  is a linear combination of the  $Y_i$ 's and each  $Y_i$  is an *independent normally* distributed random variable, then  $\hat{\beta}_1$  is also normally distributed.

We already know that  $\mathbb{E}[\hat{\beta}_1] = \beta_1$ . As  $Y_i$  are independent we also get

$$\text{Var}(\beta_1) = \sigma^2 \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{S_x^2} \right)^2 = \frac{\sigma^2}{S_x^2}.$$

This shows the claim. □

- **Variance for the estimated slope:** There are *three* aspects of the scatter plot that affect the variance of the regression slope:
  - The *spread* around the *regression line* ( $\sigma^2$ ) – less scatter around the line means the estimated slope will be more consistent from sample to sample.
  - The *spread* of the *x values*  $S_x^2/n$  – a large variance of  $x$  provides a more stable regression.
  - The *sample size*  $n$  – having a larger sample size  $n$ , gives more consistent estimates.

- **Estimated variance:** When  $\sigma^2$  is *unknown* we replace it with the

$$\tilde{\sigma}^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n - 2} = \frac{\sum_{i=1}^n e_i^2}{n - 2}.$$

Plugging this into the equation for  $\text{Var}(\hat{\beta}_1)$  we get

$$se^2(\hat{\beta}_1) = \frac{\tilde{\sigma}^2}{S_x^2}.$$

Recall: *Standard error*  $se(\hat{\theta})$  of an estimator  $\hat{\theta}$  is used to refer to an *estimate* of its *standard deviation*.

**Result:** For the normal error regression model:

$$\frac{RSS}{\sigma^2} = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sigma^2} \sim \chi_{n-2}^2,$$

and is *independent* of  $\hat{\beta}_0$  and  $\hat{\beta}_1$ .

- **(Studentized statistic:)** Since  $\hat{\beta}_1$  is *normally* distributed, the standardized statistic:

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\text{Var}(\hat{\beta}_1)}} \sim N(0, 1).$$

If we replace  $\text{Var}(\hat{\beta}_1)$  by its estimate we get the *studentized* statistic:

$$\frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} \sim t_{n-2}.$$

Recall: Suppose that  $Z \sim N(0, 1)$  and  $W \sim \chi_p^2$  where  $Z$  and  $W$  are independent. Then,

$$\frac{Z}{\sqrt{W/p}} \sim t_p,$$

the *t-distribution* with  $p$  *degrees of freedom*.

We derive this result as follows: Note that

$$\frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\text{Var}(\hat{\beta}_1)}} \cdot \sqrt{\frac{\text{Var}(\hat{\beta}_1)}{se^2(\hat{\beta}_1)}},$$

where  $\text{Var}(\hat{\beta}_1) = \sigma^2/S_x^2$ ,  $se^2(\hat{\beta}_1) = \tilde{\sigma}^2/S_x^2$  and  $\tilde{\sigma}^2(n-2)/\sigma^2$  has a  $\chi^2$ -distribution with  $n-2$  degrees of freedom. Thus

$$\begin{aligned} \frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} &= \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\text{Var}(\hat{\beta}_1)}} \cdot \sqrt{\frac{\frac{\sigma^2}{S_x^2}}{\frac{\tilde{\sigma}^2}{S_x^2}}} \\ &= \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\text{Var}(\hat{\beta}_1)}} \cdot \sqrt{\frac{\sigma^2}{\tilde{\sigma}^2}} \\ &= \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\text{Var}(\hat{\beta}_1)}} \cdot \sqrt{\frac{\sigma^2(n-2)}{\tilde{\sigma}^2(n-2)}}. \end{aligned}$$

This shows the claim.



- **Hypothesis testing:** To test

$$H_0 : \beta_1 = 0 \quad \text{versus} \quad H_1 : \beta_1 \neq 0$$

we use the *test-statistic*

$$T = \frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)}.$$

We reject  $H_0$  when the observed value of  $|T|$  i.e.,  $|t_{obs}|$ , is *large*!

Thus, given *level*  $(1 - \alpha)$ , we reject  $H_0$  if

$$|t_{obs}| > t_{n-2}(1 - \alpha/2)$$

where  $t_{n-2}(1 - \alpha/2)$  denotes the  $(1 - \alpha/2)$ -quantile of the  $t_{n-2}$ -distribution, i.e.,

$$1 - \frac{\alpha}{2} = \mathbb{P}(T \leq t_{n-2}(1 - \alpha/2)).$$

- **P-value:**  $p$ -value is the probability of obtaining a test statistic at least as extreme as the one that was actually observed, assuming that the null hypothesis is true.

The  $p$ -value depends on  $H_1$  (one-sided/two-sided).

In our case, we compute  $p$ -values using a  $t_{n-2}$ -distribution. Thus,

$$p\text{-value} = \mathbb{P}_{H_0}(|T| > |t_{obs}|).$$

If we know the  $p$ -value then we can decide to accept/reject  $H_0$  (versus  $H_1$ ) at any given  $\alpha$ .

- **Confidence interval:** A *confidence interval* (CI) is a kind of *interval estimator* of a population parameter and is used to indicate the reliability of an estimator. Using the sampling distribution of  $\hat{\beta}_1$  we can make the following probability statement:

$$\begin{aligned} \mathbb{P} \left( t_{n-2}(\alpha/2) \leq \frac{\hat{\beta}_1 - \beta_1}{\text{se}(\hat{\beta}_1)} \leq t_{n-2}(1 - \alpha/2) \right) &= 1 - \alpha \\ \mathbb{P} \left( \hat{\beta}_1 - t_{n-2}(1 - \alpha/2)\text{se}(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + t_{n-2}(\alpha/2)\text{se}(\hat{\beta}_1) \right) &= 1 - \alpha. \end{aligned}$$

Thus, a  $(1 - \alpha)$  confidence interval for  $\beta_1$  is

$$\left[ \hat{\beta}_1 - t_{n-2}(1 - \alpha/2) \cdot \text{se}(\hat{\beta}_1), \hat{\beta}_1 + t_{n-2}(1 - \alpha/2) \cdot \text{se}(\hat{\beta}_1) \right]$$

as  $t_{n-2}(1 - \alpha/2) = -t_{n-2}(\alpha/2)$ .

### 10.3.4 Sampling distribution of $\hat{\beta}_0$

The *sampling distribution* of  $\hat{\beta}_0$  is

$$N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_x^2}\right)\right).$$

This can be verified in the same way as we did for  $\hat{\beta}_1$ .

**Hypothesis testing:** In general, let  $c_0, c_1$  and  $c_*$  be specified numbers, where at least one of  $c_0$  and  $c_1$  is nonzero. Suppose that we are interested in testing the following hypotheses:

$$H_0 : c_0\beta_0 + c_1\beta_1 = c_*, \quad \text{versus} \quad H_1 : c_0\beta_0 + c_1\beta_1 \neq c_*. \quad (28)$$

We should use a scalar multiple of

$$c_0\hat{\beta}_0 + c_1\hat{\beta}_1 - c_*$$

as the test statistic. Specifically, we use

$$U_{01} = \left[ \frac{c_0^2}{n} + \frac{(c_0\bar{x} - c_1)^2}{S_x^2} \right]^{-1/2} \left( \frac{c_0\hat{\beta}_0 + c_1\hat{\beta}_1 - c_*}{\tilde{\sigma}} \right),$$

where

$$\tilde{\sigma}^2 = \frac{S^2}{n-2}, \quad S^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \sum_{i=1}^n e_i^2.$$

Note that  $\tilde{\sigma}^2$  is an unbiased estimator of  $\sigma^2$ .

For each  $\alpha \in (0, 1)$ , a level  $\alpha$  test of the hypothesis (28) is to reject  $H_0$  if

$$|U_{01}| > T_{n-2}^{-1} \left( 1 - \frac{\alpha}{2} \right).$$

The above result follows from the fact that  $c_0\hat{\beta}_0 + c_1\hat{\beta}_1 - c_*$  is normally distributed with mean  $c_0\beta_0 + c_1\beta_1 - c_*$  and variance

$$\begin{aligned} \text{Var}(c_0\hat{\beta}_0 + c_1\hat{\beta}_1 - c_*) &= c_0^2 \text{Var}(\hat{\beta}_0) + c_1^2 \text{Var}(\hat{\beta}_1) + 2c_0c_1 \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ &= c_0^2 \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_x^2} \right) + c_1^2 \sigma^2 \frac{1}{S_x^2} - 2c_0c_1 \frac{\sigma^2 \bar{x}}{S_x^2} \\ &= \sigma^2 \left[ \frac{c_0^2}{n} + \frac{c_0^2 \bar{x}^2}{S_x^2} - 2c_0c_1 \frac{\bar{x}}{S_x^2} + c_1^2 \frac{1}{S_x^2} \right] \\ &= \sigma^2 \left[ \frac{c_0^2}{n} + \frac{(c_0\bar{x} - c_1)^2}{S_x^2} \right]. \end{aligned}$$

Using the above result, we can actually test if the regression line  $y = \beta_0 + \beta_1 x$  passes through a particular point  $(x^*, y^*)$ , where  $x^* \neq 0$ . Indeed the hypotheses we are interested in now take the form:

$$\begin{aligned} H_0 : & \beta_0 + x^* \beta_1 = y^* \\ H_1 : & \beta_0 + x^* \beta_1 \neq y^*. \end{aligned}$$

This corresponds to the above test for  $c_0 = 1$  and  $c_1 = x^*$ .

**Confidence interval:** We can give a  $1 - \alpha$  confidence interval for the parameter  $c_0 \beta_0 + c_1 \beta_1$  as

$$c_0 \hat{\beta}_0 + c_1 \hat{\beta}_1 \mp \tilde{\sigma} \left[ \frac{c_0^2}{n} + \frac{(c_0 \bar{x} - c_1)^2}{S_x^2} \right]^{1/2} T_{n-2}^{-1} \left( 1 - \frac{\alpha}{2} \right).$$

### 10.3.5 Mean response

We often want to estimate the *mean* of the probability distribution of  $Y$  for some value of  $x$ .

- The *point estimator* of the mean response

$$\mathbb{E}[Y|x_h] = \beta_0 + \beta_1 x_h$$

when  $x = x_h$  is given by

$$\hat{Y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h.$$

**Lemma 10.3.** *The sampling distribution of  $\hat{Y}_h$  is given by*

$$\hat{Y}_h \sim N \left( \beta_0 + \beta_1 x_h, \sigma^2 \left( \frac{1}{n} + \frac{(x_h - \bar{x})^2}{S_x^2} \right) \right).$$

*Proof.* We first argue normality:

Both  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are *linear combinations* of independent normal random variables  $Y_i$ .

Hence,  $\hat{Y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h$  is also a linear combination of independent normally distributed random variables.

Thus,  $\hat{Y}_h$  is also normally distributed.

Now we look at the mean and variance of  $\hat{Y}_h$ :

Find the expected value of  $\hat{Y}_h$ :

$$\mathbb{E}[\hat{Y}_h] = \mathbb{E}[\hat{\beta}_0 + \hat{\beta}_1 x_h] = \mathbb{E}[\hat{\beta}_0] + \mathbb{E}[\hat{\beta}_1] x_h = \beta_0 + \beta_1 x_h.$$

Note that  $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$ , so that

$$\hat{Y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h = \bar{Y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_h = \bar{Y} + \hat{\beta}_1 (x_h - \bar{x}).$$

Note that  $\hat{\beta}_1$  and  $\bar{Y}$  are *uncorrelated*: indeed setting

$$w_i = \frac{x_i - \bar{x}}{S_x^2}$$

and noting that  $\sum_{i=1}^n w_i = 0$  we obtain

$$\text{Cov} \left( \sum_{i=1}^n w_i Y_i, \sum_{i=1}^n \frac{1}{n} Y_i \right) = \sum_{i=1}^n \frac{w_i}{n} \sigma^2 = \frac{\sigma^2}{n} \sum_{i=1}^n w_i = 0.$$

Therefore,

$$\begin{aligned} \text{Var}(\hat{Y}_h) &= \text{Var}(\bar{Y}) + (x_h - \bar{x})^2 \text{Var}(\hat{\beta}_1) \\ &= \frac{\sigma^2}{n} + (x_h - \bar{x})^2 \frac{\sigma^2}{S_x^2}. \end{aligned}$$

□

When we do not know  $\sigma^2$  we estimate it using  $\tilde{\sigma}^2$ . Thus, the *estimated variance* of  $\hat{Y}_h$  is given by

$$se^2(\hat{Y}_h) = \tilde{\sigma}^2 \left( \frac{1}{n} + \frac{(x_h - \bar{x})^2}{S_x^2} \right).$$

The variance of  $\hat{Y}_h$  is *smallest* when  $x_h = \bar{x}$ .

When  $x_h = 0$ , the variance of reduces to the variance of  $\hat{\beta}_0$ .

- The sampling distribution for the studentized statistic is given by

$$\frac{\hat{Y}_h - \mathbb{E}(\hat{Y}_h)}{se(\hat{Y}_h)} \sim t_{n-2}.$$

All inference regarding  $\mathbb{E}[\hat{Y}_h]$  are carried out using the  $t$ -distribution. A  $(1 - \alpha)$ -CI for the *mean response* when  $x = x_h$  is

$$\hat{Y}_h \mp t_{n-2}(1 - \alpha/2) se(\hat{Y}_h).$$

### 10.3.6 Prediction interval

A CI for a *future* observation is called a *prediction interval*.

Consider the prediction of a new observation  $Y$  corresponding to a given level  $x$  of the predictor.

Suppose  $x = x_h$  and the new observation is denoted  $Y_{h(new)}$ .

Note that  $\mathbb{E}[\hat{Y}_h]$  is the *mean* of the distribution of  $Y|X = x_h$ .

$Y_{h(new)}$  represents the prediction of an *individual outcome* drawn from the distribution of  $Y|X = x_h$ , i.e.,

$$Y_{h(new)} = \beta_0 + \beta_1 x_h + \epsilon_{new},$$

where  $\epsilon_{new}$  is independent of our data.

- The *point estimate* will be the *same* for both  $\mathbb{E}[Y|x_h]$  and  $Y_{h(new)}$ , namely

$$\hat{Y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h.$$

However, the variance is *larger* when predicting an individual outcome due to the *additional variation* of an individual about the mean.

- When constructing prediction limits for  $Y_{h(new)}$  we must take into consideration two sources of variation:
  - Variation in the *mean* of  $Y$ .
  - Variation around the mean.
- The *sampling* distribution of the studentized statistic:

$$\frac{Y_{h(new)} - \hat{Y}_h}{\text{se}(Y_{h(new)} - \hat{Y}_h)} \sim t_{n-2}.$$

All inference regarding  $Y_{h(new)}$  are carried out using the  $t$ -distribution:

$$\text{Var}(Y_{h(new)} - \hat{Y}_h) = \text{Var}(Y_{h(new)}) + \text{Var}(\hat{Y}_h) = \sigma^2 \left\{ 1 + \frac{1}{n} + \frac{(x_h - \bar{x})^2}{S_x^2} \right\}.$$

$$\text{Thus, } \text{se}_{pred}^2 = \text{se}^2(Y_{h(new)} - \hat{Y}_h) = \tilde{\sigma}^2 \left\{ 1 + \frac{1}{n} + \frac{(x_h - \bar{x})^2}{S_x^2} \right\}.$$

Using this result,  $(1 - \alpha)$  *prediction interval* for a new observation  $Y_{h(new)}$  is

$$\hat{Y}_h \mp t_{n-2}(1 - \alpha/2) \text{ se}_{pred}.$$

### 10.3.7 Inference about both $\beta_0$ and $\beta_1$ simultaneously

Suppose that  $\beta_0^*$  and  $\beta_1^*$  are given numbers and we are interested in testing the following hypothesis:

$$H_0 : \beta_0 = \beta_0^* \text{ and } \beta_1 = \beta_1^* \quad \text{versus} \quad H_1 : \text{at least one is different} \quad (29)$$

We shall derive the likelihood ratio test for (29).

The likelihood function (27), when maximized under the unconstrained space yields the MLEs  $\hat{\beta}_1, \hat{\beta}_0, \hat{\sigma}^2$ .

Under the constrained space,  $\beta_0$  and  $\beta_1$  are fixed at  $\beta_0^*$  and  $\beta_1^*$ , and so

$$\hat{\sigma}_0^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \beta_0^* - \beta_1^* x_i)^2.$$

The likelihood statistic reduces to

$$\begin{aligned} \Lambda(\mathbf{Y}, \mathbf{x}) &= \frac{\sup_{\sigma^2} L(\beta_0^*, \beta_1^*, \sigma^2)}{\sup_{\beta_0, \beta_1, \sigma^2} L(\beta_0, \beta_1, \sigma^2)} \\ &= \frac{\frac{1}{(2\pi\hat{\sigma}_0^2)^{n/2}} \exp\left(-\frac{1}{2\hat{\sigma}_0^2} \sum_{i=1}^n (Y_i - \beta_0^* - \beta_1^* x_i)^2\right)}{\frac{1}{(2\pi\hat{\sigma}^2)^{n/2}} \exp\left(-\frac{1}{2\hat{\sigma}^2} \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2\right)} \\ &= \left(\frac{\hat{\sigma}^2}{\hat{\sigma}_0^2}\right)^{n/2} \\ &= \left[\frac{\sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{\sum_{i=1}^n (Y_i - \beta_0^* - \beta_1^* x_i)^2}\right]^{n/2}. \end{aligned}$$

The LRT procedure specifies rejecting  $H_0$  when

$$\Lambda(\mathbf{Y}, \mathbf{x}) \leq k,$$

for some  $k$ , chosen given the level condition.

### 10.3.8 Examples

**Example 10.4.** Let us first do a sanity check if our formulas align with our intuition.

$i$	$x_i$	$Y_i$
1	1	1
2	2	2

We realise that for the sample  $(x_1, Y_1), (x_2, Y_2)$ , there exists exactly one line connecting the two points, namely  $y = 0 + x$ . On the other hand we calculate  $\bar{x} = 1.5$ ,  $\bar{y} = 1.5$ , and

$$\hat{\beta}_1 = \frac{(1 - 1.5) \cdot 1 + (2 - 1.5) \cdot 2}{(1 - 1.5)^2 + (2 - 1.5)^2} = \frac{-0.5 \cdot 1 + 0.5 \cdot 2}{0.5^2 + 0.5^2} = \frac{0.5}{0.5} = 1,$$

$$\hat{\beta}_0 = 1.5 - 1 \cdot 1.5 = 0.$$

In conclusion we also obtain the regression line  $y = \hat{\beta}_0 + \hat{\beta}_1 x = 0 + x$ .

**Example 10.5.** Take Table 11.1 from the book on page 690:

$i$	$x_i$	$Y_i$
1	1.9	0.7
2	0.8	-1.0
3	1.1	-0.2
4	0.1	-1.2
5	-0.1	-0.1
6	4.4	3.4
7	4.6	0.0
8	1.6	0.8
9	5.5	3.7
10	3.4	2.0

We find that

$$\hat{\beta}_1 = 0.685$$

$$\hat{\beta}_0 = -0.786$$

and the regression line is  $y = -0.786 + 0.685x$ .

We can also calculate

$$\tilde{\sigma}^2 = \frac{1}{8} \sum_{i=1}^{10} (Y_i - \hat{Y}_i)^2 = 1.172,$$

$$se^2(\hat{\beta}_0) = \tilde{\sigma}^2 \left( \frac{1}{8} + \frac{\bar{x}^2}{S_x^2} \right) = 0.294,$$

$$se^2(\hat{\beta}_1) = \frac{\tilde{\sigma}^2}{S_x^2} = 0.0325.$$

We can test

$$H_0 : \beta = 1 \quad \text{vs.} \quad H_1 : \beta \neq 0$$

using

$$t_{\text{obs}} = \frac{\hat{\beta}_1}{se(\hat{\beta}_1)} = \frac{0.685}{\sqrt{0.0325}} = 3.800,$$
$$P(|T| \geq |t_{\text{obs}}|) = 0.005.$$

So we reject  $H_0$  for all confidence levels less than 99.5%.

See also the corresponding R code online for an illustration.