



UNIVERSITÀ DEGLI STUDI DI TORINO  
DEPARTMENT OF MATHEMATICS

M.Sc. Stochastics and Data Science

Statistical Inference

Prof:  
Stefano Favaro

Notes by:  
Nicolas Destefanis

2018/2019



---

## Preface

These lecture notes are based on the semester-long course on Statistical Inference, given in 2018-2019 at the University of Torino.

Please, feel free to send me corrections and feedback.

*Nicolas Destefanis*

[nicolas.destefanis@edu.unito.it](mailto:nicolas.destefanis@edu.unito.it)



---

## Contents

<b>1 Week I</b>	1
1.1 26/09/18	1
1.1.1 Remarkable Inequalities	1
1.1.2 Continuous Distributions	2
1.1.3 Discrete Distributions	4
1.2 27/09/18	6
1.2.1 Matrix Algebra	6
1.2.2 Random Vectors	7
1.2.3 Multivariate Gaussian Distribution	8
<b>2 Week II</b>	11
2.1 02/10/18	11
2.2 03/10/18	15
2.2.1 Properties of a Random Sample	15
2.3 04/10/18	20
2.3.1 Concentrations of Measures	20
2.3.2 Concentrations for Sums of Independent Random Variables	22
<b>3 Week III</b>	25
3.1 09/10/18	25
3.1.1 Concentration Inequalities	25
3.2 10/10/18	31
3.2.1 Likelihood Function	31
3.3 11/10/18	33
3.3.1 Likelihood Function Method: Assumptions	33
3.3.2 The Exponential Family	35
3.3.3 The <i>Natural</i> Exponential Family	37
<b>4 Week IV</b>	39
4.1 16/10/18	39
4.1.1 Sufficiency Principle	39
4.2 17/10/18	42
4.2.1 Savage's Factorisation Theorem	42
4.2.2 Minimal Sufficient Statistic	44

## X Contents

4.3	18/10/18	46
4.3.1	Methods of Evaluating Estimators	47
<b>5</b>	<b>Week V</b>	49
5.1	23/10/18	49
5.1.1	The Cramér-Rao Lower Bound	49
5.1.2	Convergence	53
5.2	24/10/18	56
5.2.1	Maximum Likelihood Estimator	56
5.3	25/10/18	61
5.3.1	Method of Moments	61
<b>6</b>	<b>Week VI</b>	67
6.1	30/10/18	67
6.1.1	Exercises	67
6.2	31/10/18	72
6.2.1	Exercises	72
<b>7</b>	<b>Week VII</b>	75
7.1	06/11/18	75
7.1.1	Properties of the Maximum Likelihood Approach	75
7.1.2	Delta Method (a generalisation of the Central Limit Theorem for functionals)	77
7.2	07/11/18	79
7.2.1	Maximum Likelihood Approach	79
7.2.2	Asymptotic Properties of the Maximum Likelihood Estimator	82
7.3	08/11/18	84
<b>8</b>	<b>Week VIII</b>	87
8.1	13/11/18	87
8.1.1	Hypothesis Testing	87
8.2	14/11/18	91
8.2.1	Neyman-Pearson Lemma	91
8.3	15/11/18	96
<b>9</b>	<b>Week IX</b>	99
9.1	20/11/18	99
9.1.1	Rubin Test	99
9.1.2	Likelihood Ratio Test	100
9.2	21/11/18	104
<b>10</b>	<b>Week X</b>	107
10.1	27/11/18	107
10.1.1	Linear Models	107
10.1.2	Second Order Assumptions of the Linear Regression Model	109
10.2	28/11/18	111
10.2.1	Least Square Method	111
10.2.2	The Projection Matrix Operator	112
10.2.3	Estimation of $\beta$	113

## Contents XI

10.3 29/11/18	115
10.3.1 Estimation of $\sigma^2$	115
10.3.2 Sample Correlation and Coefficient of Determination	118
<b>11 Week XI</b>	121
11.1 04/12/18	121
11.1.1 Constrained Estimates	121
11.1.2 Constrained Optimisation: Lagrange Multipliers	122
11.1.3 Gaussian Theory	124
11.2 05/12/18	126
11.2.1 The Fisher-Cochran Theorem	126
<b>12 Week XII</b>	129
12.1 11/12/18	129
12.1.1 Generalised Linear Models	129
12.2 12/12/18	134
12.3 13/12/18	139
12.3.1 Solving the Likelihood Equations	139



## 1

## Week I

## 1.1 26/09/18

## 1.1.1 Remarkable Inequalities

Denote by  $X$  a random variable with finite second moment, and let  $\sigma = \sqrt{\text{Var}(X)}$ . The **Chebyshev's inequality** states that, for any  $\varepsilon > 0$

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq \varepsilon) \leq \frac{\text{Var}(X)}{\varepsilon^2}.$$

It is a particular case of the **Markov's inequality**: if  $Y$  is a non-negative random variable with finite expected value, then

$$\mathbb{P}(Y \geq t) \leq \frac{\mathbb{E}[Y]}{t}, \quad \forall t > 0. \quad (1.1)$$

Also, if  $X$  and  $Y$  are two different random variables with finite second moment, then by Hölder's inequality (taking  $p = q = 2$ ), we get

$$(\mathbb{E}[XY])^2 \leq \mathbb{E}[X^2] \mathbb{E}[Y^2]. \quad (1.2)$$

**Theorem 1.1.** *If  $X$  is a random variable taking values in some set  $I$ , with  $\mathbb{E}[X] = \mu$ , and  $f$  is a convex function on  $I$ , then we have the following inequality:*

$$f(X) \geq f(\mu) + b(X - \mu),$$

*held with probability 1 for some choice of  $b$ . By integrating both sides of the inequality with respect to the distribution of  $X$ , we obtain*

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X]),$$

*known as **Jensen's inequality**.*

2 1 Week I

### 1.1.2 Continuous Distributions

#### The Gaussian Distribution

**Definition 1.2.** A continuous random variable  $Y$  is said to have a **Gaussian distribution** with parameters  $\mu$  and  $\sigma^2$  if the density function at  $t$  is given by

$$f(t; \mu, \sigma^2) := \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(t-\mu)^2}{2\sigma^2}\right\} \mathbb{1}_{\mathbb{R}}(t).$$

If  $Y \sim \mathcal{N}(0, 1)$ , we say that  $Y$  is the **standardised Gaussian random variable**, with density and cumulative distribution functions denoted by

$$\phi(t) := \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} \mathbb{1}_{\mathbb{R}}(t) \quad \text{and} \quad \Phi(t) := \int_{-\infty}^t \phi(y) dy$$

respectively.

*Remark 1.3.* The Gaussian distribution has some interesting properties.

1. If  $Y \sim \mathcal{N}(\mu, \sigma^2)$ , its characteristic function and moment generating function are

$$\mathbb{E}[e^{itY}] = \exp\left\{it\mu - \frac{1}{2}\sigma^2 t^2\right\} \quad \text{and} \quad \mathbb{E}[e^{tY}] = \exp\left\{t\mu + \frac{1}{2}\sigma^2 t^2\right\},$$

respectively. By taking the derivatives, using the characteristic function we can obtain the *moments* of  $Y$ . Also, if we take

$$\mathbb{E}[(Y - \mu)^k] = \begin{cases} 0 & \text{for } k \text{ odd} \\ 1 \cdot 3 \dots (k-1)\sigma^k & \text{for } k \text{ even} \end{cases}.$$

2. If  $Y \sim \mathcal{N}(\mu, \sigma^2)$  and  $a, b$  are constants, then

$$a + bY \sim \mathcal{N}(a + b\mu, b^2\sigma^2).$$

This means that the entire family of distributions can be generated by linear transformations, starting from any member of the family: this property is known as **location scale**.

3. If  $Y_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  and  $Y_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$  are independent, then

$$Y_1 + Y_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2).$$

This result extends to any linear combination of Gaussian random variables, meaning if  $Y_1, \dots, Y_N \sim \mathcal{N}(\mu_i, \sigma_i^2)$  and  $Y_i \perp Y_j$  for every  $i \neq j$ , then

$$\sum_{i=1}^n (a_i Y_i + b_i) \sim \mathcal{N}\left(\sum_{i=1}^n (a_i \mu_i + b_i), \sum_{i=1}^n a_i^2 \sigma_i^2\right).$$

### The Uniform Distribution

**Definition 1.4.** A continuous random variable  $Y$  is said to have a **uniform distribution** in  $(a, b)$  if the density function at  $t$  is given by

$$f(t; a, b) := \frac{1}{b-a} \mathbb{1}_{(a,b)}(t).$$

*Remark 1.5.* An important result concerning the uniform distribution is the following: if  $Z$  is a continuous random variable with distribution function given by  $F(\cdot)$ , then the random variable  $W$  defined as  $W = F(Z)$  is uniformly distributed on  $(0, 1)$ . If we wanted to find the probability of  $W$ , we would compute

$$\mathbb{P}(W \leq t) = \mathbb{P}(F(Z) \leq t) = \mathbb{P}(Z \leq F^{-1}(t)) = F(F^{-1}(t)) = t,$$

which is the distribution function of  $\mathcal{U}(0, 1)$ .

### The Gamma Distribution

**Definition 1.6.** First of all, we define the **Gamma function** as

$$\Gamma(x) := \int_0^{+\infty} t^{x-1} e^{-t} dt.$$

Note that:

1. Using the integration by parts, we get  $\Gamma(x+1) = x\Gamma(x)$ . In particular, when  $n \in \mathbb{N}$ ,  $\Gamma(n+1) = n!$ , with  $\Gamma(1) = 1$  and  $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ .
2.  $n! \sim \sqrt{\pi n} \frac{n^{n+1}}{2} e^{-n}$  as  $n \rightarrow \infty$ .

**Definition 1.7.** A continuous random variable  $Y$  is said to have a **Gamma distribution** with shape parameter  $w$  and scale parameter  $\lambda$  if the density function at  $t$  is given by

$$f(t; w, \lambda) := \frac{\lambda^w}{\Gamma(w)} t^{w-1} e^{-\lambda t} \mathbb{1}_{(0,+\infty)}(t).$$

**Exercise 1.8.** Derive the moment formula for the Gamma distribution. Also, try to extract a Gamma distribution from a negative exponential.

**Exercise 1.9.** If  $Y_1 \sim \Gamma(w_1, \lambda)$  and  $Y_2 \sim \Gamma(w_2, \lambda)$ , with  $Y_1 \perp Y_2$ , then

$$Y_1 + Y_2 \sim \Gamma(w_1 + w_2, \lambda).$$

### The Beta Distribution

**Definition 1.10.** First of all, we define the **Beta function** as

$$B(p, q) := \int_0^1 x^{p-1} (1-x)^{q-1} dx.$$

Having said that, a continuous random variable  $Y$  is said to have a **Beta distribution** with parameter  $p$  and  $q$  if the density function at  $t$  is given by

$$f(t; p, q) := \frac{t^{p-1} (1-t)^{q-1}}{B(p, q)} \mathbb{1}_{(0,1)}(t).$$

### 1.1.3 Discrete Distributions

#### The Binomial Distribution

**Definition 1.11.** A discrete random variable  $Y$  is said to have a **Binomial distribution** with parameters  $n$  and  $p$  if the mass function at  $t$  is given by

$$\mathbb{P}(Y = t) := \binom{n}{t} p^t (1-p)^{n-t} \mathbb{1}_{\{0, \dots, n\}}(t)$$

for  $p \in (0, 1)$ . A special case of the Binomial distribution is the Bernoulli distribution, where  $n = 1$ .

#### The Hypergeometric Distribution

**Definition 1.12.** A discrete random variable  $Y$  is said to have a **Hypergeometric distribution** with parameters  $N$  and  $M$  if the mass function at  $t$  is given by

$$\mathbb{P}(Y = t) := \frac{\binom{M}{t} \binom{N-M}{n-t}}{\binom{N}{n}}$$

where  $n, N, M$  are positive integers such that  $N \geq M$  and  $N \geq n$ , and  $t$  is an integer such that

$$\max\{0, n - N + M\} \leq t \leq \min\{n, M\}.$$

*Remark 1.13.* This distribution provides the probability of throwing  $t$  success balls when  $n$  balls are drawn *without replacement* from an urn containing  $M$  success balls and  $N - M$  failure balls.

#### The Poisson Distribution

The **Poisson distribution** is a widely applied discrete distribution, and can serve as a model for a number of different types of experiments. For example, if we are modelling a phenomenon in which we are waiting for an occurrence (such as waiting for a bus, waiting for costumers to arrive in a bank), the number of occurrences in a given time interval can sometimes be modelled by the Poisson distribution. One of the basic assumptions on which the Poisson distribution is built is that, for small time intervals, the probability of an arrival is proportional to the length of waiting time. This makes it a reasonable model for situations like those indicated above. For example, it makes sense to assume that the longer we wait, the more likely it is that a customer will enter the bank.

Another area of application is in spatial distributions, where, for example, the Poisson may be used to model the distribution of bomb hits in an area or the distribution of fish in a lake.

**Definition 1.14.** A random variable  $Y$ , taking values in the non-negative integers, has a **Poisson distribution**, with a single parameter  $\lambda$  (sometimes called the *intensity parameter*), if

$$\mathbb{P}(Y = t) := \frac{e^{-\lambda} \lambda^t}{t!}, \quad t \in \mathbb{N}.$$

### The Negative Binomial Distribution

The Binomial distribution counts the number of successes in a fixed number of Bernoulli trials. Suppose that, instead, we count the number of Bernoulli trials required to get a fixed number of successes. This latter formulations leads to the Negative Binomial distribution.

**Definition 1.15.** *In a sequence of independent Bernoulli( $p$ ) trials, let the random variable  $Y$  denote the trial at which the  $r$ -th success occurs, where  $r$  is a fixed integer. Then*

$$\mathbb{P}(Y = t) := \binom{t-1}{r-1} p^r (1-p)^{t-r}, \quad t = r, r+1, \dots \quad (1.3)$$

and we say that  $Y$  has a **Negative Binomial Distribution** with parameters  $(r, p)$ .

*Remark 1.16.* The derivation of (1.3) follows quickly from the Binomial distribution. The event  $\{Y = t\}$  can occur only if there are exactly  $r-1$  successes in the first  $t-1$  trials, and a success on the  $t$ -th trial. The probability of  $r-1$  successes in  $t-1$  trials is the binomial probability

$$\binom{t-1}{r-1} p^{r-1} (1-p)^{t-r}$$

and with probability  $p$  there is a success on the  $t$ -th trial. Multiplying these probabilities, gives us (1.3).

The Negative Binomial distributions is sometimes defined in terms of the “number of failures before the  $r$ -th success”. This formulation is statistically equivalent to the one given above in terms of “trial at which the  $r$ -th success occurs”. Therefore, an alternative form of the negative binomial distribution is

$$\mathbb{P}(Y = t) := \binom{r+t-1}{t} p^r (1-p)^t, \quad t \in \mathbb{N}.$$

Unless otherwise noted, when we refer to the negative binomial distribution with parameters  $(r, p)$ , we will use this probability mass function.

In the case that  $r$  is an integer, the Negative Binomial is called a **Pascal distribution**, and it has the following interpretation: consider an infinite sequence of trials whose outcome is *success* with probability  $p$  and *failure* with probability  $1-p$ , and assume that the trials are independent. Then, the above probability mass gives us the probability of expecting  $t$  failures before achieving  $r$  successes. If  $r = 1$ , we obtain the **Geometric distribution**.

### The Geometric Distribution

**Definition 1.17.** *The **Geometric distribution** is the simplest of the waiting time distributions, and is a special case of the negative binomial distribution. If we set  $r = 1$  in (1.3), we obtain*

$$\mathbb{P}(Y = t) := p(1-p)^{t-1}, \quad t \in \mathbb{N}_+,$$

which defines the probability mass function of a Geometric random variable with success probability  $p$ . It can be interpreted as the trial at which the first success occurs, so we are “waiting for a success”.

6 1 Week I

## 1.2 27/09/18

## 1.2.1 Matrix Algebra

- Denote by  $I_k$  the *identity matrix* of order  $k$ , by  $\mathbf{1}_n$  the  $n \times 1$  vector with all elements equal to 1 and by  $\mathbf{0}_n$  the  $n \times 1$  vector with all elements equal to 0.
- For two conformable matrices  $A$  and  $B$ , we have

$$\det(AB) = \det(A) \det(B).$$

- If  $\det(A) \neq 0$ , then  $A$  is non-singular, and there exist an *inverse matrix*  $A^{-1}$  such that  $AA^{-1} = A^{-1}A = I_n$ . Moreover,

$$({}^tA)^{-1} = {}^t(A^{-1}) \quad \text{and} \quad (AB)^{-1} = B^{-1}A^{-1}.$$

- A matrix  $A$  is *positive definite* if  $A$  is symmetric and  ${}^t\mathbf{u}A\mathbf{u} > 0$  for all non-null vectors  $\mathbf{u} \in \mathbb{R}^k$ .
- For any two conformable matrices we have  $\text{tr}(AB) = \text{tr}(A)\text{tr}(B)$ .
- If  $A$  is idempotent, then  $\text{rk}(A) = \text{tr}(A)$ .
- Taking into account that the matrices  $A, B, C, D$  are conformable, we have the following identities:

$$(A + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}$$

and

$$(A + b {}^td)^{-1} = A^{-1} - \frac{A^{-1}b {}^tdA^{-1}}{1 + {}^tdA^{-1}b}.$$

- **Spectral Decomposition:** if  $A$  is a symmetric  $k \times k$  matrix, there exist  $\lambda_1, \dots, \lambda_k \in \mathbb{R}$  and an orthogonal matrix  $Q$  such that

$$A = QA {}^tQ,$$

where  $A = \text{diag}(\lambda_1, \dots, \lambda_k)$ ,  $\lambda_1, \dots, \lambda_k$  are the *eigenvalues* of  $A$  and the column of order  $j$  of  $Q$  is called the *eigenvector corresponding to*  $\lambda_j$ . Also

$$\det(A) = \det(A) = \prod_{i=1}^k \lambda_i.$$

- Given a positive semidefinite matrix  $A$ , we define the *square root of a matrix*  $A$  as the matrix  $B$  such that  $A = B {}^tB$ . In this case  $A = \text{diag}(a_1, \dots, a_k)$ , therefore it is natural to set  $B = \text{diag}(\sqrt{a_1}, \dots, \sqrt{a_k})$ . In the general case of a positive semidefinite matrix  $A$ , we put  $B = QA^{\frac{1}{2}} = Q \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_k})$  such that

$$B {}^tB = \left(QA^{\frac{1}{2}}\right) \left(A^{\frac{1}{2}} {}^tQ\right) = QA^{\frac{1}{2}} {}^tQ = A;$$

we then say that  $B$  is the square root of  $A$ . Note that if  $A > 0$

$$B^{-1} = (QA^{\frac{1}{2}})^{-1} = A^{-\frac{1}{2}} {}^tQ.$$

### 1.2.2 Random Vectors

**Definition 1.18.** If  $X_1, \dots, X_k$  are random variables defined on the same probability space, we define a **random vector** (or **multivariate random variable**) as

$$\mathbf{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_k \end{pmatrix}.$$

The **mean value** of  $\mathbf{X}$  is obtained by forming the vector of mean values with components

$$\mathbb{E}[\mathbf{X}] = \begin{pmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_k] \end{pmatrix}.$$

Similarly, we define the **variance matrix** (or **dispersion matrix**) as

$$\text{Var}(\mathbf{X}) = \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_k) \\ \vdots & \ddots & & \vdots \\ \text{Cov}(X_k, X_1) & \text{Cov}(X_k, X_2) & \dots & \text{Var}(X_k) \end{pmatrix}.$$

We also define the **correlation matrix** whose generic element is

$$\text{Cor}(X_i, X_j) := \frac{\text{Cov}(X_i, X_j)}{\sqrt{\text{Var}(X_i) \text{Var}(X_j)}}.$$

**Lemma 1.19.** If  $A = (a_{ij})$  is a  $n \times k$  matrix and  $\mathbf{b} = {}^t(b_1, \dots, b_n)$  an  $n \times 1$  vector, define

$$\mathbf{Y} = A\mathbf{X} + \mathbf{b}$$

where  $\mathbf{X} = {}^t(X_1, \dots, X_k)$ ,  $\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}$  and  $\text{Var}(\mathbf{X}) = \mathbb{V}$ . Then,

$$\mathbb{E}[\mathbf{Y}] = A\boldsymbol{\mu} + \mathbf{b} \quad \text{and} \quad \text{Var}(\mathbf{Y}) = A\mathbb{V}{}^tA.$$

**Lemma 1.20.** The variance matrix  $\mathbb{V}$  of the random vector  $\mathbf{X}$  is positive semidefinite, and is positive definite if there exist no vector such that  ${}^t\mathbf{b}\mathbf{X}$  is a degenerate random variable (meaning, a random variable which takes a single value with probability 1).

**Lemma 1.21.** If  $\mathbb{V}(\mathbf{X}) > 0$ , there exist a square matrix  $C$  such that  $\mathbf{Y} := C\mathbf{X}$  has uncorrelated components with unit variance.

**Lemma 1.22.** Let  $A = (a_{ij})$  be a square matrix of order  $k$ , and let  $\mathbf{X}$  be a random vector such that  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}]$  and  $\mathbb{V} = \text{Var}(\mathbf{X})$ . Then,

$$\mathbb{E}[\mathbf{X}A\mathbf{X}] = {}^t\boldsymbol{\mu}A\boldsymbol{\mu} + \text{tr}(A\mathbb{V}).$$

*Proof.* Now,

$$\mathbb{E}[\mathbf{X}A\mathbf{X}] = \mathbb{E} \left[ \sum_{i=1}^k \sum_{j=1}^k X_i a_{ij} X_j \right]$$

8 1 Week I

$$\begin{aligned}
&= \sum_{i=1}^k \sum_{j=1}^k a_{ij} \mathbb{E}[X_i X_j] \\
&= \sum_{i=1}^k \sum_{j=1}^k a_{ij} (\mu_i \mu_j + \mathbb{V}_{ij})
\end{aligned}$$

where  $\mu_i$  and  $\mathbb{V}_{ij}$  are the generic elements of  $\boldsymbol{\mu}$  and  $\mathbb{V}$  respectively. Therefore,

$$\begin{aligned}
\mathbb{E}[\mathbf{t} \mathbf{X} \mathbf{A} \mathbf{X}] &= \sum_{i=1}^k \sum_{j=1}^k a_{ij} \mu_i \mu_j + \sum_{i=1}^k \sum_{j=1}^k a_{ij} \mathbb{V}_{ij} \\
&= \mathbf{t} \boldsymbol{\mu} \mathbf{A} \boldsymbol{\mu} + \sum_{i=1}^k (\mathbf{A} \mathbb{V})_{ii} \\
&= \mathbf{t} \boldsymbol{\mu} \mathbf{A} \boldsymbol{\mu} + \text{tr}(\mathbf{A} \mathbb{V}),
\end{aligned}$$

where  $(\mathbf{A} \mathbb{V})_{ii}$  is the generic element of the main diagonal of  $\mathbf{A} \mathbb{V}$ . □

### 1.2.3 Multivariate Gaussian Distribution

Consider a vector  $\mathbf{Z} = {}^t(Z_1, \dots, Z_k)$  where  $Z_1, \dots, Z_k$  are independent and identically distributed standard Gaussian random variables. Now, set

$$\mathbf{Y} = \mathbf{A} \mathbf{Z} + \boldsymbol{\mu},$$

where  $\mathbf{A}$  is a non-singular  $k \times k$  matrix, and  $\boldsymbol{\mu}$  is a  $k \times 1$  vector. It is natural to think about  $\mathbf{Y}$  as a  $k$ -dimensional generalisation of the Gaussian distribution. We wish to find its distribution.

We start from the standard case

$$f_{\mathbf{Z}}(\mathbf{t}) = \frac{1}{(2\pi)^{\frac{k}{2}}} \exp\left\{-\frac{1}{2} \mathbf{t} \mathbf{t}\right\}.$$

Then since  $\mathbf{Z} = \mathbf{A}^{-1}(\mathbf{Y} - \boldsymbol{\mu})$ , we find the *Jacobian matrix* of the transformation, which is (considering all the partial derivatives)

$$\det \left| \frac{dz_i}{dy_j} \right| = \det(\mathbf{A}^{-1}) = (\det(\mathbb{V}))^{-\frac{1}{2}}.$$

Taking into account that

$$|\mathbb{V}| = |\mathbf{A} \mathbf{t} \mathbf{A}| = |\mathbf{A}|^2$$

and setting  $\mathbf{y} = \mathbf{A} \mathbf{t} + \boldsymbol{\mu}$  we obtain that

$$\mathbf{t} \mathbf{t} = {}^t[\mathbf{A}^{-1}(\mathbf{y} - \boldsymbol{\mu})][\mathbf{A}^{-1}(\mathbf{y} - \boldsymbol{\mu})] = {}^t(\mathbf{y} - \boldsymbol{\mu}) \mathbb{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}).$$

Therefore, the density of  $\mathbf{Y}$  is

$$f_{\mathbf{Y}}(\mathbf{y}) := \frac{1}{(2\pi)^{\frac{k}{2}} (\det \mathbb{V})^{\frac{1}{2}}} \exp\left\{-\frac{1}{2} {}^t(\mathbf{y} - \boldsymbol{\mu}) \mathbb{V}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right\}. \quad (1.4)$$

**Definition 1.23.** We say that the random variable  $\mathbf{Y} = {}^t(Y_1, \dots, Y_k)$  with density function  $f_{\mathbf{Y}}(\cdot)$  is a **multivariate Gaussian random variable** with mean  $\boldsymbol{\mu}$  and variance  $\mathbb{V}$ , or

$$\mathbf{Y} \sim \mathcal{N}_k(\boldsymbol{\mu}, \mathbb{V}).$$

**Exercise 1.24.** Take  $X_1, X_2$  independent random variables such that

$$X_1, X_2 \sim \mathcal{U}([0, 1]).$$

Find the distribution of  $Y = X_1 + X_2$ ,  $Z = \frac{X_2}{X_1}$  and  $W = X_1 X_2$ .

**Exercise 1.25.** Take three independent negative exponential random variables, with parameter  $\frac{1}{2}$ . Find the distribution of  $U = \frac{Y}{X}$  and  $W = X + Y + Z$ .

**Exercise 1.26.** Let  $(X, Y)$  be a bivariate random variable such that the marginal distribution of  $X \sim \mathcal{U}([-1, 1])$  and the conditional distribution

$$Y | \{X = x\} \sim \mathcal{U}((x, x + 1)).$$

Find the distribution of  $Z = -\log(Y - X)$ .

**Exercise 1.27.** Let  $A$  and  $B$  be two independent random variables such that

$$A, B \sim \mathcal{U}([0, h]).$$

Compute the probability that the equation  $Z^2 - 2AZ + B = 0$  admits complex roots.

**Exercise 1.28.** Let  $X \sim \Gamma(r, 1)$  and  $Y \sim \Gamma(s, 1)$  be two independent random variables. Find the distribution of  $\left(\frac{X}{X+Y}, X+Y\right)$ ,  $Z = \frac{X}{X+Y}$  and  $W = X+Y$ .



## 2

## Week II

## 2.1 02/10/18

## Moment Generating Function of a Multivariate Gaussian

Last week, we derived the multivariate Gaussian distribution (1.4) starting from a vector of independent and identically distributed standard Gaussian random variables. We wish to compute the moment generating function of  $\mathbf{Y}$ .

**Lemma 2.1.** Consider a positive definite  $n \times n$  matrix  $A$ , a  $n \times 1$  vector  $\mathbf{b}$  and

$$I := \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{\frac{n}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{y}^t A \mathbf{y} - 2 \mathbf{b}^t \mathbf{y}) \right\} d\mathbf{y}. \quad (2.1)$$

Then,

$$I = \frac{\exp \left\{ \frac{1}{2} \mathbf{b}^t A^{-1} \mathbf{b} \right\}}{\det(A)^{\frac{1}{2}}}. \quad (2.2)$$

*Proof.* Let  $\boldsymbol{\mu} = A^{-1} \mathbf{b}$  and expand, within the integral,  $\exp \{ \cdot \}$  by adding and subtracting

$$\frac{1}{2} \mathbf{b}^t A^{-1} \mathbf{b} :$$

$$\begin{aligned} I &= \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{\frac{n}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{y}^t A \mathbf{y} - 2 \mathbf{b}^t \mathbf{y}) \right\} d\mathbf{y} \\ &= \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{\frac{n}{2}}} \exp \left\{ -\frac{1}{2} \mathbf{y}^t A \mathbf{y} + \mathbf{b}^t \mathbf{y} + \frac{1}{2} \mathbf{b}^t A^{-1} \mathbf{b} - \frac{1}{2} \mathbf{b}^t A^{-1} \mathbf{b} \right\} d\mathbf{y} \\ &= \exp \left\{ \frac{1}{2} \mathbf{b}^t A^{-1} \mathbf{b} \right\} \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{\frac{n}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{y}^t A \mathbf{y} - 2 \mathbf{b}^t A^{-1} \mathbf{y} + \mathbf{b}^t A^{-1} \mathbf{b}) \right\} d\mathbf{y} \\ &= \det(A^{-1})^{\frac{1}{2}} \exp \left\{ \frac{1}{2} \mathbf{b}^t A^{-1} \mathbf{b} \right\} \int_{\mathbb{R}^n} g(\mathbf{y}) d\mathbf{y} \\ &= \det(A^{-1})^{\frac{1}{2}} \exp \left\{ \frac{1}{2} \mathbf{b}^t A^{-1} A A^{-1} \mathbf{b} \right\} \\ &= \frac{\exp \left\{ \frac{1}{2} \mathbf{b}^t A^{-1} \mathbf{b} \right\}}{\det(A)^{\frac{1}{2}}}, \end{aligned}$$

12 2 Week II

where  $g(\cdot)$  is the density function of the multivariate gaussian (1.4) with  $\mathbb{V}$  replaced by  $A^{-1}$ .  $\square$

Now the computation of the moment generating function of  $\mathbf{Y}$  is immediate:

$$\begin{aligned}\mathbb{E}[\exp\{\mathbf{t}\mathbf{y}\}] &= \int_{\mathbb{R}^n} \exp\{\mathbf{t}\mathbf{y}\} f_{\mathbf{Y}}(\mathbf{y}) d\mathbf{y} \\ &= \exp\left\{-\frac{1}{2} \mathbf{t}\boldsymbol{\mu}\mathbb{V}^{-1}\boldsymbol{\mu}\right\} \int_{\mathbb{R}^n} \frac{\exp\left\{-\frac{1}{2}(\mathbf{t}\mathbb{V}^{-1}\mathbf{y} - 2\mathbf{t}(\mathbb{V}^{-1}\boldsymbol{\mu} + \mathbf{t})\mathbf{y})\right\}}{(2\pi)^{\frac{n}{2}} \det(\mathbb{V})^{\frac{1}{2}}}} d\mathbf{y} \\ &= \exp\left\{-\frac{1}{2} \mathbf{t}\boldsymbol{\mu}\mathbb{V}^{-1}\boldsymbol{\mu}\right\} \exp\left\{\frac{1}{2} \mathbf{t}(\mathbb{V}^{-1}\boldsymbol{\mu} + \mathbf{t})\mathbb{V}(\mathbb{V}^{-1}\boldsymbol{\mu} + \mathbf{t})\right\} \\ &= \exp\left\{-\frac{1}{2} \mathbf{t}\boldsymbol{\mu}\mathbb{V}^{-1}\boldsymbol{\mu} + \frac{1}{2}(\mathbf{t}\boldsymbol{\mu} + \mathbf{t}\mathbb{V}\mathbf{t} + \mathbf{t}\boldsymbol{\mu}\mathbb{V}^{-1}\boldsymbol{\mu} + \mathbf{t}\boldsymbol{\mu}\mathbf{t})\right\} \\ &= \exp\left\{\frac{1}{2} \mathbf{t}\mathbb{V}\mathbf{t} + \mathbf{t}\boldsymbol{\mu}\right\},\end{aligned}$$

where we applied Lemma 2.1 with  $A = \mathbb{V}^{-1}$  and  $\mathbf{b} = \mathbb{V}^{-1}\boldsymbol{\mu} + \mathbf{t}$ . Therefore, the **characteristic function** of  $\mathbf{Y}$  is given by

$$\mathbb{E}[\exp\{i\mathbf{t}\mathbf{y}\}] = \exp\left\{i\mathbf{t}\boldsymbol{\mu} - \frac{1}{2} \mathbf{t}\mathbb{V}\mathbf{t}\right\}.$$

*Remark 2.2.* If we were interested in the *marginal distributions*, we first have to partition our vector  $\mathbf{Y}$  as

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$$

where  $Y_1 \in \mathbb{R}^r$  and  $Y_2 \in \mathbb{R}^{n-r}$  (for  $1 < r < n$ ). We also need to partition the rest of the elements, meaning

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad B = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \quad \text{and} \quad \mathbb{V} = \begin{pmatrix} V_{11} & V_{12} \\ V_{21} & V_{22} \end{pmatrix},$$

with  $V_{21} = V_{12}$ . How can we extract the marginal distribution of  $Y_1$ , for example? Simply using the characteristic function formula imposing  $t_2 = 0$ : we therefore obtain

$$\exp\left\{i\mathbf{t}_1\boldsymbol{\mu}_1 - \frac{1}{2} \mathbf{t}_1 V_{11} \mathbf{t}_1\right\}$$

which is the characteristic function of a multivariate Gaussian random variable with mean  $\boldsymbol{\mu}_1$  and variance  $V_{11}$ . Note that all marginal distributions of multivariate Gaussian are themselves multivariate Gaussians.

### The Chi-squared Distribution

If  $\mathbf{Z} = {}^t(Z_1, \dots, Z_k) \sim \mathcal{N}_k(\mathbf{0}_k, I_k)$ , define  $U_k := {}^t\mathbf{Z}\mathbf{Z}$ : this is a **central Chi-squared random variable** with parameter  $k$  (known as *degrees of freedom*), and we denote it by

$$U_k \sim \chi_k^2.$$

This random variable is closely related to the Gamma random variable, with a certain characterisation. To compute the density function of  $U_k$ , we can start from  $k = 1$ :

$$\begin{aligned} \mathbb{P}(U_1 \leq t) &= \mathbb{P}(Z_1^2 \leq t) \\ &= \mathbb{P}(-\sqrt{t} \leq Z_1 \leq \sqrt{t}) \\ &= 2\mathbb{P}(0 \leq Z_1 \leq \sqrt{t}) \\ &= 2\left(\Phi(\sqrt{t}) - \frac{1}{2}\right), \end{aligned}$$

therefore

$$f_1(t) = \frac{t^{-\frac{1}{2}} e^{-\frac{t}{2}}}{\sqrt{2\pi}} \mathbb{1}_{(0, +\infty)}(t),$$

which is the density function of a Gamma random variable with parameters  $(\frac{1}{2}, \frac{1}{2})$ . Then, by means of the additive property of the Gamma distribution, we have that  $U_k \sim \Gamma(\frac{k}{2}, \frac{1}{2})$ . Basically,

$$f_{U_k}(u) = \frac{1}{2^{\frac{k}{2}} \Gamma(\frac{k}{2})} u^{\frac{k}{2}-1} e^{-\frac{u}{2}} \mathbb{1}_{(0, \infty)}(u).$$

If  $\mathbf{Z} = {}^t(Z_1, \dots, Z_k) \sim \mathcal{N}_k(\boldsymbol{\mu}, I_k)$ , we say that  $U_k = {}^t\mathbf{Z}\mathbf{Z}$  is a **non-central Chi-squared** random variable, with  $k$  degrees of freedom and *non-centrality parameter*  $\delta = {}^t\boldsymbol{\mu}\boldsymbol{\mu}$ . We can use the following shorthand notation:  $U_k \sim \chi_k^2(\delta)$ .

### The $t$ Distribution

If  $Z \sim \mathcal{N}(0, 1)$  and  $U \sim \chi_k^2$  are independent, we say that

$$T := \frac{Z}{\sqrt{U/k}}$$

is a *Student  $t$*  random variable, with  $k$  degrees of freedom.

### The $F$ Distribution

If  $V \sim \chi_m^2$  and  $U \sim \chi_k^2$  are independent, we say that

$$F := \frac{V/m}{U/k}$$

is a *Fisher  $F$*  random variable, with  $(m, k)$  degrees of freedom.

### Multinomial Distribution

The multinomial distribution is a model for the following kind of experiment. It consists of  $n$  independent trials. Each trial results in one of  $k$  distinct possible outcomes. The probability of the  $i$ -th outcome is  $p_i$  on every trial, and  $X_i$  is the count of the number of times the  $i$ -th outcome occurred in the  $n$  trials. For  $k = 2$ , this is just a binomial experiment, in which each trial has  $k = 2$  possible outcomes and  $X_1$  counts the number of “successes” and  $X_2 = n - X_1$  counts the number of “failures” in  $n$  trials. In a general multinomial experiment, there are  $k$  possible outcome to count.

**Definition 2.3.** *The outcome of a trial may be one of  $k$  mutually exclusive events  $E_1, \dots, E_k$ , with associated probabilities given by*

$$p_j = \mathbb{P}(E_j), \quad j = 1, \dots, k \quad \sum_{j=1}^k p_j = 1.$$

If  $n$  independent replicates of the trial take place, denote by  $Y_j$  the number of times that  $E_j$  occurs. We say that the random vector  $\mathbf{Y} = (Y_1, \dots, Y_k)$  is a **multinomial** random variable with parameters  $(n, (p_1, \dots, p_k))$ . The probability mass is given by

$$\mathbb{P}(\mathbf{Y} = \mathbf{t}) := \mathbb{P}((Y_1 = t_1, \dots, Y_k = t_k)) = \frac{n!}{t_1! \dots t_k!} \prod_{i=1}^k p_i^{t_i},$$

where the factor

$$\frac{n!}{t_1! \dots t_k!}$$

is called **multinomial coefficient**. It is the number of ways the  $n$  objects can be divided into  $k$  groups, with  $t_i$  in the  $i$ -th group.

In order to find the marginal distributions, we could apply the following theorem.

**Theorem 2.4.** *Let  $n, k \in \mathbb{N}_+$  and consider a set  $A$  of vectors  $\mathbf{t} = (t_1, \dots, t_k)$  such that each  $t_i$  is a non-negative integer, and  $\sum_{i=1}^k t_i = n$ . Then, for any real numbers  $p_1, \dots, p_k$ ,*

$$(p_1 + \dots + p_k)^n = \sum_{\mathbf{t} \in A} \frac{n!}{t_1! \dots t_k!} p_1^{t_1} \dots p_k^{t_k}.$$

## 2.2 03/10/18

### 2.2.1 Properties of a Random Sample

Often, the data collected in an experiment consists of several observations on a variable of interest. Here, we present a model for data collection that is often used to describe this situation, a model referred to as *random sample*.

**Definition 2.5.** Let  $X \sim f_X(x; \vartheta)$ . Consider  $n$  independent and identically distributed (as  $X$ ) random variables  $X_1, \dots, X_n$ . A **random sample** from  $X$  is a collection of those random variables and it is denoted by

$$\mathbf{X} = (X_1, \dots, X_n),$$

with

$$X_i \stackrel{d}{=} X \quad \text{and} \quad X_i \perp X_j, \quad \forall i \neq j.$$

We denote by  $\mathbf{x} = (x_1, \dots, x_n)$  the **realisation** of the random sample.

We know that the joint probability density function, or probability mass function, of  $\mathbf{X}$  is given by

$$f_{\mathbf{X}}(\mathbf{x}; \vartheta) = \prod_{i=1}^n f_{X_i}(x_i; \vartheta) = \prod_{i=1}^n f_X(x_i; \vartheta).$$

Since  $X_1, \dots, X_n$  are identically distributed, all the marginal densities are the same function.

When a sample  $X_1, \dots, X_n$  is drawn, some summary of the values is usually computed. Any well-defined summary may be expressed mathematically as a function  $T: \mathbb{R}^n \rightarrow \mathbb{R}^m$  (with the assumption that  $m < n$ , usually  $m = 1, 2$ ) whose domain includes the sample space of the random vector  $\mathbf{X}$ . The function  $T$  may be real-valued or vector-valued; thus the summary is a random variable (or vector)  $\mathbf{Y} = T(\mathbf{X})$ . Since the random sample  $\mathbf{X}$  has a simple probabilistic structure (because the  $X_i$ 's are independent and identically distributed), the distribution of  $\mathbf{Y}$  is particularly tractable. Because this distribution is usually derived from the distribution of the variables in the random sample, it is called the *sampling distribution* of  $\mathbf{Y}$ . This distinguishes the probability distribution of  $\mathbf{Y}$  from the distribution of the population, that is, the marginal distribution of each  $X_i$ .

**Definition 2.6.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample of size  $n$  from a population and let  $T(\mathbf{X})$  be a real-valued (or vector-valued) function, whose domain includes the sample space of  $\mathbf{X}$ . Then, the random variable (or random vector)  $\mathbf{Y} = T(\mathbf{X})$  is called a **statistic**. The probability distribution of a statistic  $\mathbf{Y}$  is called the **sampling distribution** of  $\mathbf{Y}$ .

Here are some examples of statistics.

**Definition 2.7.** The **sample mean** is the arithmetic average of the values in a random sample, usually denoted by

$$\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i.$$

**Definition 2.8.** The *uncorrected sample variance* is the statistic defined by

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

while the *corrected sample variance* is

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

As is commonly done, we have suppressed the functional notation in the above definitions of these statistic. That is, we have written  $S$  rather than  $S(\mathbf{X})$ . As before, we will denote observed values of statistics with lowercase letters:  $\bar{x}$  and  $s_n^2$  will denote observed values of  $\bar{X}_n$  and  $S_n$ , respectively.

**Definition 2.9.** The *sample moments of order  $r$*  is the statistic defined by

$$M_{r,n} := \frac{1}{n} \sum_{i=1}^n X_i^r,$$

while the *sample centred moments of order  $r$*  is

$$\bar{M}_{r,n} := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^r.$$

Sample values such as the *smallest*, *largest* or *middle* observation from a random sample can provide additional summary information. These are examples of *order statistics*.

**Definition 2.10.** The *order statistics* of a random sample  $\mathbf{X}$  are the sample values placed in ascending order, and they are denoted by  $X_{(1)}, \dots, X_{(n)}$ .

*Remark 2.11.* The order statistics are random variables that satisfy  $X_{(1)} \leq \dots \leq X_{(n)}$ , in particular

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

**Definition 2.12.** The *sample median*  $M$  is a number such that approximately one-half of the observations are less than  $M$  and one-half are greater. In terms of order statistics,  $M$  is defined by

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

*Remark 2.13.* The median is a measure of location that might be considered an alternative to the sample mean. One advantage of the sample median over the sample mean is that it is less affected by extreme observations. Envision a particular set of sample values and then consider the effect of increasing the largest observation. The sample median is unchanged by this change of the sample values. The sample mean, on the other hand, increases without bound as the largest observation increases. This insensitivity to extreme observations is sometimes considered an asset of the sample median.

**Theorem 2.14.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $\mathcal{N}(\mu, \sigma^2)$ . Then,

1.  $\bar{X}_n \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$ ;
2.  $\hat{S}_n^2 \sim \frac{\sigma^2}{n-1}R$ , where  $R \sim \chi_{(n-1)}^2$ .

*Proof.* (1) Using the Moment Generating function,

$$\begin{aligned} M_{\bar{X}_n}(t) &= \mathbb{E} \left[ e^{t\bar{X}_n} \right] \\ &= \mathbb{E} \left[ \exp \left\{ t \frac{X_1 + \dots + X_n}{n} \right\} \right] \\ &= \mathbb{E} \left[ \exp \left\{ \frac{t}{n} (X_1 + \dots + X_n) \right\} \right] \\ &= \mathbb{E} \left[ \exp \left\{ \frac{t}{n} X_1 \right\} \dots \exp \left\{ \frac{t}{n} X_n \right\} \right] \\ &= \left( M_X \left( \frac{t}{n} \right) \right)^n, \end{aligned}$$

because, since  $X_1, \dots, X_n$  are identically distributed,  $M_{X_i}(t)$  is the same function for each  $i$ . Now, the moment generating function for a Gaussian distribution is given by

$$M_X(t) = \exp \left\{ \mu t + \frac{1}{2} \sigma^2 t^2 \right\},$$

therefore our case simply becomes

$$\begin{aligned} M_{\bar{X}_n}(t) &= \left( M_X \left( \frac{t}{n} \right) \right)^n \\ &= \left( \exp \left\{ \mu \frac{t}{n} + \frac{1}{2} \sigma^2 \left( \frac{t}{n} \right)^2 \right\} \right)^n \\ &= \exp \left\{ n \left( \mu \frac{t}{n} + \frac{1}{2} \sigma^2 \left( \frac{t}{n} \right)^2 \right) \right\} \\ &= \exp \left\{ \mu t + \frac{1}{2} \left( \frac{\sigma^2}{n} \right) t^2 \right\}. \end{aligned}$$

We conclude that, since the moment generating function uniquely determines the distribution,

18 2 Week II

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

(2) Consider the following statistic

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

and re-arrange the components,

$$\begin{aligned} W &= \sum_{i=1}^n \left(\frac{(X_i - \bar{X}_n) + (\bar{X}_n - \mu)}{\sigma}\right)^2 \\ &= \sum_{i=1}^n \left(\frac{X_i - \bar{X}_n}{\sigma}\right)^2 + \sum_{i=1}^n \left(\frac{\bar{X}_n - \mu}{\sigma}\right)^2 + 2\left(\frac{\bar{X}_n - \mu}{\sigma}\right) \sum_{i=1}^n (X_i - \bar{X}_n) \\ &= \sum_{i=1}^n \frac{(X_i - \bar{X}_n)^2}{\sigma^2} + \frac{n(\bar{X}_n - \mu)^2}{\sigma^2}, \end{aligned} \quad (2.3)$$

since  $\sum_{i=1}^n (X_i - \bar{X}_n) = 0$ . Now, consider the definition of the sample variance and multiply both sides by  $n - 1$ :

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

and replace it in (2.3). So far, we have that

$$W = \underbrace{\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma}\right)^2}_A = \underbrace{\frac{(n - 1)\hat{S}_n^2}{\sigma^2}}_B + \underbrace{\frac{n(\bar{X}_n - \mu)^2}{\sigma^2}}_C,$$

where

- $A \sim \chi_{(n)}^2$ . Indeed, observe that it is a sum of  $n$  independent  $\chi_{(1)}^2$  random variables: that is because we have assumed that  $X_1, \dots, X_n$  are observations of a random sample of size  $n$  from  $\mathcal{N}(\mu, \sigma^2)$ , therefore

$$\frac{X_i - \mu}{\sigma}$$

follows a standard normal distribution with one degree of freedom.

- $C \sim \chi_{(1)}^2$ . As we showed in (1), the sample mean is normally distributed with mean  $\mu$  and variance  $\frac{\sigma^2}{n}$ , therefore

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

So, if we square  $Z$ , we get a  $\chi_{(1)}^2$ .

In order to conclude the proof, we use a very important result: given three random variables  $Q \sim \chi^2_{(g)}$ ,  $Q_1 \sim \chi^2_{(g_1)}$  and  $Q_2$  such that we can write

$$Q = Q_1 + Q_2,$$

then

$$Q_2 \sim \chi^2_{(g_2)}, \quad \text{with } g_2 = g - g_1.$$

This result implies that  $B \sim \chi^2_{(n-1)}$ , meaning

$$\frac{(n-1)\hat{S}_n^2}{\sigma^2} \sim \chi^2_{(n-1)} \iff \hat{S}_n^2 \sim \frac{\sigma^2}{n-1}R,$$

where  $R \sim \chi^2_{(n-1)}$ .

□

### 2.3 04/10/18

*Remark 2.15.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \mathcal{N}(\mu, \sigma^2)$ , we know that the statistic

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \quad (2.4)$$

is distributed as a  $\mathcal{N}(0, 1)$  random variable. If we knew the value of  $\sigma$  and we measured  $\bar{X}_n$ , then we could use (2.4) as a basis for inference about  $\mu$ , since it would be the only unknown quantity. Most of the time, however,  $\sigma$  is unknown. In this case, we could consider the following statistic,

$$T_n := \frac{\bar{X}_n - \mu}{\hat{S}_n/\sqrt{n}} = \frac{\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}}{\sqrt{\frac{\hat{S}_n^2}{\sigma^2}}} = \frac{Z}{\sqrt{\frac{R}{n-1}}}, \quad (2.5)$$

where  $R \sim \chi_{(n-1)}^2$ . Now, if we do not assume independence between  $Z$  and  $R$ , this is how far we get. However, it can be proved that, given a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \mathcal{N}(\mu, \sigma^2)$ ,

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

are independent. In this case, (2.5)  $\sim t_{(n-1)}$ . Notice that also the vice versa holds.

#### 2.3.1 Concentrations of Measures

Consider a non-negative random variable  $X$  and let  $t > 0$ . Then, Markov's inequality (1.1) tells us that

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[X]}{t}.$$

Now, consider a strictly monotonically increasing function  $\phi$ , with non-negative values: then, we can write

$$\mathbb{P}(X \geq t) = \mathbb{P}(\phi(X) \geq \phi(t)) \leq \frac{\mathbb{E}[\phi(X)]}{\phi(t)}.$$

For instance, if  $\phi(x) = x^2$ , we obtain Chebyshev's inequality. More in general, we can take  $\phi(x) = x^q$  for  $x \geq 0$  and  $q > 0$ , so that we have

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq t) \leq \frac{\mathbb{E}[|X - \mathbb{E}[X]|^q]}{t^q}.$$

Notice that this is a functional in  $q$ ; also, it is possible to choose  $q$  in such a way to optimise the upper bound.

*Remark 2.16.* A related idea is at the basis of *Chernoff's bounding method*. Consider a function  $\phi(x) = e^{sx}$ , where  $s > 0$ . For every random variable  $X$  and any  $t \in \mathbb{R}$ , we have

$$\mathbb{P}(X \geq t) = \mathbb{P}(e^{sX} \geq e^{st}) \leq \frac{\mathbb{E}[e^{sX}]}{e^{st}}. \quad (2.6)$$

Our goal is to find the value of  $s > 0$  for which that upper bound is minimised. Observe that it can be proven that the upper bound given by  $q$  is better than the one obtained minimising  $s$  in here.

**Theorem 2.17.** Let  $t \geq 0$ , then

$$\mathbb{P}(X - \mathbb{E}[X] \geq t) \leq \frac{\text{Var}(X)}{\text{Var}(X) + t^2}. \quad (2.7)$$

*Proof.* For simplicity's sake, assume  $\mathbb{E}[X] = 0$ : this is not restrictive, since otherwise it is enough to rescale the final result. Now, for every  $t$  we could write

$$t = \mathbb{E}[t - X] \leq \mathbb{E}[(t - X)\mathbb{1}\{X < t\}].$$

Then, thanks to the Cauchy-Schwarz inequality (1.2),

$$\begin{aligned} t^2 &\leq \mathbb{E}[(t - X)^2] \mathbb{E}[(\mathbb{1}\{X < t\})^2] \\ &= \mathbb{E}[(t - X)^2] \mathbb{P}(X < t) \\ &= (\text{Var}(X) + t^2) \mathbb{P}(X < t), \end{aligned}$$

meaning that

$$\mathbb{P}(X < t) \geq \frac{t^2}{\text{Var}(X) + t^2}.$$

In order to conclude, we simply take the complementary:

$$\mathbb{P}(X \geq t) = 1 - \mathbb{P}(X < t).$$

□

**Theorem 2.18.** Let  $f$  and  $g$  be non-decreasing real-valued functions defined over the real line. If  $X$  is a real-valued random variable, then

$$\mathbb{E}[f(X)g(X)] \geq \mathbb{E}[f(X)] \mathbb{E}[g(X)].$$

If, however,  $f$  is non-increasing and  $g$  is non-decreasing, then

$$\mathbb{E}[f(X)g(X)] \leq \mathbb{E}[f(X)] \mathbb{E}[g(X)].$$

**Theorem 2.19.** Consider two non-decreasing functions  $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$ . Let  $X_1, \dots, X_n$  be independent real-valued random variables and define the random vector

$$\mathbf{X} = (X_1, \dots, X_n)$$

taking values in  $\mathbb{R}^n$ . Then,

$$\mathbb{E}[f(\mathbf{X})g(\mathbf{X})] \geq \mathbb{E}[f(\mathbf{X})] \mathbb{E}[g(\mathbf{X})].$$

If  $f$  is non-increasing and  $g$  is non-decreasing, then

$$\mathbb{E}[f(\mathbf{X})g(\mathbf{X})] \leq \mathbb{E}[f(\mathbf{X})] \mathbb{E}[g(\mathbf{X})].$$

### 2.3.2 Concentrations for Sums of Independent Random Variables

We would like to bound

$$\mathbb{P}(S_n - \mathbb{E}[S_n] \geq t),$$

where  $S_n = \sum_{i=1}^n X_i$  and  $X_1, \dots, X_n$  are independent real-valued random variables.

An application of Chebyshev's inequality is

$$\mathbb{P}(|S_n - \mathbb{E}[S_n]| \geq t) \leq \frac{\text{Var}(S_n)}{t^2} = \sum_{i=1}^n \frac{\text{Var}(X_i)}{t^2}.$$

If we set  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \text{Var}(X_i)$  and take  $t = n\varepsilon$ , then

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}[X_i])\right| \geq \varepsilon\right) \leq \frac{\sigma^2}{n\varepsilon^2}. \quad (2.8)$$

However, (2.8) is unsatisfactory as  $n \rightarrow \infty$  (even if it works for small values of  $n$ ). To understand why, apply the Central Limit Theorem:

$$\mathbb{P}\left(\sqrt{\frac{n}{\sigma^2}} \left|\frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}[X_i])\right| \geq y\right) \xrightarrow{n \rightarrow \infty} 1 - \phi(y) \leq \frac{1}{\sqrt{2\pi}} \frac{e^{-\frac{y^2}{2}}}{y}.$$

Indeed,

$$\begin{aligned} 1 - \phi(y) &= \int_y^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \quad x = z + y, \\ &= \int_0^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{(z+y)^2}{2}} dz \\ &= \int_0^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2} - zy - \frac{y^2}{2}} dz \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} \int_0^\infty e^{-\frac{z^2}{2} - zy} dz \\ &\leq \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} \left(-\frac{1}{y}\right) \int_0^\infty ye^{-zy} dz \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} \left[-\frac{e^{-zy}}{y}\right]_{z=0}^{z=\infty} \\ &= \frac{e^{-\frac{y^2}{2}}}{y\sqrt{2\pi}}. \end{aligned}$$

Therefore, since it does not depend on  $n$  anymore,

$$\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n(X_i - \mathbb{E}[X_i]) \geq \varepsilon\right) \approx \exp\left\{-\frac{n\varepsilon^2}{2\sigma^2}\right\}, \quad n \rightarrow \infty$$

with  $\varepsilon := \sqrt{n/\sigma^2}y$ . It is useful to compare the quantities

$$\frac{\sigma^2}{n\varepsilon^2} \quad \text{and} \quad \exp\left\{-\frac{n\varepsilon^2}{2\sigma^2}\right\} :$$

the Chebyshev's inequality applied to the sum of random variables does not get the tails' behaviour when  $n \rightarrow \infty$ . We would like to obtain a proper bound for the probability of the form provided by the Central Limit Theorem.



## 3

## Week III

## 3.1 09/10/18

## 3.1.1 Concentration Inequalities

## Hoeffding's Inequalities

The Chernoff's bounding is convenient for bounding the tail probabilities of sums of independent random variables: recalling (2.6),

$$\begin{aligned}\mathbb{P}(S_n - \mathbb{E}[S_n] \geq t) &\leq e^{-st} \mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right\} \right] \\ &= e^{-st} \prod_{i=1}^n \mathbb{E} \left[ e^{s(X_i - \mathbb{E}[X_i])} \right].\end{aligned}\quad (3.1)$$

Therefore, the problem of finding bounds for the tail probabilities reduces to the problem of finding upper bounds for moment generating function (3.1) of  $X_i - \mathbb{E}[X_i]$ .

**Lemma 3.1.** *Let  $X$  be a random variable with  $\mathbb{E}[X] = 0$  and such that  $0 \leq X \leq b$  almost everywhere. Then, for  $s > 0$ ,*

$$\mathbb{E} [e^{sX}] \leq \exp \left\{ \frac{s^2(b-a)^2}{8} \right\}.$$

*Proof.* By the convexity of the exponential function,

$$e^{sX} \leq \frac{X-a}{b-a} e^{sb} + \frac{b-X}{b-a} e^{sa}, \quad a \leq X \leq b.$$

Using the null-mean hypothesis and defining  $p := -\frac{a}{b-a}$ ,

$$\begin{aligned}\mathbb{E} [e^{sX}] &\leq \frac{b - \mathbb{E}[X]}{b-a} e^{sa} + \frac{-a + \mathbb{E}[X]}{b-a} e^{sb} \\ &= (1-p)e^{sa} + pe^{sb} \\ &= e^{sa}(1-p + pe^{s(b-a)})\end{aligned}$$

26 3 Week III

$$\begin{aligned}
 &= (1 - p + pe^{s(b-a)})e^{s(-p)(b-a)} \\
 &=: e^{\phi(u)}.
 \end{aligned}$$

Let  $u = s(b - a)$  and define  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  by taking the log of the above inequality, meaning

$$\phi(u) = -pu + \log(1 - p + pe^u).$$

Now,

$$\phi'(u) = -p + \frac{pe^u}{-p + 1 + pe^u},$$

therefore  $\phi(0) = \phi'(0) = 0$ . Moreover,

$$\begin{aligned}
 \phi''(u) &= \frac{pe^u(1 - p + pe^u) - pe^u pe^u}{(1 - p + pe^u)^2} \\
 &= \frac{pe^u(1 - p)}{(1 - p + pe^u)^2} \\
 &= \frac{pe^u}{1 - p + pe^u} \left(1 - \frac{pe^u}{1 - p + pe^u}\right) \\
 &= t(1 - t) \\
 &\leq \frac{1}{4}.
 \end{aligned}$$

By Taylor's Theorem,

$$\phi(u) = \phi(0) + u\phi'(0) + \frac{1}{2}u^2\phi''(\vartheta) \leq \frac{s^2(b-a)^2}{8}, \quad \vartheta \in [0, u],$$

which is a bound for the moment generating function of  $X$ . □

We could use Lemma [3.1](#) in [\(3.1\)](#):

$$\begin{aligned}
 \mathbb{P}(S_n - \mathbb{E}[S_n] \geq t) &\leq e^{-st} \prod_{i=1}^n \exp\left\{\frac{s^2(b_i - a_i)^2}{8}\right\} \\
 &= e^{-st} \exp\left\{\frac{s^2}{8} \sum_{i=1}^n (b_i - a_i)^2\right\} \\
 &= \exp\left\{-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right\},
 \end{aligned}$$

where

$$s = \frac{4t}{\sum_{i=1}^n (b_i - a_i)^2}.$$

**Theorem 3.2.** Let  $(X_1, \dots, X_n)$  be independent random variables such that

$$X_i \in [a_i, b_i].$$

Then, for every  $t > 0$ ,

$$\mathbb{P}(S_n - \mathbb{E}[S_n] \geq t) \leq \exp \left\{ -\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right\}$$

and

$$\mathbb{P}(S_n - \mathbb{E}[S_n] \leq -t) \leq \exp \left\{ -\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right\}.$$

These are also known as **Hoeffding inequalities**. Notice that we lost the dependence from the variance.

These inequalities have the same form as the one based on the Central Limit Theorem, except that the average variance  $\sigma^2$  is replaced by the upper bound

$$\frac{1}{4} \sum_{i=1}^n (b_i - a_i)^2.$$

### Bernstein's Inequality

Assume  $\mathbb{E}[X_i] = 0$ : we will try to bound

$$\mathbb{E}[e^{sX_i}].$$

Now, set  $\sigma_i^2 = \mathbb{E}[X_i^2]$  and define

$$F_i := \sum_{r=2}^{\infty} \frac{s^{r-2} \mathbb{E}[X_i^r]}{r! \sigma_i^2}.$$

Since

$$e^{sx} = 1 + sx + \sum_{r=2}^{\infty} \frac{s^r x^r}{r!},$$

then

$$\mathbb{E}[e^{sX_i}] = 1 + s\mathbb{E}[X_i] + \sum_{r=2}^{\infty} \frac{s^r \mathbb{E}[X_i^r]}{r!} = 1 + s^2 \sigma_i^2 F_i \leq e^{s^2 \sigma_i^2 F_i}.$$

Assuming that the  $X_i$ 's are bounded, meaning  $|X_i| \leq c$ , for every index  $r$

$$\mathbb{E}[X_i^r] \leq c^{r-2} \sigma_i^2.$$

Indeed,

$$\mathbb{E}[X_i^{r-2} X_i^2] \leq \mathbb{E}[X_i^{r-2}] \mathbb{E}[X_i^2].$$

Thus,

$$F_i \leq \sum_{r=2}^{\infty} \frac{s^{r-2} c^{r-2} \sigma_i^2}{r! \sigma_i^2} = \frac{1}{(sc)^2} \sum_{r=2}^{\infty} \frac{(sc)^r}{r!} = \frac{e^{sc} - 1 - sc}{(sc)^2}$$

and we conclude

$$\mathbb{E}[e^{sX_i}] \leq \exp \left\{ s^2 \sigma_i^2 \frac{e^{sc} - 1 - sc}{s^2 c^2} \right\}.$$

**Theorem 3.3.** Let  $(X_1, \dots, X_n)$  be independent real-valued random variables, with  $\mathbb{E}[X_i] = 0$  and  $|X_i| \leq c$ . Let

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n \text{Var}(X_i).$$

For  $t > 0$ ,

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left\{-\frac{n\sigma^2 h\left(\frac{ct}{n\sigma^2}\right)}{c^2}\right\},$$

where  $h(u) = (1+u)\log(1+u) - u$ , for  $u \geq 0$ .

**Corollary 3.4.** Since

$$h(u) \geq \frac{u^2}{2 + \frac{2u}{3}},$$

we have that

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n X_i > \varepsilon\right) \leq \exp\left\{-\frac{n\varepsilon^2}{2\sigma^2 + \frac{2c\varepsilon}{3}}\right\}.$$

Here we recover the dependence from the variance.

### Efron-Stein's Inequality

Let  $A$  be some set and let  $g : A^n \rightarrow \mathbb{R}$  be a measurable function. We consider the problem of deriving inequalities, for the variance of  $g(X_1, \dots, X_n) = Z$ , where  $X_1, \dots, X_n$  are independent random variables in  $A$ .

**Lemma 3.5.** Let  $\mathbb{E}_i[Z]$  be the expected value of  $Z$  with respect to  $X_i$ , that is

$$\mathbb{E}_i[Z] := \mathbb{E}[Z | X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n].$$

Then,

$$\text{Var}(Z) \leq \sum_{i=1}^n \mathbb{E}[(Z - \mathbb{E}_i[Z])^2].$$

*Proof.* Recall that if  $X$  and  $Y$  are arbitrary bounded random variables, then

$$\mathbb{E}[\mathbb{E}[XY | Y]] = \mathbb{E}[Y\mathbb{E}[X | Y]].$$

For simplicity's sake, we will call  $Z - \mathbb{E}[Z] = V$ . Define

$$V_i = \mathbb{E}[Z | X_1, \dots, X_i] - \mathbb{E}[Z | X_1, \dots, X_{i-1}], \quad i = 1, \dots, n.$$

Now, by definition,  $V = \sum_{i=1}^n V_i$ , therefore

$$\text{Var}(Z) = \mathbb{E}\left[\left(\sum_{i=1}^n V_i\right)^2\right],$$

since

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

Therefore,

$$\text{Var}(Z) = \mathbb{E}\left[\sum_{i=1}^n V_i^2\right] + 2\mathbb{E}\left[\sum_{i>j} V_i V_j\right] = \mathbb{E}\left[\sum_{i=1}^n V_i^2\right]$$

because, for every  $i > j$ ,

$$\mathbb{E}[V_i V_j] = \mathbb{E}[\mathbb{E}[V_i V_j | X_1, \dots, X_j]] = \mathbb{E}[V_j \mathbb{E}[V_i | X_1, \dots, X_j]] = 0,$$

since  $\mathbb{E}[V_i | X_1, \dots, X_j] = 0$ ; indeed,

$$\begin{aligned} \mathbb{E}[V_i | X_1, \dots, X_j] &= \mathbb{E}[\mathbb{E}[Z | X_1, \dots, X_i] - \mathbb{E}[Z | X_1, \dots, X_{i-1}] | X_1, \dots, X_j] \\ &= \mathbb{E}[\mathbb{E}[Z | X_1, \dots, X_i] | X_1, \dots, X_j] - \mathbb{E}[\mathbb{E}[Z | X_1, \dots, X_{i-1}] | X_1, \dots, X_j] \\ &= \mathbb{E}[\mathbb{E}[Z | X_1, \dots, X_j]] - \mathbb{E}[\mathbb{E}[Z | X_1, \dots, X_j]] \\ &= 0, \end{aligned}$$

since  $i > j$  and  $i - 1 \geq j$ .

Now, to bound  $\mathbb{E}[V_i^2]$  we apply Jensen's inequality:

$$\begin{aligned} V_i^2 &= (\mathbb{E}[Z | X_1, \dots, X_i] - \mathbb{E}[Z | X_1, \dots, X_{i-1}])^2 \\ &= (\mathbb{E}[\mathbb{E}[Z | X_1, \dots, X_n] - \mathbb{E}[Z | X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n] | X_1, \dots, X_i])^2 \\ &\leq \mathbb{E}[(\mathbb{E}[Z | X_1, \dots, X_n] - \mathbb{E}[Z | X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n])^2 | X_1, \dots, X_i] \\ &= \mathbb{E}[(Z - \mathbb{E}_i[Z])^2 | X_1, \dots, X_i]. \end{aligned}$$

Finally, taking the expected values on both sides, we conclude.  $\square$

**Theorem 3.6 (Efron-Stein).** Let  $X'_1, \dots, X'_n$  (identically distributed to  $X_1, \dots, X_n$ ) form an independent copy of  $X_1, \dots, X_n$  and define

$$Z'_i = g(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n).$$

Then, we have

$$\text{Var}(Z) \leq \frac{1}{2} \sum_{i=1}^n \mathbb{E}[(Z - Z'_i)^2].$$

Whenever

$$g(X_1, \dots, X_n) = X_1 + \dots + X_n,$$

the Efron-Stein's inequality becomes an equality.

*Proof.* Remember that, given two independent and identically random variables  $X$  and  $Y$ ,

$$\text{Var}(X) = \frac{1}{2} \mathbb{E}[(X - Y)^2].$$

Therefore, given

30 3 Week III

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

and considering  $Z$  and  $Z'_i$  (which are independent and identically distributed),

$$\mathbb{E}_i [(Z - \mathbb{E}_i [Z])^2] = \frac{1}{2} \mathbb{E}_i [(Z - Z'_i)^2],$$

then apply Lemma [3.5](#)

□

*Remark 3.7.* When  $Z = g(X_1, \dots, X_n) = \sum_{i=1}^n X_i$  is a sum of independent random variable, then the inequality is a proper equality and the bound is not improvable.

**Corollary 3.8.** *We have that*

$$\text{Var}(Z) \leq \sum_{i=1}^n \mathbb{E} [(Z - Z_i)^2].$$

*Proof.* Recall that, for every random variable  $X$ ,

$$\text{Var}(X) \leq \mathbb{E} [(X - a)^2], \quad \forall a \in \mathbb{R}.$$

Therefore, for every  $i = 1, \dots, n$  we have

$$\mathbb{E}_i [(Z - \mathbb{E}_i [Z])^2] \leq \mathbb{E}_i [(Z - Z_i)^2],$$

where  $Z_i = g_i(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$  for an arbitrary  $g_i : A^{n-1} \rightarrow \mathbb{R}$ . Taking the expectations and applying Lemma [3.5](#) we get the result.

□

## 3.2 10/10/18

### 3.2.1 Likelihood Function

The Likelihood function contains all the statistical information required to make inference.

**Definition 3.9.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim f_X(x; \vartheta)$ . Then, given an observation  $\mathbf{X} = \mathbf{x}$ , the function of  $\vartheta$  defined by

$$\mathcal{L}(\vartheta, \mathbf{x}) := f_{\mathbf{X}}(\mathbf{x}; \vartheta) = \prod_{i=1}^n f_{X_i}(x_i; \vartheta)$$

is called **Likelihood function**.

This Definition almost seems to be defining the Likelihood function to be the same as the probability density function, or the probability mass function. The only distinction between these two is which variable is considered fixed and which is varying. When we consider the probability density function, or the probability mass function,  $f(\mathbf{x}; \vartheta)$ , we are considering  $\vartheta$  as fixed and  $\mathbf{x}$  as the variable; when we consider the Likelihood function  $\mathcal{L}(\vartheta; \mathbf{x})$ , we are considering  $\mathbf{x}$  to be the observed sample point and  $\vartheta$  to be varying over all possible values. Notice that, in Probability, we are interested in computing something like  $\mathbb{P}(\mathbf{X} \in \{\cdot\})$ , and  $\vartheta$  must be known to achieve the result. In Statistics, on the other hand, observations are given and used to infer information about the parameter  $\vartheta$ .

For mathematical convenience, we shall use the **Log-Likelihood function**,

$$\log \mathcal{L}(\vartheta; \mathbf{x}) := \sum_{i=1}^n \log f_{X_i}(x_i; \vartheta) :$$

indeed, we are generally interested in knowing where the Likelihood function reaches its maximum value. Since the logarithm is a strictly increasing function, so the logarithm of a function achieves its maximum value at the same points as the function itself.

If both  $\mathbf{x}$  and  $\vartheta$  are fixed, then  $\mathcal{L}(\vartheta; \mathbf{x})$  is related to the probability that the particular value that we fixed for  $\vartheta$  has generated  $\mathbf{x}$ .

**Definition 3.10.** Consider the Likelihood function  $\mathcal{L}(\vartheta; \mathbf{x})$  and  $\vartheta_1, \vartheta_2 \in \Theta$ . Suppose we have fixed  $\vartheta_1, \vartheta_2 \in \Theta$  and we compute  $\mathcal{L}(\vartheta_i; \mathbf{x})$  for  $i = 1, 2$ . We say that  $\vartheta_1$  is **more likely** than  $\vartheta_2$  if

$$\mathcal{L}(\vartheta_1, \mathbf{x}) > \mathcal{L}(\vartheta_2; \mathbf{x}).$$

*Remark 3.11.* This means that  $\mathbf{x}$  has been more likely generated under  $\vartheta_1$ , since under  $\vartheta_1$  the probability of generating  $\mathbf{x}$  is higher than under  $\vartheta_2$ . In general, we would like to optimise the Likelihood function with respect to  $\vartheta$ : this will not always be possible through differentiation.

The statistical inference based on Likelihood functions is a consequence of two principles.

**Weak Likelihood Principle**

For a fixed parametric model  $X \sim F_X(x; \vartheta)$ , if  $\mathbf{x}$  and  $\mathbf{y}$  are two observed sample points such that

$$\mathcal{L}(\vartheta; \mathbf{x}) \propto \mathcal{L}(\vartheta; \mathbf{y}),$$

then the two samples provide the same inference on  $\vartheta$ .

**Strong Likelihood Principle**

Let  $\mathbf{X} = \mathbf{x}$  be an observed sample point under the model  $X \sim F_X(x; \vartheta)$  with Likelihood function  $\mathcal{L}_F(\mathbf{x}; \vartheta)$  and let  $\mathbf{Y} = \mathbf{y}$  be an observed sample point under the model  $Y \sim G_Y(y; \vartheta)$  with Likelihood function  $\mathcal{L}_G(\mathbf{y}; \vartheta)$ . If  $\mathcal{L}_F(\mathbf{x}; \vartheta) \propto \mathcal{L}_G(\mathbf{y}; \vartheta)$ , then the two samples provide the same inference on  $\vartheta$ .

*Example 3.12.* Suppose we have a coin and we want to determine the “success” probability. To do so, we flip the coin three times and we try to infer some information about its probability, according to the outcomes. This means that we are working with a random sample  $(X_1, X_2, X_3)$ . Compare the following cases:

$$X \sim Ber(\vartheta_1), \quad \vartheta_1 = \frac{1}{2};$$

$$X \sim Ber(\vartheta_2), \quad \vartheta_2 = \frac{1}{3};$$

$$X \sim Ber(\vartheta_3), \quad \vartheta_3 = \frac{1}{4}.$$

The results are the following:

$(x_1, x_2, x_3)$	$\vartheta_1$	$\vartheta_2$	$\vartheta_3$
(0, 0, 0)	1/8	8/27	27/64*
(0, 0, 1)	1/8	4/27*	9/64
(0, 1, 0)	1/8	4/27*	9/64
(1, 0, 0)	1/8	4/27*	9/64
(0, 1, 1)	1/8*	2/27	3/64
(1, 0, 1)	1/8*	2/27	3/67
(1, 1, 0)	1/8*	2/27	3/67
(1, 1, 1)	1/8*	1/27	1/64

According to the Likelihood Principle, we know that our coin’s probability is the one with the highest values.

### 3.3 11/10/18

#### 3.3.1 Likelihood Function Method: Assumptions

**Definition 3.13.** Consider  $X \sim f_X(x; \vartheta)$  and fix  $\vartheta \in \Theta$ , where  $\Theta$  is an open set of parameters. The following conditions are called **conditions of regularity** for  $\{f_X(\cdot, \vartheta) \mid \vartheta \in \Theta\}$ :

1. for every  $\vartheta \in \Theta$  there exist derivatives of the Likelihood function with respect to  $\vartheta$ , at least up to the third order;
2. for every  $\vartheta_0 \in \Theta$  there exist three functions, namely  $g(\cdot)$ ,  $h(\cdot)$  and  $H(\cdot)$ , that are integrable in a neighbourhood of  $\vartheta_0$  and such that

$$\begin{aligned} \left| \frac{d}{d\vartheta} f_X(x; \vartheta) \right| &\leq g(x), \\ \left| \frac{d^2}{d\vartheta^2} f_X(x; \vartheta) \right| &\leq h(x), \\ \left| \frac{d^3}{d\vartheta^3} \log(f_X(x; \vartheta)) \right| &\leq H(x); \end{aligned}$$

3. for every  $\vartheta \in \Theta$ ,

$$0 < \mathbb{E} \left[ \left( \frac{d}{d\vartheta} \log \mathcal{L}(\vartheta; \mathbf{x}) \right)^2 \right] < \infty.$$

Note that  $\log \mathcal{L}(\vartheta; \mathbf{x})$  is a well-defined random variable.

**Definition 3.14.** A statistical model is **identifiable** if, for every  $\vartheta_1, \vartheta_2 \in \Theta$ , there exist at least an event  $E$  such that

$$\mathbb{P}(X \in E; \vartheta_1) \neq \mathbb{P}(X \in E; \vartheta_2).$$

**Definition 3.15.** The random variable

$$V_n(\vartheta) := \log \mathcal{L}(\vartheta, \mathbf{x}) = \sum_{i=1}^n \log f_{X_i}(x_i; \vartheta)$$

is called **Log-Likelihood random variable**. Moreover, we call

$$V'_n(\vartheta) := \frac{d}{d\vartheta} V_n(\vartheta) = \frac{\mathcal{L}'(\vartheta; \mathbf{x})}{\mathcal{L}(\vartheta; \mathbf{x})}$$

the **score function**.

**Proposition 3.16.** The following properties are satisfied:

1.  $\mathbb{E}[V'_n(\vartheta)] = 0$ ;
2.  $\text{Var}(V'_n(\vartheta)) = \mathbb{E}[(V'_n(\vartheta))^2] = -\mathbb{E}[V''_n(\vartheta)]$ .

34 3 Week III

*Proof.* (1) First,

$$\begin{aligned}
\mathbb{E}[V_n'(\vartheta)] &= \int_{\mathbb{R}^n} V_n'(\vartheta) f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} \\
&= \int_{\mathbb{R}^n} \frac{f'_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}; \vartheta)} f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} \\
&= \int_{\mathbb{R}^n} \frac{d}{d\vartheta} f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} \\
&= \frac{d}{d\vartheta} \int_{\mathbb{R}^n} f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} \\
&= 0.
\end{aligned}$$

(2) Second,

$$\begin{aligned}
V_n''(\vartheta) &= \frac{d^2}{d\vartheta^2} \log \mathcal{L}(\vartheta; \mathbf{x}) \\
&= \frac{d}{d\vartheta} \frac{f'_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}; \vartheta)} \\
&= \frac{f''_{\mathbf{X}}(\mathbf{x}; \vartheta) f_{\mathbf{X}}(\mathbf{x}; \vartheta) - f'_{\mathbf{X}}(\mathbf{x}; \vartheta) f'_{\mathbf{X}}(\mathbf{x}; \vartheta)}{(f_{\mathbf{X}}(\mathbf{x}; \vartheta))^2} \\
&= \frac{f''_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}; \vartheta)} - \left( \frac{f'_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}; \vartheta)} \right)^2 \\
&= \frac{f''_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}; \vartheta)} - (V_n'(\vartheta))^2,
\end{aligned}$$

therefore

$$\begin{aligned}
\mathbb{E}[V_n''(\vartheta)] &= \int_{\mathbb{R}^n} \frac{f''_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}; \vartheta)} f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} - \mathbb{E}[(V_n'(\vartheta))^2] \\
&= \int_{\mathbb{R}^n} \frac{d^2}{d\vartheta^2} f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} - \mathbb{E}[(V_n'(\vartheta))^2] \\
&= \frac{d^2}{d\vartheta^2} \int_{\mathbb{R}^n} f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} - \mathbb{E}[(V_n'(\vartheta))^2] \\
&= -\mathbb{E}[(V_n'(\vartheta))^2].
\end{aligned}$$

□

**Definition 3.17.**  $I_n(\vartheta) := -\mathbb{E}[V_n''(\vartheta)]$  is called **Fisher information**.*Remark 3.18.* If  $\mathbf{X} = (X_1, \dots, X_n)$  is a random sample from  $X \sim f_X(x, \vartheta)$  regular, then

$$I_n(\vartheta) = nI_1(\vartheta).$$

*Example 3.19.* Consider a random sample  $(X_1, \dots, X_n)$  from  $X \sim \text{Ber}(\vartheta)$ . Then,

$$\begin{aligned}
\mathcal{L}(\vartheta, \mathbf{x}) &= \prod_{i=1}^n f_{X_i}(x_i; \vartheta) \\
&= \vartheta^{\sum_{i=1}^n X_i} (1 - \vartheta)^{n - \sum_{i=1}^n X_i};
\end{aligned}$$

$$\begin{aligned}
V_n(\vartheta) &= \log \vartheta \sum_{i=1}^n X_i + \log(1 - \vartheta) \left( n - \sum_{i=1}^n X_i \right); \\
V_n'(\vartheta) &= \frac{\sum_{i=1}^n X_i}{\vartheta} - \frac{n - \sum_{i=1}^n X_i}{1 - \vartheta}; \\
V_n''(\vartheta) &= \frac{-\sum_{i=1}^n X_i}{\vartheta^2} - \frac{n - \sum_{i=1}^n X_i}{(1 - \vartheta)^2}; \\
\mathbb{E}[V_n'(\vartheta)] &= \frac{1}{\vartheta} \sum_{i=1}^n \mathbb{E}[X_i] - \frac{1}{1 - \vartheta} \left( n - \sum_{i=1}^n \mathbb{E}[X_i] \right) \\
&= \frac{n\vartheta}{\vartheta} - \frac{n - n\vartheta}{1 - \vartheta} \\
&= 0; \\
\mathbb{E}[V_n''(\vartheta)] &= \frac{1}{\vartheta^2} \sum_{i=1}^n \mathbb{E}[X_i] - \frac{1}{(1 - \vartheta)^2} \left( n - \sum_{i=1}^n \mathbb{E}[X_i] \right) \\
&= \frac{n\vartheta}{\vartheta^2} - \frac{n - n\vartheta}{(1 - \vartheta)^2} \\
&= -I_n(\vartheta).
\end{aligned}$$

*Remark 3.20.* If  $\Theta$  is  $m$ -dimensional, then  $\text{Im } \boldsymbol{\vartheta}$  is an  $m \times m$  matrix.

### 3.3.2 The Exponential Family

**Definition 3.21.** A single-parameter **exponential family** (EF) is a set of probability distributions whose probability density function (or probability mass function) can be expressed in the following form: if  $X \sim EF(\vartheta)$

$$f_X(x; \vartheta) := \exp \{ Q(\vartheta)A(x) + C(x) - K(\vartheta) \}. \quad (3.2)$$

On the other hand, if we consider a sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim f_{\mathbf{X}}(\mathbf{x}; \vartheta)$ ,

$$f_{\mathbf{X}}(\mathbf{x}; \vartheta) := \exp \left\{ Q(\vartheta) \sum_{i=1}^n A(x_i) + \sum_{i=1}^n C(x_i) - K(\vartheta) \right\}.$$

*Example 3.22.* Consider  $X \sim \text{Ber}(\vartheta)$ ,  $X \in \{0, 1\}$  and  $\Theta = (0, 1)$ . We know that

$$p_X(x) = \vartheta^x (1 - \vartheta)^{1-x} \mathbf{1}_{\{0,1\}}(x).$$

Now,

$$\begin{aligned}
\exp \{ \log(\vartheta^x (1 - \vartheta)^{1-x}) \} &= \exp \{ x \log \vartheta + (1 - x) \log(1 - \vartheta) \} \\
&= \exp \{ x \log \vartheta + \log(1 - \vartheta) - x \log(1 - \vartheta) \} \\
&= \exp \left\{ x \log \left( \frac{\vartheta}{1 - \vartheta} \right) + \log(1 - \vartheta) \right\},
\end{aligned}$$

where

$$Q(\vartheta) = \log \left( \frac{\vartheta}{1 - \vartheta} \right), \quad C(x) = 0, \quad A(x) = x, \quad K(\vartheta) = -\log(1 - \vartheta).$$

**Properties of the Exponential Family****Proposition 3.23.** *Let  $X \sim EF(\vartheta)$ . Then,*

1.  $\mathbb{E}[A(x)] = \frac{K'(\vartheta)}{Q'(\vartheta)}$ ;
2.  $\text{Var}(A(x)) = \frac{K''(\vartheta)}{[Q'(\vartheta)]^2} - \mathbb{E}[A(x)] \frac{Q''(\vartheta)}{[Q'(\vartheta)]^2}$ .

*Proof.* 1. Differentiating (3.2) we obtain

$$\frac{d}{d\vartheta} f_X(x; \vartheta) = [A(x)Q'(\vartheta) - K'(\vartheta)]f_X(x; \vartheta),$$

therefore, by taking the integral on both sides of the equality,

$$\begin{aligned} \int_{\mathbb{R}} \frac{d}{d\vartheta} f_X(x; \vartheta) dx &= \int_{\mathbb{R}} [A(x)Q'(\vartheta) - K'(\vartheta)]f_X(x; \vartheta) dx \\ &= Q'(\vartheta) \int_{\mathbb{R}} A(x)f_X(x; \vartheta) dx - K'(\vartheta) \int_{\mathbb{R}} f_X(x; \vartheta) dx \end{aligned}$$

we get

$$0 = Q'(\vartheta)\mathbb{E}[A(x)] - K'(\vartheta).$$

We conclude that

$$\mathbb{E}[A(x)] = \frac{K'(\vartheta)}{Q'(\vartheta)}.$$

2. For the variance, we differentiate (3.2) twice:

$$\frac{d^2}{d\vartheta^2} f_X(x; \vartheta) = [A(x)Q''(\vartheta) - K''(\vartheta)]f_X(x; \vartheta) + [A(x)Q'(\vartheta) - K'(\vartheta)]^2 f_X(x; \vartheta),$$

therefore, by taking the integral on both sides of the equality,

$$\begin{aligned} \int_{\mathbb{R}} \frac{d^2}{d\vartheta^2} f_X(x; \vartheta) dx &= \int_{\mathbb{R}} [A(x)Q''(\vartheta) - K''(\vartheta)]f_X(x; \vartheta) dx + \\ &\quad + \int_{\mathbb{R}} [A(x)Q'(\vartheta) - K'(\vartheta)]^2 f_X(x; \vartheta) dx \end{aligned}$$

we get

$$\begin{aligned} 0 &= Q''(\vartheta)\mathbb{E}[A(x)] - K''(\vartheta) + [Q'(\vartheta)]^2 \int_{\mathbb{R}} \left[ A(x) - \frac{K'(\vartheta)}{Q'(\vartheta)} \right]^2 f_X(x; \vartheta) dx \\ &= Q''(\vartheta) \frac{K'(\vartheta)}{Q'(\vartheta)} - K''(\vartheta) + [Q'(\vartheta)]^2 \text{Var}(A(x)). \end{aligned}$$

We conclude that

$$\begin{aligned} \text{Var}(A(x)) &= \frac{K''(\vartheta) - Q''(\vartheta) \frac{K'(\vartheta)}{Q'(\vartheta)}}{[Q'(\vartheta)]^2} \\ &= \frac{K''(\vartheta)}{[Q'(\vartheta)]^2} - \mathbb{E}[A(x)] \frac{Q''(\vartheta)}{[Q'(\vartheta)]^2}. \end{aligned}$$

□

*Remark 3.24.* If  $Q(\vartheta) = \vartheta$ , then

$$\mathbb{E}[A(x)] = K'(\vartheta) \quad \text{and} \quad \text{Var}(A(x)) = K''(\vartheta).$$

### 3.3.3 The *Natural Exponential Family*

A *natural exponential family* (NEF) is a class of probability distributions that is a special case of an exponential family. Every distribution possessing a moment-generating function is a member of a natural exponential family, and the use of such distributions simplifies the theory and computation of generalized linear models.

**Definition 3.25.** The *natural exponential family* is a subset of the exponential family. It is an exponential family in which  $A(x) = x$  and  $Q(\vartheta) = \eta$ :

$$f_X(x; \eta) = \exp \{ \eta x + C(x) - K(\eta) \}.$$



---

## Week IV

### 4.1 16/10/18

An experimenter uses the information in a sample  $X_1, \dots, X_n$  to make inferences about an unknown parameter  $\vartheta$ . If the sample size  $n$  is large, then the observed sample  $x_1, \dots, x_n$  is a long list of numbers that may be hard to interpret. An experimenter might wish to summarise the information in a sample by determining a few key features of the sample values. This is usually done by computing statistics, functions of the sample. For example, the *sample mean*, the *sample variance*, the *largest observation* and the *smallest observation* are four statistics that might be used to summarise some key features of the sample.

Any statistic  $T(\mathbf{X})$  defines a form of data reduction or data summary. An experimenter who uses only the observed value of the statistic,  $T(\mathbf{X}) = T(\mathbf{x})$ , rather than the entire observed sample,  $\mathbf{X} = \mathbf{x}$ , will treat as equal two samples,  $\mathbf{X} = \mathbf{x}$  and  $\mathbf{Y} = \mathbf{y}$ , that satisfy  $T(\mathbf{x}) = T(\mathbf{y})$  even though the actual sample values may be different in some ways.

Data reduction in terms of a particular statistic can be thought of as a partition of the sample space  $\mathcal{X}$ . Let  $\mathcal{T} = \{t \mid t = T(\mathbf{x}), \mathbf{x} \in \mathcal{X}\}$  be the image of  $\mathcal{X}$  under  $T(\mathbf{x})$ . Then,  $T(\mathbf{x})$  partitions the sample space into sets  $A_t = \{\mathbf{x} \mid T(\mathbf{x}) = t\}$ . The statistic summarises the data in that, rather than reporting the entire sample  $\mathbf{x}$ , it reports only  $T(\mathbf{x}) = t$  or, equivalently,  $\mathbf{x} \in A_t$ .

We study two principles of data reduction. We are interested in methods of data reduction that do not discard important information about the unknown parameter  $\vartheta$  and methods that successfully discard information that is irrelevant as far as gaining knowledge about  $\vartheta$  is concerned. The **Sufficiency Principle** promotes a method of data reduction that does not discard information about  $\vartheta$  while achieving some summarisation of the data. The **Likelihood Principle** describes a function of the parameter, determined by the observed sample, that contains all the information about  $\vartheta$  that is available from the sample.

#### 4.1.1 Sufficiency Principle

**Definition 4.1.** A statistic  $T_n = T(\mathbf{X})$  is a **sufficient statistic** for  $\vartheta$  if the conditional distribution of the sample  $\mathbf{X}$ , given the values of  $T(\mathbf{X})$  does not depend on  $\vartheta$ .

Basically, a sufficient statistic for a parameter  $\vartheta$  is a statistic that, in a certain sense, captures all the information about  $\vartheta$  contained in the sample. Any additional information in the sample, besides the values of the sufficient statistic, does not contain any more information about  $\vartheta$ .

To use this Definition to verify that a statistic  $T_n = T(\mathbf{X})$  is a sufficient statistic for  $\vartheta$ , we must verify that for any fixed values of  $\mathbf{x}$  and  $t$ , the conditional probability

$$\mathbb{P}(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = t)$$

is the same for all values of  $\vartheta$ . Now, this probability is zero for every value of  $\vartheta$  if  $T(\mathbf{x}) \neq t$ . So, we must verify only that

$$\mathbb{P}(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x}))$$

does not depend on  $\vartheta$ . But since  $\{\mathbf{X} = \mathbf{x}\}$  is a subset of  $\{T(\mathbf{X}) = T(\mathbf{x})\}$ ,

$$\begin{aligned} \mathbb{P}(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})) &= \frac{\mathbb{P}(\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = T(\mathbf{x}))}{\mathbb{P}(T(\mathbf{X}) = T(\mathbf{x}))} \\ &= \frac{\mathbb{P}(\mathbf{X} = \mathbf{x})}{\mathbb{P}(T(\mathbf{X}) = T(\mathbf{x}))} \\ &= \frac{h_{(\mathbf{X}, T_n)}(\mathbf{x}, t; \vartheta)}{g_{T_n}(t; \vartheta)} \\ &=: \phi_{\mathbf{X} | T_n=t}(\mathbf{x}; t, \vartheta), \end{aligned} \quad (4.1)$$

where

- $h_{(\mathbf{X}, T_n)}(\mathbf{x}, t; \vartheta) := \mathbb{P}(\mathbf{X} = \mathbf{x})$  is the joint distribution of  $\mathbf{X}$  and  $T_n$ ;
- $g_{T_n}(t; \vartheta) := \mathbb{P}(T(\mathbf{X}) = T(\mathbf{x}))$  is the marginal distribution of  $T_n$ ;
- $\phi_{\mathbf{X} | T_n=t}(\mathbf{x}; t, \vartheta)$  is the conditional distribution of  $\mathbf{X}$  given  $T_n$ .

Thus,  $T_n = T(\mathbf{X})$  is a sufficient statistic for  $\vartheta$  if and only if, for every  $\mathbf{x}$ , (4.1) is a constant as a function of  $\vartheta$ . Please note that if  $\mathbf{X}$  and  $T(\mathbf{X})$  have continuous distributions, then the above conditional probabilities cannot be interpreted as we did. However, it is still appropriate to use the above criterion to determine if  $T(\mathbf{X})$  is a sufficient statistic for  $\vartheta$ , simply considering the distributions.

*Example 4.2.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \text{Ber}(\vartheta)$ . Then,  $T(\mathbf{X}) = \sum_{i=1}^n X_i$  is a sufficient statistic for  $\vartheta$ .

*Proof.* Note that  $T_n = T(\mathbf{X})$  counts the number of  $X_i$ 's that equal 1, so  $T_n = T(\mathbf{X})$  is a Binomial with parameters  $(n, \vartheta)$ . The ratio of the probability mass functions is thus

$$\begin{aligned} \frac{h_{(\mathbf{X}, T_n)}(\mathbf{x}, t; \vartheta)}{g_{T_n}(t; \vartheta)} &= \frac{\prod_{i=1}^n \vartheta^{x_i} (1 - \vartheta)^{1-x_i}}{\binom{n}{t} \vartheta^t (1 - \vartheta)^{n-t}} & t = \sum_{i=1}^n x_i \\ &= \frac{\vartheta^{\sum_{i=1}^n x_i} (1 - \vartheta)^{\sum_{i=1}^n (1-x_i)}}{\binom{n}{t} \vartheta^t (1 - \vartheta)^{n-t}} \\ &= \frac{\vartheta^t (1 - \vartheta)^{n-t}}{\binom{n}{t} \vartheta^t (1 - \vartheta)^{n-t}} \end{aligned}$$

$$= \frac{1}{\binom{n}{t}}.$$

Therefore,  $T_n = T(\mathbf{X})$  is sufficient for  $\vartheta$ .

Similarly, we can approach this problem from another direction: let us call

$$\begin{aligned} I_0 &= \left\{ (x_1, \dots, x_n) \mid x_i \in \{0, 1\} \wedge \sum_{i=1}^n x_i = 0 \right\} \\ I_1 &= \left\{ (x_1, \dots, x_n) \mid x_i \in \{0, 1\} \wedge \sum_{i=1}^n x_i = 1 \right\} \\ &\vdots \\ I_n &= \left\{ (x_1, \dots, x_n) \mid x_i \in \{0, 1\} \wedge \sum_{i=1}^n x_i = n \right\} \end{aligned}$$

such that  $I_i \cap I_j = \emptyset$  and  $\bigcup_{j=1}^n I_j = \{0, 1\}^n$ . This means that the probability is simply

$$\mathbb{P}(\mathbf{X} = \mathbf{x} \mid T_n = t) = \frac{1}{\binom{n}{t}}.$$

□

It may be unwieldy to use the definition of a sufficient statistic to find a sufficient statistic for a particular model. To use the definition, we must guess a statistic  $T_n = T(\mathbf{X})$  to be sufficient, find the probability mass function (or probability density function) of  $T(\mathbf{X})$ , and check that the ratio of the probability mass functions (or probability density functions) does not depend on  $\vartheta$ . The first step requires a good deal of intuition and the second sometimes requires some tedious analysis.

**Theorem 4.3 (Fisher Factorisation).** *If  $h_{(\mathbf{X}, T_n)}(\mathbf{x}, t; \vartheta)$  is the joint probability mass function (or probability density function) of  $\mathbf{X}$  and  $T_n$ , and  $g_{T_n}(t; \vartheta)$  is the probability density function (or probability mass function) of  $T_n$ , then  $T_n = T(\mathbf{X})$  is a sufficient statistic for  $\vartheta$  if and only if, for every  $\mathbf{x} \in \mathcal{X}$ ,  $\phi_{\mathbf{X} \mid T_n=t}(\mathbf{x}; t, \vartheta)$  is constant as a function of  $\vartheta$ .*

## 4.2 17/10/18

### 4.2.1 Savage's Factorisation Theorem

Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $f_{\mathbf{X}}(x; \vartheta)$ . The idea is to determine that a statistic  $T_n$  is sufficient for  $\vartheta$  if and only if there exist two non-negative functions  $g(\cdot)$  and  $h(\cdot)$  such that

$$\mathcal{L}(\vartheta; \mathbf{x}) = g(T(\mathbf{X}); \vartheta)h(\mathbf{x}),$$

where  $g(\cdot)$  is a function of the observable sample via  $T(\cdot)$ , while  $h(\cdot)$  is a function of the observed sample which does not depend on  $\vartheta$ .

**Theorem 4.4 (Savage's Factorisation).** *Consider a random sample  $\mathbf{X}$  from  $X \sim f_{\mathbf{X}}(x; \vartheta)$ . Let  $f_{\mathbf{X}}(\mathbf{x}; \vartheta)$  be the joint probability density function (or probability mass function) of  $\mathbf{X}$ . A statistic  $T(\mathbf{X})$  is sufficient for  $\vartheta$  if and only if there exist two functions, namely  $g(\cdot)$  and  $h(\cdot)$ , such that, for every  $\mathbf{x} \in \mathcal{X}$  and every  $\vartheta \in \Theta$ ,*

$$f_{\mathbf{X}}(\mathbf{x}; \vartheta) = g(T(\mathbf{x}); \vartheta)h(\mathbf{x}). \quad (4.2)$$

*Proof.* We give the proof only for discrete distributions.

Suppose  $T(\mathbf{X})$  is a sufficient statistic. Choose

$$g(t; \vartheta) = \mathbb{P}(T(\mathbf{x}) = t) \quad \text{and} \quad h(\mathbf{x}) = \mathbb{P}(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})).$$

Because  $T(\mathbf{X})$  is sufficient, the conditional probability defining  $h(\mathbf{x})$  does not depend on  $\vartheta$ . Thus, this choice of  $h(\mathbf{x})$  and  $g(t; \vartheta)$  is legitimate, and for this choice we have

$$\begin{aligned} f_{\mathbf{X}}(\mathbf{x}; \vartheta) &= \mathbb{P}(\mathbf{X} = \mathbf{x}) \\ &= \mathbb{P}(\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = T(\mathbf{x})) \\ &= \mathbb{P}(T(\mathbf{X}) = T(\mathbf{x})) \mathbb{P}(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})) \\ &= g(T(\mathbf{x}); \vartheta)h(\mathbf{x}). \end{aligned}$$

Then, the factorisation 4.2 has been exhibited. Notice that

$$\mathbb{P}(T(\mathbf{X}) = T(\mathbf{x})) = g(T(\mathbf{x}); \vartheta),$$

meaning that  $g(T(\mathbf{x}); \vartheta)$  is the probability mass function of  $T(\mathbf{X})$ .

Now, assume that the factorisation 4.2 exists. Let  $q(t; \vartheta)$  be the probability mass function of  $T(\mathbf{X})$ . To show that  $T(\mathbf{X})$  is sufficient, we examine the ratio

$$\frac{f_{\mathbf{X}}(\mathbf{x}; \vartheta)}{q(T(\mathbf{x}); \vartheta)}.$$

Define

$$A_{T(\mathbf{x})} := \{\mathbf{y} | T(\mathbf{y}) = T(\mathbf{x})\};$$

then,

$$\begin{aligned} \frac{f_{\mathbf{X}}(\mathbf{x}; \vartheta)}{q(T(\mathbf{x}); \vartheta)} &= \frac{g(T(\mathbf{x}); \vartheta)h(\mathbf{x})}{q(T(\mathbf{x}); \vartheta)} \\ &= \frac{g(T(\mathbf{x}); \vartheta)h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} g(\mathbf{y}; \vartheta)h(\mathbf{y})} \end{aligned}$$

$$\begin{aligned}
&= \frac{g(T(\mathbf{x}); \vartheta)h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} g(T(\mathbf{x}); \vartheta)h(\mathbf{y})} \\
&= \frac{g(T(\mathbf{x}); \vartheta)h(\mathbf{x})}{g(T(\mathbf{x}); \vartheta) \sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})} \\
&= \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})}
\end{aligned}$$

since  $T$  is a constant on  $A_{T(\mathbf{x})}$ . Therefore, because the ratio does not depend on  $\vartheta$ , we conclude thanks to Theorem 4.3 that  $T(\mathbf{X})$  is a sufficient statistic for  $\vartheta$ .  $\square$

*Remark 4.5.* To use Theorem 4.4 to find a sufficient statistic, we factor the probability density function (or probability mass function) into two parts, with one part not depending on  $\vartheta$ . The part that does not depend on  $\vartheta$  constitutes  $h(\mathbf{x})$ . The other one, the one that does depend on  $\vartheta$ , usually depends on the sample  $\mathbf{x}$  only through some function  $T(\mathbf{x})$ , and this function is a sufficient statistic for  $\vartheta$ .

*Example 4.6.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \text{Ber}(\vartheta)$ . Then,

$$\mathcal{L}(\vartheta; \mathbf{x}) = \prod_{i=1}^n \vartheta^{x_i} (1 - \vartheta)^{1-x_i} = \vartheta^{\sum_{i=1}^n x_i} (1 - \vartheta)^{n - \sum_{i=1}^n x_i}.$$

Here

$$h(T(\mathbf{x})) = 1 \quad \text{and} \quad g(\mathbf{x}) = g\left(\sum_{i=1}^n x_i\right) = \vartheta^{\sum_{i=1}^n x_i} (1 - \vartheta)^{n - \sum_{i=1}^n x_i}.$$

*Example 4.7.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \mathcal{N}(0, 1)$ . Then,

$$\begin{aligned}
\mathcal{L}(\vartheta; \mathbf{x}) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(x_i - \vartheta)^2\right\} \\
&= (2\pi)^{-n/2} \exp\left\{-\frac{1}{2} \sum_{i=1}^n (x_i - \vartheta)^2\right\} \\
&= (2\pi)^{-n/2} \exp\left\{-\frac{1}{2} \sum_{i=1}^n x_i^2 - \frac{n\vartheta^2}{2} + \vartheta \sum_{i=1}^n x_i\right\} \\
&= (2\pi)^{-n/2} \exp\left\{-\frac{1}{2} \sum_{i=1}^n x_i^2\right\} \exp\left\{-\frac{n\vartheta^2}{2} + \vartheta \sum_{i=1}^n x_i\right\} \\
&= h(\mathbf{x})g(T(\mathbf{x}); \vartheta).
\end{aligned}$$

Here

$$\begin{aligned}
h(\mathbf{x}) &= (2\pi)^{-n/2} \exp\left\{-\frac{1}{2} \sum_{i=1}^n x_i^2\right\} \\
g(T(\mathbf{x}); \vartheta) &= \exp\left\{-\frac{n\vartheta^2}{2} + \vartheta \sum_{i=1}^n x_i\right\}.
\end{aligned}$$

Consider now a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \mathcal{N}(\mu, \sigma^2)$ , with  $\mu$  unknown. Then,

$$\begin{aligned} f_{\mathbf{X}}(\mathbf{x}; \mu, \sigma^2) &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n + \bar{x}_n - \mu)^2 \right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \left[ \sum_{i=1}^n (x_i - \bar{x}_n)^2 + n(\bar{x}_n - \mu)^2 \right] \right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \right\} \exp \left\{ -\frac{n(\bar{x}_n - \mu)^2}{2\sigma^2} \right\}. \end{aligned}$$

We can define

$$h(\mathbf{x}) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \right\}$$

which does not depend on the unknown parameter  $\mu$ . The factor that contains  $\mu$  depends on the sample  $\mathbf{x}$  only through the sample mean  $T(\mathbf{x}) = \bar{x}_n$ . So we have

$$g(t; \mu, \sigma^2) = \exp \left\{ -\frac{n(\bar{x}_n - \mu)^2}{2\sigma^2} \right\},$$

and note that

$$f_{\mathbf{X}}(\mathbf{x}; \mu, \sigma^2) = h(\mathbf{x})g(T(\mathbf{x}); \mu).$$

Thus, by Theorem 4.4,  $T(\mathbf{X}) = \bar{X}_n$  is a sufficient statistic for  $(\mu, \sigma^2)$ .

#### 4.2.2 Minimal Sufficient Statistic

In the preceding section, we found one sufficient statistic for each model considered. In any problem there are, however, many sufficient statistics.

It is always true that the complete sample  $\mathbf{X}$  is a sufficient statistic. We can factor the probability density function (or probability mass function) of  $\mathbf{X}$  as

$$f_{\mathbf{X}}(\mathbf{x}; \vartheta) = g(T(\mathbf{x}); \vartheta)h(\mathbf{x}),$$

where

$$g(T(\mathbf{x}); \vartheta) = \mathcal{L}(\vartheta; \mathbf{x}) \quad \text{and} \quad h(\mathbf{x}) = 1, \quad \forall \mathbf{x} \in \mathcal{X}.$$

By Theorem 4.4,  $T(\mathbf{X}) = \mathbf{X}$  is a sufficient statistic.

Also, it follows that any one-to-one function of a sufficient statistic is a sufficient statistic. Suppose  $T(\mathbf{X})$  is a sufficient statistic and define

$$T^*(\mathbf{x}) = r(T(\mathbf{x})), \quad \forall \mathbf{x} \in \mathcal{X},$$

where  $r(\cdot)$  is a one-to-one function with inverse  $r^{-1}$ . Then, by Theorem 4.4, there exist  $g$  and  $h$  such that

$$\mathcal{L}(\vartheta; \mathbf{x}) = g(T(\mathbf{x}); \vartheta)h(\mathbf{x}) = g(r^{-1}(T^*(\mathbf{x}); \vartheta))h(\mathbf{x}).$$

Defining  $g^*(t; \vartheta) = g(r^{-1}(t); \vartheta)$ , we see that

$$\mathcal{L}(\vartheta; \mathbf{x}) = g^*(T^*(\mathbf{x}); \vartheta)h(\mathbf{x}).$$

So, by Theorem 4.4,  $T^*(\mathbf{X})$  is a sufficient statistic for  $\vartheta$ .

Because of the numerous sufficient statistics in a problem, we might ask whether one sufficient statistic is any better than another. Recall that the purpose of a sufficient statistic is to achieve data reduction without loss of information about the parameter  $\vartheta$ ; thus, a statistic that achieves the most data reduction while still retaining all the information about  $\vartheta$  might be considered preferable.

**Definition 4.8.** A sufficient statistic  $T(\mathbf{X})$  is called a **minimal sufficient statistic** if, for any other sufficient statistic  $T'(\mathbf{X})$ ,  $T(\mathbf{x})$  is a function of  $T'(\mathbf{x})$ .

To say that  $T(\mathbf{x})$  is a function of  $T'(\mathbf{x})$  simply means that if  $T'(\mathbf{x}) = T'(\mathbf{y})$ , then  $T(\mathbf{x}) = T(\mathbf{y})$ . In terms of the partition sets, if  $\{B_{t'} \mid t' \in \mathcal{J}'\}$  are partition sets for  $T'(\mathbf{x})$  and  $\{A_t \mid t \in \mathcal{J}\}$  are the partition sets for  $T(\mathbf{x})$ , then Definition 4.8 states that every  $B_{t'}$  is a subset of some  $A_t$ . Thus, the partition associated with a minimal sufficient statistic is the “largest” possible partition for a sufficient statistic and a minimal sufficient statistic achieves the greatest possible data reduction for a sufficient statistic.

## 4.3 18/10/18

Using Definition 4.8 to find a minimal sufficient statistic is impractical, as was using Definition 4.1 to find sufficient statistics. We would need to guess that  $T(\mathbf{X})$  was a minimal sufficient statistic and then verify the condition in the definition. Fortunately, the following result gives an easier way to find a minimal sufficient statistic.

**Theorem 4.9 (Lehmann-Scheffé).** *Let  $f_{\mathbf{X}}(\mathbf{x}; \vartheta)$  be the probability density function (or probability mass function) of a random sample  $\mathbf{X}$ . Suppose there exist a function  $T(\mathbf{x})$  such that, for two sample points  $\mathbf{x}$  and  $\mathbf{y}$ , the ratio*

$$\frac{f_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{y}; \vartheta)}$$

*is constant as a function of  $\vartheta$  if and only if  $T(\mathbf{x}) = T(\mathbf{y})$ . Then,  $T(\mathbf{X})$  is a minimal sufficient statistic for  $\vartheta$ .*

*Proof.* In order to simplify the proof, we assume  $f_{\mathbf{X}}(\mathbf{x}; \vartheta) > 0$  for every  $\mathbf{x} \in \mathcal{X}$ .

First, we show that  $T(\mathbf{X})$  is a sufficient statistic. Let

$$\mathcal{T} := \{t \mid t = T(\mathbf{x}), \mathbf{x} \in \mathcal{X}\}$$

be the image of  $\mathcal{X}$  under  $T(\cdot)$ . Define the *partition sets induced by  $T(\cdot)$*  as

$$\mathcal{A}_t = \{\mathbf{x} \mid T(\mathbf{x}) = t\}.$$

For each  $\mathcal{A}_t$ , choose and fix one element  $\mathbf{x} \in \mathcal{A}_t$ . For any  $\mathbf{x} \in \mathcal{X}$ ,  $\mathbf{x}_{T(\mathbf{x})}$  is the fixed element that is in the same set,  $\mathcal{A}_t$ , as  $\mathbf{x}$ . Since  $\mathbf{x}$  and  $\mathbf{x}_{T(\mathbf{x})}$  are in the same set  $\mathcal{A}_t$ ,

$$T(\mathbf{x}) = T(\mathbf{x}_{T(\mathbf{x})})$$

and hence

$$\frac{f_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}_{T(\mathbf{x})}; \vartheta)}$$

is a constant as a function of  $\vartheta$  by assumption. Thus, we can define a function on  $\mathcal{X}$  by

$$h(\mathbf{x}) = \frac{f_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}_{T(\mathbf{x})}; \vartheta)}$$

and observe  $h$  does not depend on  $\vartheta$ . Now, define a function on  $\mathcal{T}$  by

$$g(t; \vartheta) = f_{\mathbf{X}}(\mathbf{x}_t; \vartheta).$$

Then, it can be seen that

$$f_{\mathbf{X}}(\mathbf{x}; \vartheta) = \frac{f_{\mathbf{X}}(\mathbf{x}_{T(\mathbf{x})}; \vartheta) f_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}_{T(\mathbf{x})}; \vartheta)} = g(T(\mathbf{x}); \vartheta) h(\mathbf{x})$$

and, by Theorem 4.4  $T(\mathbf{X})$  is a sufficient statistic for  $\vartheta$ .

To show that  $T(\mathbf{X})$  is minimal, let  $T'(\mathbf{X})$  be any other sufficient statistic. By Theorem 4.4 there exist functions  $g'$  and  $h'$  such that

$$f_{\mathbf{X}}(\mathbf{x}; \vartheta) = g'(T'(\mathbf{x}); \vartheta)h'(\mathbf{x}).$$

Let  $\mathbf{x}$  and  $\mathbf{y}$  be any two sample points with  $T'(\mathbf{x}) = T'(\mathbf{y})$ . Then,

$$\frac{f_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{y}; \vartheta)} = \frac{g'(T'(\mathbf{x}); \vartheta)h'(\mathbf{x})}{g'(T'(\mathbf{y}); \vartheta)h'(\mathbf{y})} = \frac{h'(\mathbf{x})}{h'(\mathbf{y})}.$$

Since this ratio does not depend on  $\vartheta$ , the assumptions of the theorem imply that  $T(\mathbf{x}) = T(\mathbf{y})$ . Thus,  $T(\mathbf{x})$  is a function of  $T'(\mathbf{x})$  and we conclude that  $T(\mathbf{x})$  is minimal.  $\square$

*Remark 4.10.* The notion of *minimal sufficiency* is good, because it allows us to choose only statistics with the less number of block in the partition. However, this is still not enough, since there are many sufficient statistics and we need further properties to distinguish them.

#### 4.3.1 Methods of Evaluating Estimators

The methods discussed in the previous section have outlined reasonable techniques for finding point estimators of parameters. A difficulty that arises, however, is that since we can usually apply more than one of these methods in a particular situation, we are often faced with the task of choosing between estimators. Of course, it is possible that different methods of finding estimators will yield the same answer, which makes evaluation a bit easier, but, in many cases, different methods will lead to different estimators.

**Definition 4.11.** The **bias** of a point estimator  $T_n$ , of a parameter  $\vartheta$ , is the difference between the expected value of  $T_n$  and  $\vartheta$ , that is

$$b(T_n) := \mathbb{E}[T_n] - \vartheta.$$

An estimator  $T_n$  whose bias is identically equal to zero is called **unbiased** for  $\vartheta$  and satisfies  $\mathbb{E}[T_n] = \vartheta$ .

*Remark 4.12.* This is a good property, because it is preferable for a statistic who is estimating  $\vartheta$  to be centred in  $\vartheta$ . However, for completeness, we are looking for a measure of dispersion around the expected value, which is the variance of  $T_n$ . Such value measures the goodness of the estimator only in the case it is unbiased, since otherwise the variance would measure the dispersion around something that is not of interest. Indeed,

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

For a more general approach, consider Chebyshev's Inequality (taking the complementary)

$$\mathbb{P}(|T_n - \vartheta| < k) \geq 1 - \frac{\mathbb{E}[(T_n - \vartheta)^2]}{k^2},$$

which suggests to control the following quantity.

**Definition 4.13.** The **Mean Squared Error** (MSE) of an estimator  $T_n$  of a parameter  $\vartheta$  is the function of  $\vartheta$  described by

$$\text{MSE}(T_n) := \mathbb{E}[(T_n - \vartheta)^2].$$

Notice that the Mean Squared Error measures the average squared difference between the estimator  $T_n$  and the parameter  $\vartheta$ , a somewhat reasonable measure of performance for a point estimator. In general, any increasing function of the absolute distance  $|T_n - \vartheta|$  would serve to measure the goodness of an estimator (*mean absolute error*,  $\mathbb{E}[|T_n - \vartheta|]$ , is a reasonable alternative), but the Mean Squared Error has at least two advantages over the distance measures. First, it is quite tractable analytically and, second, it has the interpretation

$$\begin{aligned} \text{MSE}(T_n) &= \mathbb{E}[(T_n - \vartheta)^2] \\ &= \mathbb{E}[(T_n - \mathbb{E}[T_n] + \mathbb{E}[T_n] - \vartheta)^2] \\ &= \mathbb{E}[(T_n - \mathbb{E}[T_n])^2] + \mathbb{E}[(T_n - \vartheta)^2] + \underbrace{2\mathbb{E}[(T_n - \mathbb{E}[T_n])(\mathbb{E}[T_n] - \vartheta)]}_{=0} \\ &= \text{Var}(T_n) + (\mathbb{E}[T_n] - \vartheta)^2 \\ &= \text{Var}(T_n) + (b(T_n))^2 \end{aligned}$$

*Remark 4.14.* Thus, the Mean Squared Error incorporates two components, one measuring the variability of the estimator (*precision*), and the other measuring its bias (*accuracy*). An estimator that has good Mean Squared Error properties has small combined variance and bias. To find one, we need to find estimator that control both variance and bias. Clearly, unbiased estimators do a good job of controlling bias. For an unbiased estimator we have

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

and so, if an estimator is unbiased, its Mean Squared Error is equal to its variance.

Although many unbiased estimators are also reasonable from the standpoint of Mean Squared Error, be aware that controlling bias does not guarantee that Mean Squared Error is controlled. In particular, it is sometimes the case that a trade-off occurs between variance and bias in such a way that a small increase in bias can be traded for a larger decrease in variance, resulting in an improvement in Mean Squared Error.

**Definition 4.15.** Consider two estimators  $T'_n$  and  $T''_n$  for  $\vartheta$ . We say that  $T'_n$  is **more efficient** than  $T''_n$  if  $\text{MSE}(T'_n) < \text{MSE}(T''_n)$ .

## 5

## Week V

## 5.1 23/10/18

## 5.1.1 The Cramér-Rao Lower Bound

**Theorem 5.1 (Cramér- Rao).** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from a regular parametric model  $X \sim f_X(x; \vartheta)$ . Then, for any estimator  $T_n$  of  $\vartheta$ , we have that

$$\text{Var}(T_n) \geq \frac{(1 + b'(T_n))^2}{I_n(\vartheta)}, \quad (5.1)$$

where  $b(T_n)$  is the bias of  $T_n$  and  $I_n$  is the Fisher information.

*Proof.* The proof of this theorem is elegantly simple, and is a clever application of the Cauchy-Schwarz inequality, or stated statistically, the fact that any two random variables  $X$  and  $Y$ ,

$$\begin{aligned} (\text{Cov}(X, Y))^2 &= (\mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])])^2 \\ &\leq \mathbb{E}[(X - \mathbb{E}[X])^2] \mathbb{E}[(Y - \mathbb{E}[Y])^2] \\ &= \text{Var}(X) \text{Var}(Y). \end{aligned}$$

If we rearrange this, we can get a lower bound on the variance of  $X$ ,

$$\text{Var}(X) \geq \frac{(\text{Cov}(X, Y))^2}{\text{Var}(Y)}.$$

Observe that

$$\mathbb{E}[T_n] = \vartheta + b(T_n) \quad \text{and} \quad \frac{d}{d\vartheta} \mathbb{E}[T_n] = 1 + b'(T_n).$$

We know that, from Proposition 3.16,

$$\mathbb{E}[V_n'(\vartheta)] = 0,$$

therefore

$$\text{Cov}(T_n, V_n'(\vartheta)) = \mathbb{E}[T_n V_n'(\vartheta)] - \mathbb{E}[T_n] \mathbb{E}[V_n'(\vartheta)] = \mathbb{E}[T_n V_n'(\vartheta)].$$

Continuing our evaluation, we obtain

$$\begin{aligned}
\mathbb{E}[T_n V_n'(\vartheta)] &= \int_{\mathbb{R}^n} T_n V_n'(\vartheta) f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} \\
&= \int_{\mathbb{R}^n} T_n \frac{f'_{\mathbf{X}}(\mathbf{x}; \vartheta)}{f_{\mathbf{X}}(\mathbf{x}; \vartheta)} f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} \\
&= \int_{\mathbb{R}^n} T_n \frac{d}{d\vartheta} f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} \\
&= \frac{d}{d\vartheta} \int_{\mathbb{R}^n} T_n f_{\mathbf{X}}(\mathbf{x}; \vartheta) d\mathbf{x} \\
&= \frac{d}{d\vartheta} \mathbb{E}[T_n] \\
&= 1 + b'(T_n).
\end{aligned}$$

Since  $\mathbb{E}[V_n'(\vartheta)] = 0$ , it follows that

$$\text{Var}(V_n'(\vartheta)) = \mathbb{E}[(V_n'(\vartheta))^2] = -\mathbb{E}[V_n''(\vartheta)] = I_n(\vartheta).$$

Then, by Cauchy-Schwarz inequality, we conclude that

$$\text{Var}(T_n) \geq \frac{(\text{Cov}(T_n, V_n'(\vartheta)))^2}{\text{Var}(V_n'(\vartheta))} = \frac{(1 + b'(T_n))^2}{I_n(\vartheta)}.$$

□

*Remark 5.2.* Remember that, since

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

we obtain

$$\text{MSE}(T_n) \geq \frac{(1 + b'(\vartheta))^2}{I_n(\vartheta)} + b^2(T_n).$$

Also, the variance of an unbiased estimator  $T_n(\vartheta)$  cannot be smaller than the Cramér-Rao lower bound (5.1). In general, however, we do not know if there exist an estimator whose variance is equal to the Cramér-Rao lower bound. The proper lower bound involves the Mean Squared Error

$$\text{MSE}(T_n) = \frac{(1 + b'(\vartheta))^2}{I_n(\vartheta)} + b^2(\vartheta).$$

If we add the assumption of independent samples, then the calculation of the lower bound is simplified. The expectation in the denominator becomes a univariate calculation, as the following corollary shows.

**Corollary 5.3.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from a regular parametric model  $X \sim f_X(x; \vartheta)$ . If there exist an unbiased estimator  $T_n$  for  $\vartheta$  whose variance is equal to the Cramér-Rao lower bound (5.1), then the estimator is unique.

*Proof.* Consider two unbiased estimators  $T_{1n}$  and  $T_{2n}$  for  $\vartheta$  such that

$$\text{Var}(T_{1n}) = \text{Var}(T_{2n}) = \frac{1}{I_n(\vartheta)} =: v.$$

Define another estimator  $T_n$  as follows:

$$T_n := \frac{T_{1n} + T_{2n}}{2}.$$

This is unbiased, since

$$\mathbb{E}[T_n] = \mathbb{E}\left[\frac{T_{1n} + T_{2n}}{2}\right] = \frac{2\vartheta}{2} = \vartheta.$$

Now,

$$\begin{aligned} \text{Var}(T_n) &= \text{Var}\left(\frac{T_{1n} + T_{2n}}{2}\right) \\ &= \frac{1}{4} [\text{Var}(T_{1n}) + \text{Var}(T_{2n}) + 2 \text{Cov}(T_{1n}, T_{2n})] \\ &= \frac{(\text{Var}(T_{1n}) \text{Var}(T_{2n}))^{1/2}}{4} \frac{1}{(\text{Var}(T_{1n}) \text{Var}(T_{2n}))^{1/2}} [\text{Var}(T_{1n}) + \text{Var}(T_{2n}) + 2 \text{Cov}(T_{1n}, T_{2n})] \\ &= \frac{1}{2} \left( v + \frac{\text{Cov}(T_{1n}, T_{2n})}{(\text{Var}(T_{1n}) \text{Var}(T_{2n}))^{1/2}} v \right) \\ &= \frac{1}{2} (1 + \rho) v, \end{aligned}$$

where  $\rho$  is the correlation coefficient. Now, if  $\rho < 1$ , we get that

$$\text{Var}(T_n) = \frac{1}{2} (1 + \rho) v$$

is less than the Cramér-Rao lower bound [\(5.1\)](#), which is impossible. Therefore,  $\rho = 1$ , meaning

$$\text{Cor}(T_{1n}, T_{2n}) = 1 \quad T_{2n} = a + bT_{1n},$$

from which we obtain

$$\mathbb{E}[T_{2n}] = \mathbb{E}[a + bT_{1n}] \implies \vartheta = a + b\vartheta \Leftrightarrow a = 0 \wedge b = 1.$$

This implies that  $T_{1n} = T_{2n}$ . Therefore, it exist a unique unbiased estimator whose variance is equal to the Cramér-Rao lower bound [\(5.1\)](#).  $\square$

**Definition 5.4.** Consider a regular model  $X \sim f_X(x; \vartheta)$ . We say that an estimator is **efficient** if its variance is equal to the Cramér-Rao lower bound [\(5.1\)](#). Moreover, we say that an unbiased efficient estimator  $T_n$  for  $\vartheta$  is **absolute efficient**, i.e.

$$\text{Var}(T_n) = \frac{1}{I_n(\vartheta)}.$$

Finally, the **efficiency** of  $T_n$  is defined as

$$\text{eff}(T_n) := \frac{1}{\text{Var}(T_n) I_n(\vartheta)} \in [0, 1].$$

*Remark 5.5.* We introduced the absolute efficiency at the cost of assuming regularity for the parametric model.

**Proposition 5.6.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from a regular parametric model  $X \sim f_X(x; \vartheta)$ . Let  $T_n$  be an unbiased estimator for  $\vartheta$ . Then,  $T_n$  is efficient for  $\vartheta$  if and only if

$$V'_n(\vartheta) = I_n(\vartheta)(T_n - \vartheta). \quad (5.2)$$

*Proof.* Cramér-Rao inequality (5.1) can be expressed as

$$(\text{Cov}(T_n, V'_n(\vartheta)))^2 \leq \text{Var}(T_n) I_n(\vartheta).$$

Suppose we have efficiency, then

$$\text{Var}(T_n) = \frac{1}{I_n(\vartheta)}$$

if and only if  $\text{Cor}(T_n, V'_n(\vartheta)) = 1$ , if and only if  $T_n$  is proportional to  $V'_n(\vartheta)$ , meaning

$$V'_n(\vartheta) = aT_n + b. \quad (5.3)$$

Now,

$$\mathbb{E}[V'_n(\vartheta)] = a\mathbb{E}[T_n] + b \implies 0 = a\vartheta + b$$

hence  $b = -a\vartheta$ . This means that

$$I_n(\vartheta) = \text{Var}(V'_n(\vartheta)) = a^2 \text{Var}(T_n) = \frac{a^2}{I_n(\vartheta)} \implies a = I_n(\vartheta),$$

( $a$  cannot be equal to  $-I_n(\vartheta)$ , because in this case we would have  $\text{Cor}(T_n, V'_n(\vartheta)) = -1$ ) therefore (5.3) becomes

$$V'_n(\vartheta) = I_n(\vartheta)(T_n - \vartheta).$$

Conversely, imagine that the equality (5.2) holds: now,

$$\begin{aligned} V'_n(\vartheta) &= I_n(T_n - \vartheta) \\ T_n &= \frac{V'_n(\vartheta)}{I_n} + \vartheta \end{aligned}$$

and we simply evaluate its variance, meaning

$$\text{Var}(T_n) = \text{Var}\left(\frac{V'_n(\vartheta)}{I_n} - \vartheta\right) = \frac{1}{I_n^2} \text{Var}(V'_n(\vartheta)) = \frac{1}{I_n}.$$

Therefore,  $T_n$  is (absolutely) efficient by definition.  $\square$

**Theorem 5.7 (Rao-Blackwell).** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from a parametric model  $X \sim f_X(x; \vartheta)$ . Let  $T_{1n}$  be a sufficient estimator for  $\vartheta$  and  $T_{2n}$  be an unbiased estimator for  $\vartheta$ . Define

$$T_n := \mathbb{E}[T_{2n} | T_{1n}].$$

Then,

1.  $T_n$  is a function of  $T_{1n}$ ;
2.  $\mathbb{E}[T_n] = \vartheta$ ;
3.  $\text{Var}(T_n) < \text{Var}(T_{2n})$ .

### 5.1.2 Convergence

This section treats the somewhat fanciful idea of allowing the sample size to approach infinity and investigates the behaviour of certain sample quantities as this happens. Although the notion of an infinite sample size is a theoretical artefact, it can often provide us with some useful approximations for the finite-sample case, since it usually happens that expressions become simplified in the limit.

We are mainly concerned with three types of convergence, and treat them in varying amount of detail. But first, let us recall some useful definitions.

**Definition 5.8.** A sequence of random variables  $X_1, X_2, \dots$  is said to **converge in distribution** to a random variable  $X$  if

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$$

at all points  $x$  where  $F_X(x)$  is continuous.

**Definition 5.9.** A sequence of random variables  $X_1, X_2, \dots$  is said to **converge in probability** to a random variable  $X$  if, for every  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| \geq \varepsilon) = 0,$$

or equivalently

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| < \varepsilon) = 1$$

**Definition 5.10.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim f_X(x; \vartheta)$ . We say that  $T_n$  is **asymptotically unbiased** for  $\vartheta$  if

$$\lim_{n \rightarrow \infty} \mathbb{E}[T_n] = \vartheta,$$

or equivalently

$$\lim_{n \rightarrow \infty} b(T_n) = 0.$$

We say that an estimator  $T_n$  is **consistent in mean-square** for  $\vartheta$  if

$$\lim_{n \rightarrow \infty} \text{MSE}(T_n) = 0. \quad (5.4)$$

Since  $\text{MSE}(T_n) = \text{Var}(T_n) + b^2(T_n)$ , (5.4) is equivalent to say that

$$\lim_{n \rightarrow \infty} \text{Var}(T_n) = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} b^2(T_n) = 0.$$

A consistent estimator (in mean-square) is also asymptotically unbiased.

Finally, we say that  $T_n$  is **consistent in probability** for  $\vartheta$  if, for every  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|T_n - \vartheta| < \varepsilon) = 1, \quad \text{or} \quad \lim_{n \rightarrow \infty} \mathbb{P}(|T_n - \vartheta| \geq \varepsilon) = 0.$$

**Remark 5.11.** The consistency in mean-square implies the consistency in probability:

$$\mathbb{P}(|T_n - \vartheta| < \varepsilon) \geq 1 - \frac{\text{MSE}(T_n)}{\varepsilon^2}$$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|T_n - \vartheta| < \varepsilon) \geq 1 - \lim_{n \rightarrow \infty} \frac{\text{MSE}(T_n)}{\varepsilon^2}.$$

**Exercise 5.12.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \mathcal{N}(\mu, \sigma^2)$ . For

1. the corrected sample variance

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

2. the uncorrelated sample variance

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

show that (1) is unbiased for  $\text{Var}(X) = \sigma^2$  and (2) is biased for  $\text{Var}(X) = \sigma^2$ .

*Proof.* (1) First,

$$\begin{aligned} \mathbb{E}[\hat{S}_n^2] &= \mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2\right] \\ &= \mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n X_i^2 - \frac{2\bar{X}_n}{n-1} \sum_{i=1}^n X_i + \frac{n}{n-1} \mathbb{E}[\bar{X}_n^2]\right] \\ &= \frac{1}{n-1} \sum_{i=1}^n \mathbb{E}[X_i^2] - \frac{2n}{n-1} \mathbb{E}[\bar{X}_n^2] + \frac{n}{n-1} \mathbb{E}[\bar{X}_n^2] \\ &= \frac{1}{n-1} \sum_{i=1}^n \mathbb{E}[X_i^2] - \frac{n}{n-1} \mathbb{E}[\bar{X}_n^2]. \end{aligned} \quad (5.5)$$

Now,

$$\begin{aligned} \mathbb{E}[X_i^2] &= \text{Var}(X_i) + (\mathbb{E}[X_i])^2 = \sigma^2 + \mu^2, \\ \mathbb{E}[\bar{X}_n^2] &= \text{Var}(\bar{X}_n) + (\mathbb{E}[\bar{X}_n])^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) + \mu^2 \\ &= \frac{\sigma^2}{n} + \mu^2, \end{aligned}$$

therefore

$$\begin{aligned} (5.5) &= \frac{1}{n-1} \sum_{i=1}^n (\sigma^2 + \mu^2) - \frac{n}{n-1} \left(\frac{\sigma^2}{n} + \mu^2\right) \\ &= \frac{n}{n-1} (\sigma^2 + \mu^2) - \frac{n}{n-1} \left(\frac{\sigma^2}{n} + \mu^2\right) \\ &= \frac{n}{n-1} \sigma^2 - \frac{\sigma^2}{n-1} \\ &= \sigma^2. \end{aligned}$$

This shows that the corrected sample variance is an unbiased estimator for the variance of a model.  $\square$

**Exercise 5.13.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim f_X(x; \vartheta)$  such that  $\mathbb{E}[X] = \vartheta$  and  $\text{Var}(X) = 1$ . Check if the following estimators are unbiased:

1.  $T_{n1} = \frac{1}{n} \sum_{i=1}^n X_i$ ;
2.  $T_{n2} = \frac{2}{n(n+1)} \sum_{i=1}^n i X_i$ ;
3.  $T_{n3} = \sum_{i=1}^n (-1)^{i+1} X_i$ ;
4.  $T_{n4} = \frac{1}{2^n - 1} \sum_{i=1}^n \binom{n}{i} X_i$ .

*Proof.* (1) *Unbiased:*

$$\mathbb{E}[T_{1n}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = \frac{n\vartheta}{n} = \vartheta.$$

(2) *Unbiased:*

$$\mathbb{E}[T_{n2}] = \frac{2}{n(n+1)} \sum_{i=1}^n i \mathbb{E}[X_i] = \frac{2\vartheta}{n(n+1)} \sum_{i=1}^n i = \frac{2\vartheta}{n(n+1)} \frac{n(n+1)}{2} = \vartheta.$$

(3) *Unbiased, when  $n$  is odd:*

$$\mathbb{E}[T_{n3}] = \sum_{i=1}^n \mathbb{E}[(-1)^{i+1} X_i] = \vartheta \sum_{i=1}^n (-1)^{i+1} = \vartheta$$

(4) Observe that

$$\mathbb{P}(X = i) = \binom{n}{i} \vartheta^i (1 - \vartheta)^{n-i}$$

and

$$\sum_{i=0}^n \binom{n}{i} \vartheta^i (1 - \vartheta)^{n-i} = 1.$$

If we take  $\vartheta = 1/2$ ,

$$\left(\frac{1}{2}\right)^n \sum_{i=0}^n \binom{n}{i} = 1 \implies \sum_{i=1}^n \binom{n}{i} = 2^n - 1.$$

Therefore,  $T_{n4}$  is unbiased. □

## 5.2 24/10/18

### 5.2.1 Maximum Likelihood Estimator

Let  $X_1, \dots, X_n$  be an independent and identically distributed random sample from a population  $X \sim f_X(x; \vartheta)$ . As we have seen, the Likelihood function is given by

$$\mathcal{L}(\vartheta; \mathbf{x}) = \prod_{i=1}^n f(x_i; \vartheta_1, \dots, \vartheta_n).$$

**Definition 5.14.** For each sample point  $\mathbf{X} = \mathbf{x}$ , let  $\hat{\vartheta}(\mathbf{x})$  be a parameter value at which  $\mathcal{L}(\vartheta, \mathbf{x})$  attains its maximum as a function of  $\vartheta$ , with  $\mathbf{x}$  fixed. A **maximum Likelihood estimator** of the parameter  $\vartheta$ , based on a sample  $\mathbf{X}$  is

$$\arg \max_{\vartheta \in \Theta} \log \mathcal{L}(\vartheta, \mathbf{x}). \quad (5.6)$$

This can be done in the following situations:

1. when the model is *regular*: this allows us to solve the Maximum Likelihood equation

$$V'_n(\vartheta) = 0; \quad (5.7)$$

2. when the parameter  $\vartheta$  is a (positive) integer, take the ratio

$$\frac{f_{n+1}(x)}{f_n(x)}.$$

**Exercise 5.15.** A computer takes an exponential time to solve a task. Assume that each time is modelled as a negative exponential random variable  $T$ , meaning

$$X \sim f_T(t; \vartheta) = \vartheta e^{-\vartheta t} \mathbf{1}_{\mathbb{R}_+}(t).$$

Suppose we have  $n$  independent tasks are assigned to the computer, and let  $T_i \stackrel{d}{=} T$  be the times, for  $i = 1, \dots, n$ .

1. Show that

$$\bar{T}_n = \frac{1}{n} \sum_{i=1}^n T_i$$

is an unbiased estimator for  $1/\vartheta$ , where  $\vartheta$  is the parameter of the negative exponential distribution.

2. Find the distribution of  $\bar{T}_n$ : can you say that  $\tilde{T}_n = (\bar{T}_n)^{-1}$  is an unbiased estimator for  $\vartheta$ ?
3. Find an unbiased estimator for  $\vartheta$  and compute the corresponding Mean Squared Error.
4. Are these estimators consistent for  $\vartheta$  and  $1/\vartheta$ ?

*Proof.* (1) First of all,  $\bar{T}_n$  is clearly unbiased, since

$$\mathbb{E}[\bar{T}_n] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n T_i\right] = \frac{1}{n} \sum_{i=1}^n \frac{1}{\vartheta} = \frac{1}{\vartheta}.$$

Also,

$$\text{Var}(\bar{T}_n) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(T_i) = \frac{1}{n^2} \frac{n}{\vartheta^2} = \frac{1}{n\vartheta^2}.$$

(2) We know that  $\sum_{i=1}^n T_i \sim \Gamma(n, \vartheta)$ . Define  $X := \sum_{i=1}^n T_i$ : therefore, our goal is to find the distribution of  $\frac{X}{n}$ :

$$F_{X/n}(t) = \mathbb{P}\left(\frac{X}{n} \leq t\right) = \mathbb{P}(X \leq nt) = F_X(nt),$$

hence,

$$f_{X/n}(t) = nf_X(nt) \implies \frac{T_i}{n} \sim \mathcal{E}(n\vartheta) = \Gamma(1, n\vartheta).$$

Therefore,

$$\frac{1}{n} \sum_{i=1}^n T_i = \sum_{i=1}^n \frac{T_i}{n} \sim \Gamma(n, n\vartheta).$$

Now,

$$\begin{aligned} \mathbb{E}[\tilde{T}_n] &= \mathbb{E}\left[\frac{n}{\sum_{i=1}^n T_i}\right] \\ &= \int_0^\infty \frac{n}{t} \frac{\vartheta^n}{\Gamma(n)} t^{n-1} e^{-\vartheta t} dt \\ &= \frac{\vartheta^n n}{\Gamma(n)} \int_0^\infty t^{n-2} e^{-\vartheta t} dt \\ &= \frac{\vartheta^n n}{\Gamma(n)} \frac{\Gamma(n-1)}{\vartheta^{n-1}} \underbrace{\int_0^\infty \frac{\vartheta^{n-1}}{\Gamma(n-1)} t^{n-2} e^{-\vartheta t} dt}_{\|1} \\ &= \frac{\vartheta^n n \Gamma(n-1)}{\Gamma(n) \vartheta^{n-1}} \\ &= \frac{\vartheta n}{n-1}, \end{aligned}$$

therefore we conclude that it is indeed biased, but asymptotically unbiased.

(3) We want to find an unbiased estimator for  $\vartheta$ ;  $S_n$  must be unbiased for  $\vartheta$ . Indeed,

$$\mathbb{E}[S_n] = \mathbb{E}\left[\frac{n-1}{n} \hat{S}_n\right] = \frac{n-1}{n} \mathbb{E}[\tilde{T}_n] = \frac{n-1}{n} \frac{n\vartheta}{n-1} = \vartheta,$$

therefore it is unbiased. The corresponding Mean Squared Error is

$$\begin{aligned} \text{MSE}(S_n) &= \text{Var}(S_n) + b^2(S_n) \\ &= \text{Var}(S_n) \\ &= \mathbb{E}[S_n^2] - \vartheta^2 \\ &= \mathbb{E}\left[\left(\frac{n-1}{\sum_{i=1}^n T_i}\right)^2\right] - \vartheta^2 \end{aligned}$$

$$\begin{aligned}
&= (n-1)^2 \int_0^{+\infty} \frac{1}{t^2} \frac{\vartheta^n}{\Gamma(n)} t^{n-3} e^{-\vartheta t} dt - \vartheta^2 \\
&= \frac{(n-1)^2 \vartheta^n}{\Gamma(n)} \frac{\Gamma(n-2)}{\vartheta^{n-2}} - \vartheta^2 \\
&= \vartheta^2 (n-1)^2 \frac{\Gamma(n-2)}{\Gamma(n)} - \vartheta^2 \\
&= \vartheta^2 (n-1)^2 \frac{(n-3)!}{(n-1)(n-2)(n-3)!} - \vartheta^2 \\
&= \frac{\vartheta^2}{n-2}.
\end{aligned}$$

(4) Is

$$\bar{T}_n = \frac{1}{n} \sum_{i=1}^n T_i$$

consistent for  $1/\vartheta$ ? Yes, by the Weak Law of Large Numbers, since

$$\bar{T}_n \rightarrow \frac{1}{\vartheta}.$$

Is  $S_n$  consistent for  $\vartheta$ ? Also, yes. Indeed

$$S_n = \frac{n-1}{\sum_{i=1}^n T_i} = \frac{(n-1)/n}{(\sum_{i=1}^n T_i)/n} \xrightarrow{n \rightarrow \infty} \frac{1}{\vartheta^{-1}} = \vartheta. \quad \square$$

**Exercise 5.16.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim f_X(x; \vartheta)$  such that

$$f_X(x; \vartheta) = \vartheta x^{\vartheta-1} \mathbb{1}_{(0,1)}(x).$$

1. Find a sufficient statistic and a Maximum Likelihood estimator for  $\vartheta$ .
2. Is the Maximum Likelihood estimator unbiased for  $\vartheta$ ? Compute the Mean Squared Error.
3. Is the Maximum Likelihood estimator consistent for  $\vartheta$ ?
4. Propose another estimator for  $\vartheta$ .

*Proof.* (1) We have that

$$\begin{aligned}
\mathcal{L}(\vartheta; \mathbf{x}) &= \prod_{i=1}^n \vartheta x_i^{\vartheta-1} \\
&= \vartheta^n \left( \prod_{i=1}^n x_i \right)^{\vartheta-1} \\
&= \underbrace{\vartheta^n \left( \prod_{i=1}^n x_i \right)^{\vartheta}}_{g(T(\mathbf{x}), \vartheta)} \cdot \underbrace{\left( \prod_{i=1}^n x_i \right)^{-1}}_{h(\mathbf{x})},
\end{aligned}$$

and

$$\begin{aligned}\log \mathcal{L}(\vartheta; x) &= n \log \vartheta + (\vartheta - 1) \log \left( \prod_{i=1}^n x_i \right) \\ &= n \log \vartheta + (\vartheta - 1) \left( \sum_{i=1}^n \log x_i \right).\end{aligned}\quad (5.8)$$

Therefore, thanks to the Factorisation Theorem,  $\prod_{i=1}^n X_i$  is a sufficient statistic for  $\vartheta$ , whilst the Maximum Likelihood estimator is obtained deriving (5.8):

$$0 = \frac{n}{\vartheta} + \sum_{i=1}^n \log x_i \implies \hat{\vartheta}_n = -\frac{n}{\sum_{i=1}^n \log x_i}.$$

(2) First of all, we compute the distribution of  $Y := -\log X$ :

$$\begin{aligned}\mathbb{P}(Y \leq y) &= \mathbb{P}(-\log X \leq y) \\ &= \mathbb{P}(X \geq e^{-y}) \\ &= 1 - \mathbb{P}(X < e^{-y}) \\ &= 1 - e^{-y\vartheta},\end{aligned}$$

therefore  $Y \sim \mathcal{E}_n(\vartheta)$ , and consequently  $-\sum_{i=1}^n \log x_i \sim \Gamma(n, \vartheta)$ . This implies that

$$\mathbb{E}[\hat{\vartheta}_n] = \int_0^\infty \frac{n}{t} \frac{\vartheta^n}{\Gamma(n)} t^{n-1} e^{-\vartheta t} dt = \frac{n\vartheta}{n-1}$$

is biased for  $\vartheta$ , but asymptotically unbiased.

(3) Since

$$\hat{\vartheta}_n = -\frac{n}{\sum_{i=1}^n \log x_i},$$

by the Weak Law of Large Numbers

$$-\frac{1}{n} \sum_{i=1}^n \log x_i \xrightarrow{n \rightarrow \infty} \frac{1}{\hat{\vartheta}_n}$$

and

$$\frac{1}{\hat{\vartheta}_n} \xrightarrow{n \rightarrow \infty} \vartheta.$$

This implies that the Maximum Likelihood estimator is (weakly) consistent.

(4) We could use the Method of Moments. We know that the theoretical moment of order 1 of a parametric model is

$$\mathbb{E}[X] = \int_0^1 x \vartheta x^{\vartheta-1} dx = \frac{\vartheta}{\vartheta+1}.$$

We compare the theoretical moment of order 1 with the empirical model of the same order, which is the sample mean:

60 5 Week V

$$\frac{1}{n} \sum_{i=1}^n X_i = \frac{\vartheta}{\vartheta + 1} \implies \bar{\vartheta}_n = \frac{\bar{X}_n}{1 - \bar{X}_n}.$$

Therefore,

$$\bar{\vartheta}_n = \frac{\bar{X}_n}{1 - \bar{X}_n} \rightarrow \frac{\vartheta/(\vartheta + 1)}{1 - \vartheta/(\vartheta + 1)} = \vartheta.$$

□

### 5.3 25/10/18

When sampling from a population described by a probability density function (or probability mass function)  $f_X(x; \vartheta)$ , knowledge of  $\vartheta$  yields knowledge of the entire population. Hence, it is natural to seek a method of finding a good estimator of the point  $\vartheta$ , that is, a good *point estimator*. It is also the case that the parameter  $\vartheta$  has a meaningful physical interpretation (as in the case of a population mean) so there is the direct interest in obtaining a good point estimate of  $\vartheta$ . It may also be the case that some function of  $\vartheta$  is of interest.

**Definition 5.17.** A *point estimator* is any function of a sample. That is, any statistic is a point estimator.

Notice that the definition makes no mention of any correspondence between the estimator and the parameter it is to estimate. While it might be argued that such statement should be included in the definition, such a statement would restrict the available set of estimators. Also, there is no mention in the definition of the range of the statistic.

There is one distinction that must be made clear, the difference between an estimate and an estimator. An *estimator* is a function of the sample, while an *estimate* is the realised value of an estimator (that is, a number) that is obtained when a sample is actually taken. Notationally, when a sample is taken, an estimator is a function of the random variables  $X_1, \dots, X_n$ , while an estimate is a function of the realised values  $x_1, \dots, x_n$ .

In many cases, there will be an obvious or natural candidate for a point estimator of a particular parameter. For example, the sample mean is a natural candidate for a point estimator of the population mean. However, when we leave a simple case like this, intuition may not only desert us, it may also lead us astray. Therefore, it is useful to have some techniques that will at least give us a some reasonable candidates for consideration.

#### 5.3.1 Method of Moments

The Method of Moments is, perhaps, the oldest method of finding point estimators. It has the virtue of being quite simple to use and almost always yields some sort of estimate. In many cases, unfortunately, this method yields estimators that may be improved upon. However, it is a good place to start when other methods prove intractable.

**Definition 5.18.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from a parametric model  $X \sim f_X(x; \vartheta)$ . Assume that  $\vartheta \in \Theta$  such that  $\dim \Theta = n$ . Method of Moments estimators are found by equating the first  $k$  sample moments to the corresponding  $k$  population moments, and solving the resulting system of simultaneous equations. More precisely, define

$$m_j = \frac{1}{n} \sum_{i=1}^n X_i^j, \quad \mu_j = \mathbb{E}[X^j], \quad j = 1, \dots, k$$

The population moment  $\mu_j$  will typically be a function of  $\vartheta_1, \dots, \vartheta_k$ . The **Method of Moments** estimator  $(\tilde{\vartheta}_1, \dots, \tilde{\vartheta}_k)$  of  $(\vartheta_1, \dots, \vartheta_k)$  is obtained by solving the following system of equations of  $(\vartheta_1, \dots, \vartheta_k)$  in terms of  $(m_1, \dots, m_k)$ :

$$\begin{cases} m_1 = \mu_1(\vartheta_1, \dots, \vartheta_k) \\ \vdots \\ m_k = \mu_k(\vartheta_1, \dots, \vartheta_k) \end{cases} \quad (5.9)$$

*Example 5.19.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \mathcal{N}(\mu, \sigma^2)$ . In the preceding notation,  $\vartheta_1 = \mu$  and  $\vartheta_2 = \sigma^2$ . We have that

$$m_1 = \bar{X}_n, \quad m_2 = \frac{1}{n} \sum_{i=1}^n X_i^2, \quad \mu_1 = \mu, \quad \mu_2 = \mu^2 + \sigma^2$$

and hence we must solve

$$\bar{X}_n = \mu, \quad \frac{1}{n} \sum_{i=1}^n X_i^2 = \mu^2 + \sigma^2.$$

Solving for  $\mu$  and  $\sigma^2$  yields the Method of Moments estimators

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

In this simple example, the Method of Moments solutions coincides with our intuition, and perhaps give some credence to both. The method is somewhat more helpful, however, when no obvious estimator suggests itself.

**Exercise 5.20.** We want to estimate the proportion  $\vartheta$  of individuals in a population, for which a certain feature  $X$  takes value in a set  $\mathcal{A}$ . We take a sample of the population of size  $n$  and we measure the feature  $X$ .

Let  $Z_n$  be the number of individuals with feature  $X$  in  $\mathcal{A}$ .

1. Is  $\frac{Z_n}{n}$  unbiased for  $\vartheta$ ? Is it consistent? Is it asymptotically Gaussian?
2. Find a Maximum Likelihood estimator for  $\vartheta$ .

*Proof.* (1) Denote

$$Y_i = \begin{cases} 1 & \text{if } X \text{ takes values in } \mathcal{A} \text{ for an individual } i \\ 0 & \text{otherwise} \end{cases}.$$

Obviously,  $Y_i \sim \text{Ber}(\vartheta)$ . Therefore

$$Z_n = \sum_{i=1}^n Y_i, \quad \frac{Z_n}{n} = \frac{1}{n} \sum_{i=1}^n Y_i$$

and

$$\mathbb{E} \left[ \frac{Z_n}{n} \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} [Y_i] = \vartheta,$$

meaning that it is unbiased.

Also, it is consistent, since

$$\frac{Z_n}{n} \rightarrow \vartheta, \quad n \rightarrow \infty,$$

and asymptotically Gaussian: indeed

$$\sqrt{n} \left( \frac{Z_n}{n} - \vartheta \right) = \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n Y_i - \vartheta \right) = \sqrt{n} (\bar{Y}_n - \vartheta) = \frac{\bar{Y}_n - \vartheta}{\frac{\sigma}{\sqrt{n}}} \rightarrow \mathcal{N}(0, \vartheta(1 - \vartheta)).$$

(2) Take  $Y_1, \dots, Y_n$  from  $Y \sim \text{Ber}(\vartheta)$  and consider

$$\begin{aligned} \log \mathcal{L}(\vartheta; \mathbf{x}) &= \log \left( \prod_{i=1}^n \vartheta^{x_i} (1 - \vartheta)^{n - x_i} \right) \\ &= \log \left( \vartheta^{\sum_{i=1}^n x_i} (1 - \vartheta)^{n - \sum_{i=1}^n x_i} \right) \\ &= \sum_{i=1}^n x_i \log \vartheta + \left( n - \sum_{i=1}^n x_i \right) \log(1 - \vartheta), \end{aligned}$$

so that

$$\frac{d}{d\vartheta} \log \mathcal{L}(\vartheta; \mathbf{x}) = \frac{\sum_{i=1}^n x_i}{\vartheta} - \frac{n - \sum_{i=1}^n x_i}{1 - \vartheta} = 0,$$

and therefore

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

□

**Exercise 5.21.** Consider a population such that only 40% of individuals survive after one week. After one week we check the population and we find  $k$  individuals. We want to estimate the size of the initial population.

*Proof.* The direct way to compute the number of survivals after one week knowing the parameter  $\vartheta$  is represented by a binomial distribution

$$\mathbb{P}(K = k) = \binom{N}{k} \vartheta^k (1 - \vartheta)^{N - k} = \mathcal{L}(N; k, \vartheta).$$

To find an estimate for  $N$  consider

$$\begin{aligned} \frac{\mathcal{L}(N + 1; k, \vartheta)}{\mathcal{L}(N; k, \vartheta)} &= \binom{N + 1}{k} \binom{N}{k}^{-1} \frac{\vartheta^k (1 - \vartheta)^{N - k + 1}}{\vartheta^k (1 - \vartheta)^{N - k}} \\ &= \frac{(N + 1)!}{(N + 1 - k)! k!} \frac{k! (N - k)!}{N!} (1 - \vartheta) \\ &= \frac{N + 1}{N + 1 - k} (1 - \vartheta) \end{aligned} \tag{5.10}$$

The ratio 5.10 is greater or equal than one for  $k = N, N + 1$  and until  $N \leq \frac{k}{\vartheta} - 1$ . It is less or equal than one for any  $N \geq \frac{k}{\vartheta} - 1$ . Then we can take as an estimator of  $N$  smallest integer larger than  $\frac{k}{\vartheta} - 1$ .

□

**Exercise 5.22.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim f_X(x; \vartheta)$ , where

$$f_X(x; \vartheta) = 2\vartheta x e^{-\vartheta x^2} \mathbf{1}_{\mathbb{R}_+}(x), \quad \vartheta > 0.$$

1. Find the Maximum Likelihood estimator for  $\vartheta$ .
2. Find a sufficient statistic for  $\vartheta$ .
3. Find the Moment estimator for  $\vartheta$ .
4. Compare the two estimators.

*Proof.* (1) The Likelihood Function is

$$\mathcal{L}(\vartheta; \mathbf{x}) = \prod_{i=1}^n 2\vartheta x_i e^{-\vartheta x_i^2} = (2\vartheta)^n \prod_{i=1}^n x_i \exp\left\{-\vartheta \sum_{i=1}^n x_i^2\right\},$$

and if we take the log

$$\begin{aligned} \mathcal{L}(\vartheta; \mathbf{x}) &= n \log 2\vartheta + \sum_{i=1}^n \log x_i - \vartheta \sum_{i=1}^n x_i^2 \\ &\propto n \log \vartheta - \vartheta \sum_{i=1}^n x_i^2. \end{aligned}$$

Now,

$$\frac{d}{d\vartheta} \log \mathcal{L}(\vartheta; \mathbf{x}) = n \frac{1}{\vartheta} - \sum_{i=1}^n x_i^2 = 0 \iff \hat{\vartheta}_n = \frac{n}{\sum_{i=1}^n x_i^2}.$$

(2) We can apply the Fisher Factorisation Theorem:

$$\mathcal{L}(\vartheta; \mathbf{x}) = \underbrace{\prod_{i=1}^n x_i}_{h(\mathbf{x})} \cdot \underbrace{(2\vartheta)^n \exp\left\{-\vartheta \sum_{i=1}^n x_i^2\right\}}_{g(T(\mathbf{x}), \vartheta)}$$

and conclude that  $\sum_{i=1}^n X_i^2$  is sufficient for  $\vartheta$ .

(3) For the Moment based estimator, remember that

$$m_1 = \mathbb{E}[X] = \frac{1}{n} \sum_{i=1}^n X_i.$$

Also,

$$\begin{aligned} \mu_1 &= \mathbb{E}[X] \\ &= \int_{\mathbb{R}} x f_X(x) dx \\ &= \int_0^{+\infty} 2\vartheta x^2 \exp\{-\vartheta x^2\} dx \\ &= \int_0^{+\infty} y^{\frac{1}{2}} \vartheta^{-\frac{1}{2}} e^{-y} dy, \quad y = \vartheta x^2 \end{aligned}$$

$$\begin{aligned} &= \vartheta^{-\frac{1}{2}} \int_0^{+\infty} y^{\frac{3}{2}-1} e^{-y} dy, & \Gamma(\alpha) &= \int_0^{+\infty} x^{\alpha-1} e^{-x} dx \\ &= \vartheta^{-\frac{1}{2}} \Gamma\left(\frac{3}{2}\right) \\ &= \vartheta^{-\frac{1}{2}} \Gamma\left(1 + \frac{1}{2}\right), & \Gamma(\alpha + 1) &= \alpha \Gamma(\alpha) \\ &= \frac{\vartheta^{-\frac{1}{2}}}{2} \Gamma\left(\frac{1}{2}\right) \\ &= \frac{\sqrt{\pi}}{2} \vartheta^{-\frac{1}{2}}, \end{aligned}$$

therefore

$$\frac{\sqrt{\pi}}{2} \vartheta^{-\frac{1}{2}} = \frac{1}{n} \sum_{i=1}^n X_i \implies \hat{\vartheta}_n = \frac{n^2 \pi}{4 (\sum_{i=1}^n X_i)^2}$$

and it is not a function of the sufficient statistic.

□



## 6

## Week VI

## 6.1 30/10/18

## 6.1.1 Exercises

**Exercise 6.1.** The connection of a user to a server is accepted with probability  $\vartheta \in (0, 1)$ . Assume that the trials for connecting are independent and denote by  $X_1, \dots, X_n$  the number of trials that are necessary to get the connection in  $n$  different days. Set  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ .

1. Is  $\bar{X}_n$  an unbiased estimator for  $\vartheta^{-1}$ ?
2. Study the convergence of

$$\sqrt{n} \left( \bar{X}_n - \frac{1}{\vartheta} \right), \quad X_i \in \{1, 2, \dots\}.$$

*Proof.* First of all  $(X_i) \sim \text{Geom}(\vartheta)$  if and only if

$$\mathbb{P}(X_i = k) = \vartheta(1 - \vartheta)^{k-1}, \quad k = 1, 2, \dots$$

Now,

$$\mathbb{E}[X_i] = \frac{1}{\vartheta},$$

and

$$\mathbb{E}[\bar{X}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = \frac{1}{\vartheta},$$

therefore the estimator is unbiased for  $\vartheta^{-1}$ .

(2) Note that

$$\sqrt{n} \left( \bar{X}_n - \frac{1}{\vartheta} \right) = \frac{\sum_{i=1}^n X_i - n/\vartheta}{\sqrt{\frac{1-\vartheta}{\vartheta^2}} \sqrt{n}} \cdot \sqrt{\frac{1-\vartheta}{\vartheta^2}} \xrightarrow{n \rightarrow \infty} N \left( 0, \frac{1-\vartheta}{\vartheta^2} \right),$$

therefore the estimator is asymptotically Gaussian.  $\square$

**Exercise 6.2.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim f_X(x; \vartheta)$ , where

$$f_X(x; \vartheta) = \frac{4x^2}{\vartheta^3 \sqrt{\pi}} \exp\left\{-\frac{x^2}{\vartheta^2}\right\} \mathbb{1}_{\mathbb{R}_+}(x), \quad \vartheta > 0.$$

1. Find the Maximum Likelihood estimator for  $\vartheta$ .
2. Find a sufficient statistic for  $\vartheta$ .
3. Can you find an unbiased estimator for  $\vartheta$ ?

*Proof.* (1) The Likelihood function is

$$\begin{aligned} \mathcal{L}(\vartheta; \mathbf{x}) &= \prod_{i=1}^n \frac{4x_i^2}{\vartheta^3 \sqrt{\pi}} \exp\left\{-\frac{x_i^2}{\vartheta^2}\right\} \\ &= \left(\frac{4}{\vartheta^3 \sqrt{\pi}}\right)^n \prod_{i=1}^n x_i^2 \exp\left\{-\sum_{i=1}^n \frac{x_i^2}{\vartheta^2}\right\}, \end{aligned}$$

and

$$\log \mathcal{L}(\vartheta, \mathbf{x}) = n \log\left(\frac{4}{\vartheta^3 \sqrt{\pi}}\right) + \sum_{i=1}^n \log x_i^2 - \frac{1}{\vartheta^2} \sum_{i=1}^n x_i^2.$$

Now,

$$\frac{d}{d\vartheta} \log \mathcal{L}(\vartheta; \mathbf{x}) = n \frac{\vartheta^3 \sqrt{\pi}}{4} \frac{4}{\sqrt{\pi}} (-3) \frac{1}{\vartheta^4} + \frac{2}{\vartheta^3} \sum_{i=1}^n x_i^2.$$

Therefore,

$$\frac{3n}{\vartheta^2} = \frac{1}{\vartheta^3} \sum_{i=1}^n x_i^2 \implies \hat{\vartheta}_n = \frac{1}{3n} \sum_{i=1}^n x_i^2$$

is a Maximum Likelihood estimator for  $\vartheta$ .

(2)  $\sum_{i=1}^n X_i^2$  is a sufficient statistic for  $\vartheta$ . Indeed,

$$\mathcal{L}(\vartheta; \mathbf{x}) = \left(\frac{4}{\vartheta^3 \sqrt{\pi}}\right)^n \exp\left\{\frac{1}{\vartheta} \sum_{i=1}^n x_i^2\right\} \prod_{i=1}^n x_i^2.$$

(3) To check unbiasedness, first compute

$$m_1 = \mu_1 \implies \frac{1}{n} \sum_{i=1}^n X_i = \mathbb{E}[X].$$

$$\mathbb{E}[\hat{\vartheta}_n] = \mathbb{E}\left[\frac{1}{3n} \sum_{i=1}^n X_i^2\right] = \frac{1}{3n} \sum_{i=1}^n \mathbb{E}[X_i^2] = \frac{1}{3} \mathbb{E}[X^2].$$

The second moment of  $X$  is

$$\begin{aligned} \mathbb{E}[X^2] &= \int_{\mathbb{R}} x^2 f_X(x) dx \\ &= \int_0^{+\infty} x^2 \frac{4x^2}{\vartheta^3 \sqrt{\pi}} \exp\left\{-\frac{x^2}{\vartheta^2}\right\} dx \end{aligned}$$

$$\begin{aligned}
&= \int_0^{+\infty} \frac{4x^4}{\sqrt{\pi}\vartheta^2} \exp\left\{-\frac{x^2}{\vartheta^2}\right\} dx \\
&= \int_0^{+\infty} \frac{2x}{\vartheta^2} \cdot \frac{2x^3}{\sqrt{\pi}\vartheta} \exp\left\{-\frac{x^2}{\vartheta^2}\right\} dx, \quad t = \frac{x^2}{\vartheta^2} \\
&= \int_0^{+\infty} \frac{1}{\sqrt{\pi}} t^{\frac{3}{2}} \frac{\vartheta^3}{\vartheta} e^{-t} dt \\
&= \frac{2}{\sqrt{\pi}} \vartheta^2 \int_0^{+\infty} t^{\frac{3}{2}-1} e^{-t} dt \\
&= \frac{2}{\sqrt{\pi}} \vartheta^2 \Gamma\left(\frac{5}{2}\right) \\
&= \frac{2}{\sqrt{\pi}} \vartheta^2 \Gamma\left(1 + \frac{3}{2}\right) \\
&= \frac{2}{\sqrt{\pi}} \vartheta^2 \Gamma\left(\frac{3}{2}\right) \frac{3}{2} \\
&= \frac{2}{\sqrt{\pi}} \vartheta^2 \frac{\sqrt{\pi}}{2} \frac{3}{2} \\
&= \frac{3}{2} \vartheta^2,
\end{aligned}$$

therefore

$$\mathbb{E}[\hat{\vartheta}_n] = \frac{\vartheta^2}{2}.$$

□

**Exercise 6.3.** A server gets requests of jobs from users. Every job takes an exponential time with parameter  $\lambda$ , which is known. Assuming that the server has been active for a time  $X = x$ , find the Maximum Likelihood estimator of the number  $n$  of requests.

*Proof.* Since we are dealing with sums of exponentials, we have that  $X \sim \Gamma(n, \lambda)$ . Easily, we can compute

$$\mathbb{E}[X] = \frac{n}{\lambda} \iff \hat{n} = \lfloor x\lambda \rfloor,$$

which could be an estimator. As we did last time, consider

$$f_n(x) = \frac{\lambda^n}{\Gamma(n)} x^{n-1} e^{-\lambda x}$$

and compute

$$\frac{f_{n+1}(x)}{f_n(x)} = \frac{\lambda x}{n}.$$

Observe that the mapping  $n \mapsto \frac{f_{n+1}(x)}{f_n(x)}$  is increasing until  $\lambda x \geq n$  and then decreasing. Therefore, we have indeed that

$$\hat{n} = \lfloor x\lambda \rfloor.$$

□

**Exercise 6.4.** For estimating the number of individuals in a population we can use the following method: we take a first sample of  $k$  individuals, all distinct, and we mark them. After that, we insert them again in the population. We take another sample of size  $k$  (all distinct). Let  $N$  be the number of the population.

1. What is the probability that in the second group there are  $r$  individuals from the first group?
2. Find an estimator for  $N$ .

*Proof.* (1) The distribution in question is clearly an Hypergeometric:

$$\mathbb{P}(X = r) = \frac{\binom{K}{r} \binom{N-k}{k-r}}{\binom{N}{k}}.$$

(2) For an estimator of  $N$ , consider

$$\begin{aligned} \mathbb{P}(X = r) &= g(N) \\ &= c \cdot \frac{(N-k)!}{\frac{N!}{(N-k)!} (N-2k+r)!} \\ &= c \cdot \frac{[(N-k)!]^2}{N!(N-2k+r)!} \end{aligned}$$

where  $c$  is a constant that does not depend on  $N$ . Finally, we consider the the usual ratio

$$\frac{g(N+1)}{g(N)} = \dots = \frac{k^2 - r}{r}.$$

□

**Exercise 6.5.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim U([a-b], [b+a])$ .

1. Find the Moment based estimator of  $a$  and  $b$ .
2. Find the Maximum Likelihood estimator of  $a$  and  $b$ .
3. Compare the estimators.

*Proof.* (1) Compute

$$\begin{aligned} \mathbb{E}[X] &= \int_{\mathbb{R}} \frac{x}{a+b-(a-b)} dx \\ &= \frac{1}{2b} \frac{x^2}{2} \Big|_{a-b}^{a+b} \\ &= \frac{1}{4b} (a^2 + b^2 + 2ab - a^2 - b^2 + 2ab) \\ &= a, \\ \mathbb{E}[X^2] &= \int_{a-b}^{a+b} \frac{x^2}{2b} dx \\ &= \frac{x^3}{6b} \Big|_{a-b}^{a+b} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{6b}(a^3 + b^3 + 3a^2b + 3ab^2 - (a^3 - b^3 - 3a^2b + 3ab^2)) \\
&= \frac{1}{6b}(2b^3 + 6ba^2b) \\
&= a^2 + \frac{b^2}{3}
\end{aligned}$$

so that

$$\mathbb{E}[X] = a = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{and} \quad \mathbb{E}[X^2] = a^2 + \frac{b^2}{3} = \frac{1}{n} \sum_{i=1}^n X_i^2.$$

Therefore

$$\bar{T}_n(a) = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{and} \quad \bar{T}_n(b) = \left[ \frac{3}{n} \sum_{i=1}^n X_i^2 - 3 \left( \frac{1}{n} \sum_{i=1}^n X_i \right)^2 \right]^{\frac{1}{2}}.$$

(2) For the Maximum Likelihood estimator, first compute

$$\mathcal{L}(a, b; \mathbf{x}) = \left( \frac{1}{2b} \right)^n \prod_{i=1}^n \mathbb{1}_{[a-b, b+a]}(x_i), \quad x_i.$$

We have that

$$x_i \in [a-b, b+a] \iff a-b \leq \min_i x_i \wedge \max_i x_i \leq a+b. \quad (6.1)$$

We want  $b$  to be as small as possible, because it is at the denominator: indeed, we are looking for the smallest  $b$  such that (6.1) is satisfied for some  $a$ :

$$b \geq \frac{\max_i x_i - \min_i x_i}{2} \implies b = \frac{\max_i x_i - \min_i x_i}{2},$$

where (6.1) becomes an equality if

$$a = \frac{\max_i x_i + \min_i x_i}{2}.$$

Therefore the Maximum Likelihood estimators for  $a$  and  $b$  are

$$\begin{aligned}
\hat{T}_n(a) &= \frac{\max_i x_i + \min_i x_i}{2} \\
\hat{T}_n(b) &= \frac{\max_i x_i - \min_i x_i}{2}
\end{aligned}$$

and

$$\mathcal{L}(a, b; \mathbf{x}) = \left( \frac{1}{2b} \right)^n \mathbb{1}_A(x_i),$$

where  $A := a-b \leq \min_i x_i \wedge \max_i x_i \leq a+b$ .

□

## 6.2 31/10/18

## 6.2.1 Exercises

**Exercise 6.6.** Suppose you want to achieve knowledge about a private information. Instead of asking directly, a commonly used technique is to ask the sample to flip a coin. Each person would then check either 0 or 1. Then the interviewer may ask whether the answer they should have given to a direct question corresponds to the result of the coin flipped.

Consider the experiment of flipping a the coin:

YES with probability  $p$ , NO with probability  $1 - p$ ,

with  $p$  known. We would like to estimate the proportion  $\vartheta$  of population that have a feature  $A$  without asking directly, but through this randomisation process.

*Proof.* For example, consider an individual  $i$ : he/she either has or does not have the feature  $A$ . Assume he/she has the feature  $A$  and flips the coin. With probability  $p$  he obtains heads. Therefore, the answer to the direct and to the randomised question coincide.

Now, define an indicator random variable for the two answers being the same:

$$X_i = \begin{cases} 1 & \text{if YES} \\ 0 & \text{if NO} \end{cases}.$$

Therefore,

$$\begin{aligned} \mathbb{P}(X_i = 1) &= \mathbb{P}(X_i = 1 | A) \mathbb{P}(A) + \mathbb{P}(X_i = 1 | \neg A) \mathbb{P}(\neg A) \\ &= p\vartheta + (1-p)(1-\vartheta) \\ &= p\vartheta + 1 - p - \vartheta + p\vartheta \\ &= (1-p) - \vartheta(1-2p). \end{aligned}$$

Notice that  $\mathbf{X} = (X_1, \dots, X_n)$  is a random sample form  $X \sim Ber((1-p) - \vartheta(1-2p))$ . If we would like to estimate  $\vartheta$  using the Likelihood function, then set

$$\hat{T} = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i,$$

and therefore

$$\hat{\vartheta} = \frac{\frac{1}{n} \sum_{i=1}^n X_i - (1-p)}{2p-1}, \quad p \neq \frac{1}{2}.$$

Observe that it is unbiased, since

$$\mathbb{E}[\hat{\vartheta}] = \frac{\frac{1}{n} n(\vartheta p + (1-p)(1-\vartheta)) - (1-p)}{2p-1} = \vartheta$$

it is consistent. Also, thanks to the Weak Law of Large Numbers,

$$\text{Var}(\hat{\vartheta}) = \frac{\frac{1}{n} \sum_{i=1}^n X_i - (1-p)}{2p-1} \xrightarrow{n \rightarrow \infty} \vartheta.$$

□

**Exercise 6.7.** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$ , where  $X_i$  is the survival time of a lamp, meaning  $X_i \sim \mathcal{E}_n(\vartheta)$ . We already know that one possible unbiased estimator for  $\vartheta$  is

$$\frac{n-1}{\sum_{i=1}^n X_i}.$$

Fix a time  $T$  and denote by  $Z$  the number of lamps survived at the time  $T$ .

Find the Maximum Likelihood estimator for  $\vartheta$  based on the information about  $Z$ .

*Proof.* Being  $Z \sim \mathcal{B}(n, p)$ , the probability that  $z$  lamps are still working at time  $T$  is

$$\mathbb{P}(Z = z) = \binom{n}{z} p^z (1-p)^{n-z}, \quad p = \mathbb{P}(X_i \geq T) = e^{-\vartheta T},$$

where instead  $p$  is the probability that a lamp is still working after time  $T$ , being  $X_i \sim \mathcal{E}_n(p)$ . Now, combining our knowledge of a Maximum Likelihood estimator for the binomial coefficient

$$\hat{p} = \frac{z}{n}$$

and

$$e^{-\vartheta T} = \frac{z}{n},$$

we have that

$$\hat{\vartheta} = -\frac{\log\left(\frac{z}{n}\right)}{T}.$$

We have faced problems of maximising the Log-Likelihood function with respect to  $\vartheta$ . This is particularly difficult when

- the parameter of interest is in  $\mathbb{N}$ ;
- the parameter of interest defines the support of the model, for example the Maximum Likelihood estimator function of  $\max\{a, b\}$  or  $\min\{a, b\}$ ;
- Laplace model, where

$$f_X(x; \vartheta) \propto \exp\{-|x - \vartheta|\};$$

- when the model is not regular or the Likelihood equation is unsolvable. Suppose for example to have a regular model  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim f_X(x; a, b)$ , meaning that we can differentiate with respect to both  $a$  and  $b$ . Also, suppose that we cannot solve the Likelihood equation explicitly. For that matter, consider a  $\Gamma(\alpha, \beta)$ :

$$\mathcal{L}(\alpha, \beta; \mathbf{x}) = \prod_{i=1}^n \frac{\beta^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-\beta x_i} = \left(\frac{\beta^\alpha}{\Gamma(\alpha)}\right)^n \left(\prod_{i=1}^n x_i\right)^{\alpha-1} \exp\left\{-\beta \sum_{i=1}^n x_i\right\}$$

and

$$\log \mathcal{L}(\alpha, \beta; \mathbf{x}) = n\alpha \log \beta - n \log \Gamma(\alpha) + (\alpha - 1) \sum_{i=1}^n \log x_i - \beta \sum_{i=1}^n x_i.$$

Taking the derivatives with respect to  $\alpha$ ,

$$n \log \beta - \frac{n\Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_{i=1}^n \log x_i = 0$$

74 6 Week VI

and  $\beta$

$$n\alpha \frac{1}{\beta} - \sum_{i=1}^n x_i = 0.$$

These are numerically approachable thanks to the Newton-Rapson method.

---

## Week VII

### 7.1 06/11/18

#### 7.1.1 Properties of the Maximum Likelihood Approach

As we have stated before, concerning the Maximum Likelihood estimator, we would like to solve the equation (5.6). If the aforementioned cases are not relevant for the treated situation, we must implement different approaches. For example,

1. when the parameter defines the support of the model;
2. when the model is *regular*, but we cannot solve explicitly the Maximum Likelihood equation

$$V'_n(\vartheta) = 0.$$

If we were in this case, we could follow two different approaches.

- a) The *substitution approach*, if we had some idea of the value we are looking for.
- b) **Newton-Rapson Method.**

Let  $\mathcal{L}(\vartheta; \mathbf{x})$  be the Likelihood function, with  $\vartheta \in \Theta$ . First of all, consider the *gradient*

$$g(\vartheta) = V'_n(\vartheta)$$

and the *hessian*

$$H(\vartheta) = V''_n(\vartheta).$$

Then, fix a value  $\vartheta_0 \in \Theta$  and consider the Taylor expansion for  $g(\vartheta)$  around  $\vartheta_0$ ,

$$g(\vartheta) = g(\vartheta_0) + (\vartheta - \vartheta_0)H(\vartheta_0) + \dots$$

Next, replace  $\vartheta$  with the Maximum Likelihood estimator  $T_n$  for  $\vartheta$ :

$$g(T_n) = g(\vartheta_0) + H(\vartheta_0)(T_n - \vartheta_0) + \dots$$

This must be null, since the Maximum Likelihood estimator maximises the Likelihood function. Now,

$$T_n - \vartheta_0 = -\frac{g(\vartheta_0)}{H(\vartheta_0)} + \dots$$

Therefore, we can write an approximate solution for the Maximum Likelihood estimator of  $\vartheta$ , meaning

$$T_n = \vartheta_0 - \frac{g(\vartheta_0)}{H(\vartheta_0)^{-1}} \implies T_n = \vartheta_0 - \frac{V'_n(\vartheta)}{V''_n(\vartheta_0)^{-1}}.$$

Finally, write down the iteration

$$t^{(i)} = t^{(i-1)} - \frac{V'_n(t^{(i-1)})}{V''_n(t^{(i-1)})}, \quad \forall i \in \mathbb{N}_+.$$

Note that the iterative procedure naturally stops when  $t^{(i)} = t_{(i+1)}$ , that is when  $V'_n(\vartheta) = 0$ . However, there are two problems to solve. First,  $t^{(0)}$  must be specified, and second  $[V''_n(\cdot)]^{-1}$  must exist. How do we solve the existence of that inverse function? We can use a variation of the Newton-Rapson method, which is called **Score algorithm**. This consists in replacing  $V'_n(t^{(i)})$  with  $\mathbb{E}[V''_n(t^{(i)})]$ , which is always positive definite.

*Example 7.1.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X$ , where  $X$  is distributed as a zero-truncated Poisson distribution with parameter  $\vartheta$ . Now, for the generic Poisson,

$$\mathbb{P}(Y = x) = \frac{e^{-\vartheta} \vartheta^x}{x!} \mathbb{1}_{\mathbb{N}}(x),$$

where  $e^{-\vartheta}$  is the constant of normalisation. For the truncated Poisson, we would like to change the constant in order to “truncate” the Poisson at 0:

$$\mathbb{P}(X = x) = \frac{\vartheta^x e^{-\vartheta}}{x!(1 - e^{-\vartheta})} \mathbb{1}_{\mathbb{N}_+}(x),$$

since

$$\sum_{i=0}^{\infty} \frac{e^{-\vartheta} \vartheta^i}{i!} = 1 \implies \sum_{i=1}^{\infty} \frac{e^{-\vartheta} \vartheta^i}{i!} = 1 - e^{-\vartheta}.$$

Now,

$$\mathcal{L}(\vartheta, x_i) = \prod_{i=1}^n \frac{\vartheta^{x_i} e^{-\vartheta}}{x_i!(1 - e^{-\vartheta})} \mathbb{1}_{\mathbb{N}_+}(x_i) = \frac{e^{-n\vartheta} \vartheta^{\sum_{i=1}^n x_i}}{(\prod_{i=1}^n x_i!) (1 - e^{-\vartheta})^n}$$

and

$$\log \mathcal{L}(\vartheta, x_i) \propto \log \vartheta \sum_{i=1}^n x_i - n\vartheta - n \log(1 - e^{-\vartheta});$$

therefore

$$\frac{d}{d\vartheta} \log \mathcal{L}(\vartheta, x_i) = -n + \frac{\sum_{i=1}^n x_i}{\vartheta} - \frac{ne^{-\vartheta}}{1 - e^{-\vartheta}} = \frac{n\bar{X}_n}{\vartheta} - \frac{n}{1 - e^{-\vartheta}}$$

and we have

$$\frac{n\bar{X}_n}{\vartheta} - \frac{n}{1 - e^{-\vartheta}} = 0 \implies T_n = \bar{X}_n (1 - e^{-T_n}) \implies t^{(i)} = \bar{X}_n \left(1 - e^{-t^{(i-1)}}\right),$$

for  $i = 1, 2, \dots$ . If we wanted to apply the *Newton-Rapson* method, first we compute

$$V_n'(\vartheta) = \frac{n\bar{X}_n}{\vartheta} - \frac{n}{1 - e^{-\vartheta}}, \quad V_n''(\vartheta) = -\frac{n\bar{X}_n}{\vartheta^2} + \frac{ne^{-\vartheta}}{(1 - e^{-\vartheta})^2},$$

then simply calculate

$$t^{(i)} = t^{(i-1)} - \underbrace{\left[ \frac{n\bar{X}_n}{t^{(i-1)}} - \frac{n}{1 - e^{-t^{(i-1)}}} \right]}_{V_n'(t^{(i-1)})} \cdot \underbrace{\left[ -\frac{n\bar{X}_n}{(t^{(i-1)})^2} + \frac{ne^{-t^{(i-1)}}}{(1 - e^{-t^{(i-1)}})^2} \right]^{-1}}_{(V_n''(t^{(i-1)}))}$$

### 7.1.2 Delta Method (a generalisation of the Central Limit Theorem for functionals)

Consider  $X_1, \dots, X_n$  independent Bernoulli random variables, with parameter  $p$ . The typical parameter of interest is  $p$ , but sometimes we also want to look out for

$$\frac{p}{1-p}. \quad (7.1)$$

Since a valid estimator for  $p$  is  $\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i$ , one may consider

$$\frac{\hat{p}}{1-\hat{p}}$$

as an estimator for  $\frac{p}{1-p}$ . What are the properties of this estimator? Can we say something about his variance? And what about its distribution? In order to answer to this questions, we could apply the *Delta method*.

Let  $T_1, \dots, T_n$  be random variables, with means  $\vartheta_1, \dots, \vartheta_n$  respectively, and consider

$$\mathbf{T} = (T_1, \dots, T_n), \quad \boldsymbol{\vartheta} = (\vartheta_1, \dots, \vartheta_n).$$

Suppose there is a differentiable function  $g(\mathbf{T})$  for which we want to approximate *mean* and *variance*. Define

$$g_i'(\boldsymbol{\vartheta}) := \left. \frac{d}{dt_i} g(\mathbf{t}) \right|_{t_1=\vartheta_1, \dots, t_n=\vartheta_n}$$

and consider the first order Taylor expansion

$$g(\mathbf{t}) = g(\boldsymbol{\vartheta}) + \sum_{i=1}^n g_i'(\boldsymbol{\vartheta})(t_i - \vartheta_i) + R.$$

Therefore,

$$g(\mathbf{t}) \approx g(\boldsymbol{\vartheta}) + \sum_{i=1}^n g_i'(\boldsymbol{\vartheta})(t_i - \vartheta_i)$$

$$g(\mathbf{T}) \approx g(\boldsymbol{\vartheta}) + \sum_{i=1}^n g_i'(\boldsymbol{\vartheta})(T_i - \vartheta_i)$$

78 7 Week VII

$$\mathbb{E}[g(\mathbf{T})] \approx g(\boldsymbol{\vartheta}) + \sum_{i=1}^n g'_i(\boldsymbol{\vartheta}) \mathbb{E}[T_i - \vartheta_i].$$

By similar calculation, we can obtain an approximation of the variance of  $g(\mathbf{T})$ :

$$\text{Var}(g(\mathbf{T})) \approx \sum_{i=1}^n (g'_i(\boldsymbol{\vartheta}))^2 \text{Var}(T_i) + 2 \sum_{i>j} g'_i(\boldsymbol{\vartheta}) g'_j(\boldsymbol{\vartheta}) \text{Cov}(T_i, T_j).$$

Then, by using the Taylor approximation of the mean and variance of  $g(\mathbf{T})$ , we can state the following generalisation of the Central Limit Theorem, which is known as *Delta method*.

**Theorem 7.2 (Delta Method).** *Let  $(Y_n)_{n \geq 1}$  be a sequence of random variables such that*

$$\sqrt{n}(Y_n - \vartheta) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, \sigma^2).$$

*For a given function  $g(\cdot)$  and a specific value of  $\vartheta$ , assume that  $g'(\vartheta)$  exists and it is not zero. Then,*

$$\sqrt{n}(g(Y_n) - g(\vartheta)) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, \sigma^2 (g'(\vartheta))^2).$$

*Proof.* Consider the Taylor expansion of  $g(Y_n)$  around  $\vartheta$ :

$$g(Y_n) = g(\vartheta) + g'(\vartheta)(Y_n - \vartheta) + o(Y_n - \vartheta).$$

Now,

$$\sqrt{n}(g(Y_n) - g(\vartheta)) = \sqrt{n}g'(\vartheta)(Y_n - \vartheta) + \sqrt{n}o(Y_n - \vartheta) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, \sigma^2 (g'(\vartheta))^2),$$

since

$$\sqrt{n}o(Y_n - \vartheta) = \sqrt{n}(Y_n - \vartheta) \cdot \frac{o(Y_n - \vartheta)}{(Y_n - \vartheta)} \xrightarrow{n \rightarrow \infty} 0.$$

□

## 7.2 07/11/18

### 7.2.1 Maximum Likelihood Approach

A useful property of the Maximum Likelihood estimator is the *invariance property*. Suppose we have a parametric family indexed by a parameter  $\vartheta$ , but we are interested in estimating  $\tau(\vartheta)$ . The invariance property of Maximum Likelihood estimators says that, if  $\hat{\vartheta}$  is the Maximum Likelihood estimator for  $\vartheta$ , then  $\tau(\hat{\vartheta})$  is the Maximum Likelihood estimator for  $\tau(\vartheta)$ .

In general, the invariance property of the Maximum Likelihood estimator is true for any functional of  $\vartheta$ . To show this for an arbitrary function  $\tau(\cdot)$ , we define

$$\mathcal{L}^*(\eta; \mathbf{x}) := \sup_{\{\vartheta \in \Theta \mid \tau(\vartheta) = \eta\}} \mathcal{L}(\vartheta; \mathbf{x}). \quad (7.2)$$

The value  $\hat{\eta}$  that maximises (7.2) is called the **Maximum Likelihood estimator** of  $\eta = \tau(\vartheta)$ , and is defined as

$$\mathcal{L}^*(\hat{\eta}; \mathbf{x}) = \sup_{\eta} \mathcal{L}^*(\eta; \mathbf{x}).$$

**Theorem 7.3 (Invariance Property of Maximum Likelihood estimators).** *If  $\hat{\vartheta}$  is the Maximum Likelihood estimator for  $\vartheta$ , then, for every other function  $\tau(\vartheta)$ , the Maximum Likelihood estimator of  $\tau(\vartheta)$  is  $\tau(\hat{\vartheta})$ .*

*Proof.* Let  $\hat{\eta}$  be the value that maximises the function  $\mathcal{L}^*(\eta; \mathbf{x})$ . We must show that

$$\mathcal{L}^*(\tau(\hat{\vartheta}); \mathbf{x}) = \sup_{\eta} \mathcal{L}^*(\eta; \mathbf{x}).$$

Therefore, we have that

$$\begin{aligned} \mathcal{L}^*(\hat{\eta}; \mathbf{x}) &= \sup_{\eta} \left( \sup_{\{\vartheta \in \Theta \mid \tau(\vartheta) = \eta\}} \mathcal{L}(\vartheta; \mathbf{x}) \right) \\ &= \sup_{\vartheta} \mathcal{L}(\vartheta; \mathbf{x}) \\ &= \mathcal{L}(\hat{\vartheta}; \mathbf{x}) \\ &= \sup_{\{\vartheta \in \Theta \mid \tau(\vartheta) = \tau(\hat{\vartheta})\}} \mathcal{L}(\vartheta; \mathbf{x}) \\ &= \mathcal{L}^*(\tau(\hat{\vartheta}); \mathbf{x}). \end{aligned}$$

Therefore,  $\tau(\hat{\vartheta})$  is the Maximum Likelihood estimator of  $\tau(\vartheta)$ , where  $\hat{\vartheta}$  is the Maximum Likelihood estimator of  $\vartheta$ .  $\square$

*Remark 7.4.* Theorem 7.3 holds also in the context of *multidimensional* parametric spaces.

*Example 7.5.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \mathcal{N}(\mu, \sigma^2)$ . Find the Maximum Likelihood estimator for  $(\mu, \sigma^2)$ .

80 7 Week VII

*Proof.* We know that

$$\begin{aligned}\mathcal{L}(\mu, \sigma^2; \mathbf{x}) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_i - \mu)^2}{2\sigma^2}\right\} \\ &= \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right\},\end{aligned}$$

therefore

$$\begin{aligned}\log \mathcal{L}(\mu, \sigma^2; \mathbf{x}) &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \\ &\propto -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2.\end{aligned}$$

Next, we impose

$$\begin{aligned}\frac{d}{d\mu} \log \mathcal{L}(\mu, \sigma^2; \mathbf{x}) &= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0 \\ \frac{d}{d\sigma^2} \log \mathcal{L}(\mu, \sigma^2; \mathbf{x}) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0\end{aligned}$$

to obtain

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2.$$

We want to verify that  $(\hat{\mu}, \hat{\sigma}^2)$  is a maximum. For every  $\mu \neq \bar{x}_n$ , we have

$$\sum_{i=1}^n (x_i - \hat{\mu})^2 > \sum_{i=1}^n (x_i - \bar{x}_n)^2 :$$

indeed,

$$\begin{aligned}\sum_{i=1}^n (x_i - \hat{\mu})^2 &= \sum_{i=1}^n [(x_i - \bar{x}_n) + (\bar{x}_n - \hat{\mu})]^2 \\ &= \sum_{i=1}^n (x_i - \bar{x}_n)^2 + \sum_{i=1}^n (\bar{x}_n - \hat{\mu})^2 + 2 \underbrace{\sum_{i=1}^n (x_i - \bar{x}_n)(\bar{x}_n - \hat{\mu})}_{=0} \\ &= \sum_{i=1}^n (x_i - \bar{x}_n)^2 + \sum_{i=1}^n (\bar{x}_n - \hat{\mu})^2 \\ &> \sum_{i=1}^n (x_i - \bar{x}_n)^2.\end{aligned}$$

By using this inequality, we have

$$\left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right\} \geq \left(\frac{1}{2\pi\hat{\sigma}^2}\right)^{\frac{n}{2}} \exp\left\{-\frac{1}{2\hat{\sigma}^2} \sum_{i=1}^n (x_i - \hat{\mu})^2\right\},$$

therefore, the probability of verifying that  $(\hat{\mu}, \hat{\sigma}^2)$  is a maximum reduces to a one-dimensional problem: verifying that

$$(\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right\}$$

achieves its global maximum at

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

Indeed,

$$\begin{aligned} 0 &= \frac{d}{d\sigma^2} \left( (\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right\} \right) \\ &= -\frac{n}{2} (\sigma^2)^{-\frac{n}{2}-1} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right\} + \\ &\quad + (\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right\} \left( -\frac{(-1)}{2\sigma^4} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \right) \\ &= \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right\} \left( -\frac{n}{2} (\sigma^2)^{-\frac{n+2}{2}} + (\sigma^2)^{-\frac{n}{2}} \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \right) \end{aligned}$$

if and only if

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

□

From a computational level, the invariance property is a good one, but it is also the reason why the Maximum Likelihood estimator can be biased for a finite sample size.

**Proposition 7.6.** *Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from a parametric model  $X \sim f_X(x; \vartheta)$ . If the Maximum Likelihood estimator  $T_n$  for  $\vartheta$  exists and is unique, then  $T_n$  is a **function of the sufficient statistic** for  $\vartheta$ .*

Indeed, if a sufficient statistic for  $\vartheta$  exists, then the Likelihood function  $\mathcal{L}(\vartheta, \mathbf{x})$  can be factorised as follows:

$$\mathcal{L}(\vartheta, \mathbf{x}) = g(T(\mathbf{x}), \vartheta)h(\mathbf{x}),$$

and therefore

$$\log \mathcal{L}(\vartheta, \mathbf{x}) = \log(g(T(\mathbf{x}), \vartheta)) + \log h(\mathbf{x}) \propto \log(g(T(\mathbf{x}), \vartheta)).$$

### 7.2.2 Asymptotic Properties of the Maximum Likelihood Estimator

**Theorem 7.7.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from a parametric model  $X \sim f_X(x; \vartheta)$  and consider the Likelihood function  $\mathcal{L}(\vartheta; \mathbf{x})$ . Let  $\hat{\vartheta}_n$  be the Maximum Likelihood estimator for  $\vartheta$  and let  $\tau(\cdot)$  be a continuous function of  $\vartheta$ . Under the usual assumptions of regularity, if, for every  $\varepsilon > 0$  and every  $\vartheta \in \Theta$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left( \left| \tau(\hat{\vartheta}_n) - \tau(\vartheta) \right| \geq \varepsilon \right) = 0,$$

then the Maximum Likelihood estimators are consistent in probability.

**Theorem 7.8.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from a parametric model  $X \sim f_X(x; \vartheta)$  and consider the Likelihood function  $\mathcal{L}(\vartheta; \mathbf{x})$ . Let  $\hat{\vartheta}_n$  be the Maximum Likelihood estimator for  $\vartheta$  and let  $\tau(\cdot)$  be a continuous function of  $\vartheta$ . Then,

$$\sqrt{n}[\tau(\hat{\vartheta}_n) - \tau(\vartheta)] \xrightarrow{n \rightarrow \infty} \mathcal{N} \left( 0, \frac{1}{I_1(\vartheta)} \right),$$

where  $I_1(\vartheta)$  is the Fisher information of the parametric model, then the Maximum Likelihood estimators are

- consistent (weakly);
- asymptotically unbiased;
- asymptotically Gaussian;
- asymptotically efficient.

*Proof.* We will prove the result for  $\vartheta$ . The proof for  $\tau(\vartheta)$  is simply an application of the Delta method.

We know that

$$\mathcal{L}(\vartheta; \mathbf{x}) = \sum_{i=1}^n \log f_{X_i}(x_i; \vartheta).$$

Denote by  $\mathcal{L}'$  and  $\mathcal{L}''$  the first and second derivatives of  $\mathcal{L}$  with respect to  $\vartheta$ , respectively. Expand the first derivative around  $\vartheta_0$

$$\mathcal{L}'(\vartheta; \mathbf{x}) = \mathcal{L}'(\vartheta_0; \mathbf{x}) + (\vartheta - \vartheta_0)\mathcal{L}''(\vartheta_0; \mathbf{x}) + \dots$$

and then substitute  $\vartheta$  with its Maximum Likelihood estimator:

$$0 = \mathcal{L}'(\vartheta_0; \mathbf{x}) + (\hat{\vartheta}_n - \vartheta_0)\mathcal{L}''(\vartheta_0; \mathbf{x}) + \dots$$

Therefore, we would like to study the convergence of

$$\sqrt{n}(\hat{\vartheta}_n - \vartheta_0) = -\sqrt{n} \frac{\mathcal{L}'(\vartheta_0; \mathbf{x})}{\mathcal{L}''(\vartheta_0; \mathbf{x})} = -\frac{\frac{1}{\sqrt{n}}\mathcal{L}'(\vartheta_0; \mathbf{x})}{\frac{1}{n}\mathcal{L}''(\vartheta_0; \mathbf{x})}.$$

Now, we can apply the Central Limit Theorem to the numerator and the Weak Law of Large Numbers to the denominator. Indeed,

$$\frac{1}{\sqrt{n}}\mathcal{L}'(\vartheta_0; \mathbf{x}) = \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n w_i \right),$$

where

$$w_i = \frac{\frac{d}{d\vartheta} f_{X_i}(x_i; \vartheta)}{f_{X_i}(x_i; \vartheta)}$$

with  $\mathbb{E}[w_i] = 0$  and  $\text{Var}(w_i) = I_n(\vartheta_0) = nI_1(\vartheta_0)$ . By the Central Limit Theorem,

$$\frac{1}{\sqrt{n}} \mathcal{L}'(\vartheta_0; \mathbf{x}) \longrightarrow \mathcal{N}(0, I_1(\vartheta_0)).$$

For the denominator,

$$-\frac{1}{n} \mathcal{L}''(\vartheta_0; \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n w_i^2 - \frac{1}{n} \sum_{i=1}^n \frac{\frac{d^2}{d\vartheta^2} f_{X_i}(x_i; \vartheta)}{f_{X_i}(x_i; \vartheta)},$$

with

$$\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n w_i^2 \right] = I_1(\vartheta_0)$$

and

$$\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \frac{\frac{d^2}{d\vartheta^2} f_{X_i}(x_i; \vartheta)}{f_{X_i}(x_i; \vartheta)} \right] = 0.$$

Therefore, by the Weak Law of Large Numbers,

$$-\frac{1}{n} \mathcal{L}''(\vartheta_0; \mathbf{x}) \longrightarrow I_1(\vartheta_0).$$

In conclusion, we have proved that

$$\sqrt{n}(\hat{\vartheta}_n - \vartheta) \longrightarrow \frac{X}{I_1(\vartheta_0)},$$

where  $X \sim \mathcal{N}(0, I_1(\vartheta_0))$ , meaning that

$$\sqrt{n}(\hat{\vartheta}_n - \vartheta) \longrightarrow \mathcal{N} \left( 0, \frac{1}{I_1(\vartheta_0)} \right).$$

□

**7.3 08/11/18**

*Example 7.9.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \mathcal{N}(\vartheta, \sigma^2)$ . We already showed that, by the Factorisation Theorem,

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

is a sufficient statistic for  $\vartheta$ . Is it also a minimal statistic?

*Proof.* Compute

$$\begin{aligned} \frac{\mathcal{L}(\vartheta, \mathbf{x})}{\mathcal{L}(\vartheta, \mathbf{y})} &= \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2 \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2 + \frac{2\mu}{2\sigma^2} \sum_{i=1}^n x_i - \frac{1}{2\sigma^2} n\mu^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2 - \frac{2\mu}{2\sigma^2} \sum_{i=1}^n y_i + \frac{1}{2\sigma^2} n\mu^2 \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2 + \frac{\mu}{\sigma^2} n\bar{x}_n - \frac{1}{2\sigma^2} n\mu^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2 - \frac{\mu}{\sigma^2} n\bar{y}_n + \frac{1}{2\sigma^2} n\mu^2 \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n [(x_i - \bar{x}_n)^2 - (y_i - \bar{y}_n)^2] \right\} \exp \{ n(\bar{x}_n - \bar{y}_n)(-2\mu) + (\bar{x}_n^2 - \bar{y}_n^2) \}, \end{aligned}$$

where

$$\begin{aligned} \sum_{i=1}^n (x_i - \mu)^2 &= \sum_{i=1}^n (x_i - \bar{x}_n)^2 + n(\bar{x}_n - \mu)^2 \\ \sum_{i=1}^n (y_i - \mu)^2 &= \sum_{i=1}^n (y_i - \bar{y}_n)^2 + n(\bar{y}_n - \mu)^2. \end{aligned}$$

If  $\bar{x}_n = \bar{y}_n$ , then

$$\frac{\mathcal{L}(\vartheta; \mathbf{x})}{\mathcal{L}(\vartheta; \mathbf{y})}$$

does not depend on  $\mu$ . Moreover,

$$\sum_{i=1}^n (x_i - \bar{x}_n)^2 = \sum_{i=1}^n (y_i - \bar{y}_n)^2,$$

therefore

$$\frac{\mathcal{L}(\vartheta; \mathbf{x})}{\mathcal{L}(\vartheta; \mathbf{y})}$$

does not depend on  $\sigma^2$  either. This means that

$$\left( \bar{X}_n; \sum_{i=1}^n (X_i - \bar{X}_n)^2 \right)$$

is a minimal sufficient statistic for  $\vartheta$ , otherwise

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

is a minimal sufficient statistic for  $\vartheta$ , since it is a bijection of it. Being minimal sufficient, we have the best possible predictor.  $\square$

*Example 7.10.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim f_X(x; \vartheta)$ , where

$$f_X(\mathbf{x}; \vartheta) = \frac{1}{\pi} \cdot \frac{1}{1 + (\mathbf{x} - \vartheta)^2},$$

a Cauchy distribution. Now,

$$\mathcal{L}(\vartheta; \mathbf{x}) = \prod_{i=1}^n \frac{1}{\pi} \frac{1}{1 + (x_i - \vartheta)^2} = \underbrace{\left( \frac{1}{\pi} \right)^n}_{h(\mathbf{x})} \cdot \underbrace{\prod_{i=1}^n \frac{1}{1 + (x_i - \vartheta)^2}}_{g(T(\mathbf{x}), \vartheta)},$$

meaning that the random sample is the sufficient statistic, according to the Factorisation Theorem.



## 8

## Week VIII

## 8.1 13/11/18

## 8.1.1 Hypothesis Testing

**Definition 8.1.** We define a **statistical hypothesis** as any sentence that specifies completely (or partially) the statistical model

$$\{f_X(x; \vartheta) \mid \vartheta \in \Theta\}.$$

If the specification is complete, then the hypothesis is called **simple**. If the specification is partial, then the statistical hypothesis is **composite**.

*Example 8.2.* For example,  $\vartheta = \vartheta^*$  is a simple hypothesis,  $\vartheta > \vartheta^*$  is a composite unidirectional hypothesis.

The definition of a hypothesis is rather general, but the important point is that a hypothesis makes a statement about the population. The goal of a hypothesis test is to decide, based on a sample from the population, which of the two complementary hypotheses is true.

**Definition 8.3.** The two complementary hypotheses in a hypothesis testing problem are called the **null hypothesis** and the **alternative hypothesis**, denoted by  $H_0$  and  $H_1$ , respectively.

We denote by  $\Theta_0 \subset \Theta$  the set of values of  $\Theta$  that are specified by the null hypothesis  $H_0$ . If  $\vartheta$  denotes a population parameter, the general format of the null and alternative hypotheses is

$$H_0 : \vartheta \in \Theta_0 \quad \text{and} \quad H_1 : \vartheta \notin \Theta_0.$$

The null statistical hypothesis is something that exists before the experiment. It is a hypothesis that exists until proven otherwise. The alternative of  $H_0$  (denoted by  $H_1$ , as we said) is the complement of  $H_0$ . In a hypothesis testing problem, after observing the sample, the experimenter must decide either to accept  $H_0$  as true or to reject  $H_0$  as false and decide  $H_1$  is true.

The decision rule about rejecting (or accepting)  $H_0$  has different interpretations. If we decide (based on the data) to reject  $H_0$ , this has the unique consequence that  $H_0$  is false. Otherwise, if we decide (based on the data) to accept  $H_0$ , this does not imply the support of  $H_0$ .

The hypothesis test is a decision rule based on the sample space, and this decision rule is characterised by a function defined on the set of all the possible values of  $\mathbf{X} = (X_1, \dots, X_n)$ . There will be some points  $\mathbf{x} \in \mathcal{R}_0 \subset \mathbb{R}^n$  (the *sample space*) such that the decision rule leads to rejecting  $H_0$ , and some other points  $\mathbf{x} \notin \mathcal{R}_0$  such that the decision rule leads to *not* rejecting  $H_0$ .

**Definition 8.4.** A *hypothesis testing procedure of hypothesis test* is a rule that specifies:

1. for which sample values the decision is made to accept  $H_0$  as true;
2. for which sample values  $H_0$  is rejected and  $H_1$  is accepted as true.

The subset of the sample space for which  $H_0$  will be rejected is called the **rejection region** (or **critical region**), defined as  $\mathcal{R}_0$ . The complement of the rejection region is called the **acceptance region**.

*Remark 8.5.* On a “philosophical” level, some people worry about the distinction between *rejection*  $H_0$  and *accepting*  $H_1$ . In the first case, there is nothing implied about what the experimenter *is* accepting, only that the state defined by  $H_0$  is being rejected. Similarly, a distinction can be made between *accepting*  $H_0$  and *not rejecting*  $H_0$ . The first phrase implies that the experimenter is willing to assert the state of nature specified by  $H_0$ , while the second phrase implies that the experimenter really does not believe  $H_0$ , but does not have evidence to reject it. For the most part, we will not be concerned with these issues. We view a hypothesis testing problem as a problem in which one of two actions is going to be taken (the actions being the assertion of  $H_0$  and  $H_1$ ).

At this stage, the decision rule provides a partition of the sample space in two subsets. However, we can work with such a partition only for values like  $n = 1, 2, 3$  for example, which of course is not realistic since  $n$  is the sample’s dimension. We must therefore reduce the dimensionality of the decision rule (or the dimensionality of the partition of the sample space). We can do that by using a *sufficient statistic*

$$T(X_1, \dots, X_n) = T(\mathbf{X}) = T_n,$$

a function  $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$  defined on the set of all possible values of  $\mathbf{X}$ . There will be some points  $\mathbf{X} \in \mathcal{R}_0 \subseteq \mathbb{R}^n$  such that the decision rule leads to reject  $H_0$ , and some other points  $\mathbf{X} \notin \mathcal{R}_0 \subseteq \mathbb{R}^n$  such that the decision rule leads to not reject  $H_0$ . The decision rule, therefore, is determined by the following criterion:

- if  $\mathbf{X} \in \mathcal{R}_0$ , then we reject  $H_0$ ;
- if  $\mathbf{X} \notin \mathcal{R}_0$ , then we do not reject  $H_0$ .

This allows us to “translate” the original decision rule into the *new* decision rule, which is

- if  $T_n \in C_0$ , then we reject  $H_0$ ;
- if  $T_n \notin C_0$ , then we do not reject  $H_0$ ,

where  $C_0 \subset T(\mathbb{R}_n)$  is the *rejection region* obtained by transforming  $\mathcal{R}_0$  through the statistic  $T_n$ .

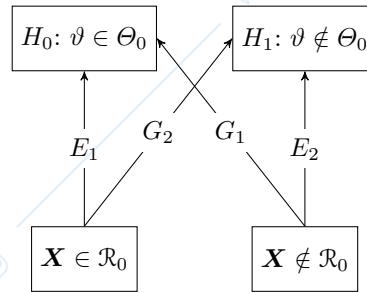
Therefore, briefly, what we obtained is the following: the **statistical hypothesis** defines a partition of the parametric space  $\Theta$ , meaning

$$\begin{cases} H_0 : \vartheta \in \Theta_0 \\ H_1 : \vartheta \notin \Theta_0 \end{cases} \quad (8.1)$$

The decision rule, on the other hand, defines a partition of the sample space, or a lower dimensional counterpart, meaning

$$\begin{cases} \text{if } T_n \in C_0, \text{ then we reject } H_0 \\ \text{if } T_n \notin C_0, \text{ then we do not reject } H_0 \end{cases} \quad (8.2)$$

A test is a link between this two partitions. The possible links are:



Now, consider the statistical hypothesis (8.1) and the decision rule (8.2).

- $G_1$  is the decision of accepting  $H_0$  (on the basis of the sample), where  $H_0$  is true: *correct decision*.
- $G_2$  is the decision of rejecting  $H_0$  (on the basis of the sample), where  $H_0$  is false: *correct decision*.
- $E_1$  is the decision of rejecting  $H_0$  (on the basis of the sample), despite  $H_0$  is true: *wrong decision*.
- $E_2$  is the decision of accepting  $H_0$  (on the basis of the sample), where  $H_0$  is false: *wrong decision*.

If we reject  $H_0$ , then we can take a correct or a wrong decision (the same is valid if we do not reject  $H_0$ ).

**Definition 8.6.** Define

- $\alpha := \mathbb{P}(E_1) = \mathbb{P}(\mathbf{X} \in \mathcal{R}_0, \vartheta \in \Theta_0)$ ;
- $\beta := \mathbb{P}(E_2) = \mathbb{P}(\mathbf{X} \notin \mathcal{R}_0, \vartheta \notin \Theta_0)$ ;
- $1 - \alpha := \mathbb{P}(G_1) = \mathbb{P}(\mathbf{X} \notin \mathcal{R}_0, \vartheta \in \Theta_0)$ ;
- $1 - \beta := \mathbb{P}(G_2) = \mathbb{P}(\mathbf{X} \in \mathcal{R}_0, \vartheta \notin \Theta_0)$ ,

where  $\alpha$  is the probability of the **first type error** and  $\beta$  is the probability of the **second type error**. Also,  $\gamma := 1 - \beta$  is known as the **power of the test**. Also, the **power function** will be defined as

$$\gamma(\vartheta^*) := \mathbb{P}(\mathbf{X} \in \mathcal{R}_0, \vartheta = \vartheta^*).$$

90 8 Week VIII

Our goal is to define a rejection region  $\mathcal{R}_0$  such that  $\alpha$  and  $\beta$  are as small as they can get. Fix a rejection region  $C_0$  and define

$$\alpha = \alpha(C_0) \quad \beta = \beta(C_0).$$

Suppose that  $C'_0$  is a different rejection region such that

$$\alpha(C'_0) < \alpha(C_0) \implies \beta(C'_0) > \beta(C_0). \quad (8.3)$$

Indeed, (8.3) is possible if we consider  $C'_0 \subset C_0$  such that

$$\alpha(C'_0) = \mathbb{P}(\mathbf{X} \in C'_0 | H_0) < \mathbb{P}(\mathbf{X} \in C_0 | H_0) = \alpha(C_0),$$

therefore  $(C'_0)^c \subset (C_0)^c$  implies

$$\begin{aligned} \beta(C'_0) &= \mathbb{P}(\mathbf{X} \notin C'_0 | H_1) \\ &= \mathbb{P}(\mathbf{X} \in (C'_0)^c | H_1) \\ &> \mathbb{P}(\mathbf{X} \in (C_0)^c | H_1) \\ &= \mathbb{P}(\mathbf{X} \notin C_0 | H_1) \\ &= \beta(C_0), \end{aligned}$$

meaning that  $\alpha$  and  $\beta$  are negatively correlated. Our way to proceed is the following: fix  $\alpha$  and look for the rejection region with the minimum  $\beta$ .

## 8.2 14/11/18

*Example 8.7.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_{30})$  from  $X \sim \text{Ber}(\vartheta)$ . We know that

$$\hat{\vartheta} = \frac{1}{30} \sum_{i=1}^{30} X_i.$$

Consider the following test

$$\begin{cases} H_0 : & \vartheta = \frac{1}{2} \\ H_1 : & \vartheta = \frac{2}{3} \end{cases} . \quad (8.4)$$

Of course,

$$\sum_{i=1}^{30} X_i \sim \mathcal{B}(30, \vartheta).$$

Now, take

1.  $C'_0$ : we reject  $H_0$  if the number of heads is greater than 25, meaning

$$\sum_{i=1}^{30} X_i > 25;$$

2.  $C''_0$ : we reject  $H_0$  if the number of heads is greater than 17

$$\sum_{i=1}^{30} X_i > 17.$$

(1) and (2) are two possible rejection regions; if we use  $C'_0$  as a decision rule, we would expect not to reject much, therefore it is likely to have a higher value of  $\beta$ . Otherwise, if we use  $C''_0$ , the value of  $\alpha$  would be greater than the one we would obtain adopting  $C'_0$ . Now,

$$\alpha(C'_0) = \mathbb{P} \left( \sum_{i=1}^{30} X_i > 25 \mid \vartheta = \frac{1}{2} \right) = \sum_{r=26}^{30} \binom{30}{r} \left( \frac{1}{2} \right)^r \left( \frac{1}{2} \right)^{30-r} = 0,000029$$

$$\beta(C'_0) = \mathbb{P} \left( \sum_{i=1}^{30} X_i \leq 25 \mid \vartheta = \frac{2}{3} \right) = \sum_{r=0}^{25} \binom{30}{r} \left( \frac{2}{3} \right)^r \left( \frac{1}{3} \right)^{30-r} = 0,9877$$

$$\alpha(C''_0) = \mathbb{P} \left( \sum_{i=1}^{30} X_i > 17 \mid \vartheta = \frac{1}{2} \right) = \sum_{r=18}^{30} \binom{30}{r} \left( \frac{1}{2} \right)^r \left( \frac{1}{2} \right)^{30-r} = 0,18$$

$$\beta(C''_0) = \mathbb{P} \left( \sum_{i=1}^{30} X_i \leq 17 \mid \vartheta = \frac{2}{3} \right) = \sum_{r=0}^{17} \binom{30}{r} \left( \frac{2}{3} \right)^r \left( \frac{1}{3} \right)^{30-r} = 0,166.$$

### 8.2.1 Neyman-Pearson Lemma

For a fixed sample size, it is usually impossible to make both types of error probabilities arbitrarily small.

**Definition 8.8.** For  $0 \leq \alpha \leq 1$ , a test with power function  $\gamma(\vartheta)$  is a **size  $\alpha$  test** if

$$\sup_{\vartheta \in \Theta_0} \gamma(\vartheta) = \alpha.$$

**Definition 8.9.** For  $0 \leq \alpha \leq 1$ , a test with power function  $\gamma(\vartheta)$  is a **level  $\alpha$  test** if

$$\sup_{\vartheta \in \Theta_0} \gamma(\vartheta) \leq \alpha.$$

According to our definitions, the set of level  $\alpha$  tests contains the set of size  $\alpha$  tests. Moreover, the distinction becomes important in complicated models and complicated testing situations, where it is often computationally impossible to construct a size  $\alpha$  test. In such situations, an experimenter must be satisfied with a level  $\alpha$  test, realising that some compromises may be made.

**Definition 8.10.** Let  $\mathcal{C}$  be a class of tests for testing

$$H_0 : \vartheta \in \Theta_0$$

versus

$$H_1 : \vartheta \in \Theta_0^c.$$

A test in the class  $\mathcal{C}$ , with power function  $\gamma(\vartheta)$ , is a **uniformly most powerful (UMP) class  $\mathcal{C}$  test** if

$$\gamma(\vartheta) \geq \gamma'(\vartheta), \quad \forall \vartheta \in \Theta_0^c$$

and every  $\gamma'(\cdot)$  that is a power function of a test in class  $\mathcal{C}$ .

In this section, the class  $\mathcal{C}$  will be the class of all level  $\alpha$  tests. Then, the test in Definition 8.10 is called a UMP level  $\alpha$  test. However, the requirements in Definition 8.10 are so strong that UMP tests do not exist in many realistic problems. But in problems that have UMP tests, a UMP test might well be considered the best test in the class. Thus, we would like to be able to identify UMP tests if they exist.

The following famous Theorem clearly describes which tests are UMP level  $\alpha$  tests in the the situation where the null and alternative hypotheses both consists of only one probability distribution for the sample (that is, when both  $H_0$  and  $H_1$  are simple hypotheses).

**Theorem 8.11 (Neyman-Pearson Lemma).** Let  $f_{\mathbf{X}}(\mathbf{x}; \vartheta_i)$  be the probability density function (or probability mass function) of  $\mathbf{X} = (X_1, \dots, X_n)$  corresponding to  $\vartheta_i$ . Consider testing

$$H_0 : \vartheta = \vartheta_0$$

versus

$$H_1 : \vartheta = \vartheta_1$$

using a test with rejection region  $\mathcal{R}_0$  that satisfies

$$\begin{cases} \mathbf{x} \in \mathcal{R}_0 & \text{if } f_{\mathbf{X}}(\mathbf{x}; \vartheta_1) > k f_{\mathbf{X}}(\mathbf{x}; \vartheta_0) \\ \mathbf{x} \in \mathcal{R}_0^c & \text{if } f_{\mathbf{X}}(\mathbf{x}; \vartheta_1) < k f_{\mathbf{X}}(\mathbf{x}; \vartheta_0), \end{cases} \quad (8.5)$$

the *shape condition of the rejection region*, for some  $k \geq 0$  and

$$\alpha = \mathbb{P}(\mathbf{X} \in \mathcal{R}_0; \vartheta = \vartheta_0), \quad (8.6)$$

the *size condition of the rejection region*. Then,

1. **sufficiency**: any test that satisfies (8.5) and (8.6) is a UMP level  $\alpha$  test;
2. **necessity**: if there exist a test satisfying (8.5) and (8.6) with  $k > 0$ , then every UMP level  $\alpha$  test is a size  $\alpha$  test - satisfies (8.6) - and every UMP level  $\alpha$  test satisfies (8.5) except perhaps on a set  $\mathcal{A}$  satisfying  $\mathbb{P}(\mathbf{X} \in \mathcal{A}; \vartheta_0) = \mathbb{P}(\mathbf{X} \in \mathcal{A}; \vartheta_1) = 0$ .

*Proof.* We will prove the Theorem for the case that  $f_{\mathbf{X}}(\mathbf{x}; \vartheta_0)$  and  $f_{\mathbf{X}}(\mathbf{x}; \vartheta_1)$  are probability density functions of continuous random variables. The proof for discrete random variables can be accomplished by replacing integrals with sums.

Note first that any test that satisfies (8.6) is a size  $\alpha$  test and, hence, a level  $\alpha$  test: indeed,

$$\sup_{\vartheta \in \Theta_0} \mathbb{P}(\mathbf{X} \in \mathcal{R}_0; \vartheta) = \mathbb{P}(\mathbf{X} \in \mathcal{R}_0; \vartheta_0) = \alpha,$$

since  $\Theta_0$  has only one point. Define

$$\phi(\mathbf{x}) := \begin{cases} 1 & \text{if } \mathbf{x} \in \mathcal{R}_0 \\ 0 & \text{if } \mathbf{x} \in \mathcal{R}_0^c \end{cases},$$

so that

- $\phi(\cdot)$  be a test function that satisfies (8.5) and (8.6);
- $\phi'(\cdot)$  be a test function of any other level  $\alpha$  test;
- $\gamma(\vartheta)$  and  $\gamma'(\vartheta)$  be the power functions corresponding to the test  $\phi(\cdot)$  and  $\phi'(\cdot)$ , respectively.

Now, (8.5) implies that

$$[\phi(\mathbf{x}) - \phi'(\mathbf{x})] (f_{\mathbf{X}}(\mathbf{x}; \vartheta_1) - k f_{\mathbf{X}}(\mathbf{x}; \vartheta_0)) \geq 0, \quad \forall \mathbf{x}$$

since  $\phi(\cdot) = 1$  if  $f_{\mathbf{X}}(\mathbf{x}; \vartheta_1) > k f_{\mathbf{X}}(\mathbf{x}; \vartheta_0)$  and  $\phi(\cdot) = 0$  if  $f_{\mathbf{X}}(\mathbf{x}; \vartheta_1) < k f_{\mathbf{X}}(\mathbf{x}; \vartheta_0)$ . Thus, for every  $\vartheta \in \Theta$ ,

$$\begin{aligned} 0 &\leq \int [\phi(\mathbf{x}) - \phi'(\mathbf{x})] (f_{\mathbf{X}}(\mathbf{x}; \vartheta_1) - k f_{\mathbf{X}}(\mathbf{x}; \vartheta_0)) dx \\ &= \gamma(\vartheta_1) - \gamma'(\vartheta_1) - k [\gamma(\vartheta_0) - \gamma'(\vartheta_0)]. \end{aligned} \quad (8.7)$$

Statement (1) is proved by noting that, since  $\phi'$  is a level  $\alpha$  test and  $\phi$  is a size  $\alpha$  test,

$$\gamma(\vartheta_0) - \gamma'(\vartheta_0) = \alpha - \gamma'(\vartheta_0) \geq 0.$$

Thus (8.7) and  $k \geq 0$  imply that

$$0 \leq \gamma(\vartheta_1) - \gamma'(\vartheta_1) - k [\gamma(\vartheta_0) - \gamma'(\vartheta_0)] \leq \gamma(\vartheta_1) - \gamma'(\vartheta_1),$$

showing that  $\gamma(\vartheta_1) \geq \gamma'(\vartheta_1)$ , and hence  $\phi(\cdot)$  has greater power than  $\phi'(\cdot)$ . Since  $\phi'(\cdot)$  was an arbitrary level  $\alpha$  test and  $\vartheta_1$  is the only point in  $\Theta_0^c$ ,  $\phi(\cdot)$  is a UMP level  $\alpha$  test.

To prove statement (2), let  $\phi'(\cdot)$  now be a test function for any UMP level  $\alpha$  test. By the previous part, the test  $\phi(\cdot)$  satisfying (8.5) and (8.6) is also a UMP level  $\alpha$  test, which implies that

$$\gamma(\vartheta_1) = \gamma'(\vartheta_1).$$

This fact, (8.7) and  $k > 0$  imply

$$\alpha - \gamma'(\vartheta_0) = \gamma(\vartheta_0) - \gamma'(\vartheta_0) \leq 0.$$

Now, since  $\phi'$  is a level  $\alpha$  test, we have that  $\gamma'(\vartheta_0) \leq \alpha$ , thus  $\gamma'(\vartheta_0) = \alpha$ . This implies that  $\phi'(\cdot)$  is a size  $\alpha$  test; this also implies that (8.7) is an equality in this case. But the non-negative integrand

$$[\phi(\mathbf{x}) - \phi'(\mathbf{x})][f_{\mathbf{X}}(\mathbf{x}; \vartheta_1) - kf_{\mathbf{X}}(\mathbf{x}; \vartheta_0)]$$

will have a zero integral only if  $\phi'(\cdot)$  satisfies (8.5) except perhaps on a set  $\mathcal{A}$  with

$$\int_{\mathcal{A}} f_{\mathbf{X}}(\mathbf{x}; \vartheta_i) d\mathbf{x} = 0, \quad i = 0, 1.$$

This implies that the last assertion in statement (2) is true.  $\square$

*Example 8.12.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \mathcal{N}(\vartheta, \sigma^2)$ , where  $\sigma^2$  is supposed to be known. We would like to test

$$H_0 : \vartheta = \vartheta_0$$

versus

$$H_1 : \vartheta = \vartheta_1,$$

where  $\vartheta_1 > \vartheta_0$ . For a specified value of  $\alpha$ , a level  $\alpha$  test in this problem is any test that satisfies

$$\mathbb{P}(\mathbf{X} \in \mathcal{R}_0; \vartheta_0).$$

We would like to compute

$$\frac{\mathcal{L}(\vartheta_1; \mathbf{x})}{\mathcal{L}(\vartheta_0; \mathbf{x})} \geq c,$$

so

$$\begin{aligned} & \frac{\exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \vartheta_1)^2\right\}}{\exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \vartheta_0)^2\right\}} \geq c \\ & \exp\left\{-\frac{1}{2\sigma^2} \left[\sum_{i=1}^n (x_i - \vartheta_1)^2 - \sum_{i=1}^n (x_i - \vartheta_0)^2\right]\right\} \geq c \\ & \exp\left\{-\frac{1}{2\sigma^2} \left[\sum_{i=1}^n x_i^2 + n\vartheta_1^2 - 2\vartheta_1 \sum_{i=1}^n x_i - \sum_{i=1}^n x_i^2 - n\vartheta_0^2 + 2\vartheta_0 \sum_{i=1}^n x_i\right]\right\} \geq c \end{aligned}$$

$$\begin{aligned}
& \exp \left\{ -\frac{1}{2\sigma^2} \left[ n(\vartheta_1^2 - \vartheta_0^2) - 2(\vartheta_1 - \vartheta_0) \sum_{i=1}^n x_i \right] \right\} \geq c \\
& \exp \left\{ -\frac{1}{2\sigma^2} (\vartheta_1 - \vartheta_0) \left[ n(\vartheta_1 + \vartheta_0) - 2 \sum_{i=1}^n x_i \right] \right\} \geq c \\
& -(\vartheta_1 - \vartheta_0) \left[ n(\vartheta_1 - \vartheta_0) - 2 \sum_{i=1}^n x_i \right] \geq 2\sigma^2 \log c \\
& n(\vartheta_1 - \vartheta_0) - 2 \sum_{i=1}^n x_i \leq -\frac{2\sigma^2 \log c}{\vartheta_1 - \vartheta_0} \\
& -2 \sum_{i=1}^n x_i \leq -\frac{2\sigma^2 \log c}{\vartheta_1 - \vartheta_0} - n(\vartheta_1 - \vartheta_0) \\
& \sum_{i=1}^n x_i \geq \frac{\sigma^2 \log c}{\vartheta_1 - \vartheta_0} + \frac{n}{2}(\vartheta_1 - \vartheta_0)
\end{aligned}$$

is the shape of the rejection region.

*Example 8.13.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim \text{Ber}(\vartheta)$ . We would like to test

$$H_0 : \vartheta = \vartheta_0$$

versus

$$H_1 : \vartheta = \vartheta_1,$$

where  $\vartheta_1 > \vartheta_0$ . Recall the alternative statement (??): we would like to compute

$$\frac{\mathcal{L}(\vartheta_0; \mathbf{x})}{\mathcal{L}(\vartheta_1; \mathbf{x})} \leq k,$$

so

$$\begin{aligned}
& \frac{\vartheta_0^{\sum_{i=1}^n x_i} (1 - \vartheta_0)^{n - \sum_{i=1}^n x_i}}{\vartheta_1^{\sum_{i=1}^n x_i} (1 - \vartheta_1)^{n - \sum_{i=1}^n x_i}} \leq k \\
& \left( \frac{\vartheta_0}{\vartheta_1} \right)^{\sum_{i=1}^n x_i} \left( \frac{1 - \vartheta_0}{1 - \vartheta_1} \right)^{n - \sum_{i=1}^n x_i} \leq k \\
& \sum_{i=1}^n x_i \log \left( \frac{\vartheta_0}{\vartheta_1} \right) + \left( n - \sum_{i=1}^n x_i \right) \log \left( \frac{1 - \vartheta_0}{1 - \vartheta_1} \right) \leq \log k \\
& \sum_{i=1}^n x_i \left[ \log \left( \frac{\vartheta_0}{\vartheta_1} \right) - \log \left( \frac{1 - \vartheta_0}{1 - \vartheta_1} \right) \right] \leq \log k - n \log \left( \frac{1 - \vartheta_0}{1 - \vartheta_1} \right) \\
& \sum_{i=1}^n x_i \geq k'.
\end{aligned}$$

Here,

$$\alpha = \mathbb{P} \left( \sum_{i=1}^n x_i \geq k'; \vartheta = \vartheta_0 \right) = \sum_{r=k'}^n \binom{n}{r} \vartheta_0^r (1 - \vartheta_0)^{n-r}.$$

96 8 Week VIII

We are looking for a value of  $k'$  such that a very close  $\alpha$  to the  $\alpha$  chosen is provided, that is because  $k \in \{0, \dots, n\}$ , and in general it does not exist  $k$  that determines the exact value of  $\alpha$ .

### 8.3 15/11/18

*Example 8.14.* Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $X \sim EF(\cdot)$ , the exponential family, where

$$f_{\mathbf{X}}(\mathbf{x}; \vartheta) = \exp\{A(\mathbf{x})Q(\vartheta) + C(\mathbf{x}) - K(\vartheta)\},$$

and  $Q$  is monotone increasing. We want to test

$$H_0 : \vartheta = \vartheta_0, \quad H_1 : \vartheta = \vartheta_1.$$

To do so, we can apply the Neyman-Pearson Lemma:

$$\begin{aligned} \frac{\mathcal{L}(\vartheta_0; \mathbf{x})}{\mathcal{L}(\vartheta_1; \mathbf{x})} &\leq k \\ \frac{\prod_{i=1}^n \exp\{A(x_i)Q(\vartheta_0) + C(x_i) - K(\vartheta_0)\}}{\prod_{i=1}^n \exp\{A(x_i)Q(\vartheta_1) + C(x_i) - K(\vartheta_1)\}} &\leq k \\ \exp\left\{Q(\vartheta_0) \sum_{i=1}^n A(x_i) + \sum_{i=1}^n C(x_i) - nK(\vartheta_0) - Q(\vartheta_1) \sum_{i=1}^n A(x_i) - \sum_{i=1}^n C(x_i) + nK(\vartheta_1)\right\} &\leq k \\ \exp\left\{[Q(\vartheta_0) - Q(\vartheta_1)] \sum_{i=1}^n A(x_i) - n[K(\vartheta_0) - K(\vartheta_1)]\right\} &\leq k \\ [Q(\vartheta_0) - Q(\vartheta_1)] \sum_{i=1}^n A(x_i) - n[K(\vartheta_0) - K(\vartheta_1)] &\leq \log k \\ \sum_{i=1}^n A(x_i) &\geq \frac{\log k + n[K(\vartheta_0) - K(\vartheta_1)]}{Q(\vartheta_0) - Q(\vartheta_1)}. \end{aligned}$$

We conclude that

$$\sum_{i=1}^n A(x_i) \geq k'$$

is the shape of the rejection region for an arbitrary exponential family with  $Q(\cdot)$  monotone increasing, meaning that  $Q(\vartheta_0) < Q(\vartheta_1)$ .

We will see an extension of the Neyman-Pearson Lemma for

1. *unidirectional hypothesis (Rubin Test)*;
2. *arbitrary hypothesis (Likelihood Ratio Test)*.

**Corollary 8.15.** Consider the hypotheses of the Neyman-Pearson Lemma, and suppose that  $T(\mathbf{X})$  is a sufficient statistic for the parameter  $\vartheta$ . Denote by  $f_{\mathbf{X}}(\mathbf{x}; \vartheta_i)$ , for  $i = 0, 1$ , the joint density function of the random sample, and by  $g_T(t; \vartheta_i)$ , for  $i = 0, 1$ , the density function of the sufficient statistic. Then, any test based on  $T$  with rejection region  $S$  (a subset of the sample space defined by  $T$ ) is a UMP level  $\alpha$  test if it satisfies

$$\begin{aligned} g_T(t; \vartheta_1) &> k g_T(t; \vartheta_0), \quad t \in S \\ g_T(t; \vartheta_1) &< k g_T(t; \vartheta_0), \quad t \in S^c \end{aligned} \tag{8.8}$$

for some  $k \geq 0$ , where

$$\alpha = \mathbb{P}(T(\mathbf{X}) \in S; \vartheta_0). \tag{8.9}$$

*Proof.* In terms of the random sample  $\mathbf{X} = (X_1, \dots, X_n)$ , the test based on the sufficient statistic  $T$  has the following rejection region:

$$\mathcal{R}_0 = \{\mathbf{x} \mid T(\mathbf{x}) \in S\}.$$

By the Factorisation Theorem, the probability density function (or probability mass function) of the random sample  $\mathbf{X} = (X_1, \dots, X_n)$  can be written as

$$f_{\mathbf{X}}(\mathbf{x}; \vartheta_i) = g_T(T(\mathbf{x}); \vartheta_i)h(\mathbf{x}), \quad i = 0, 1$$

for some non-negative function  $h(\cdot)$ . Multiplying the inequalities in (8.8) by this non-negative function, we see that  $\mathcal{R}_0$  satisfies

$$\mathbf{x} \in \mathcal{R}_0 \text{ if } f_{\mathbf{X}}(\mathbf{x}; \vartheta_1) = g_T(T(\mathbf{x}); \vartheta_1)h(\mathbf{x}) > k g_T(T(\mathbf{x}); \vartheta_0)h(\mathbf{x}) = k f_{\mathbf{X}}(\mathbf{x}; \vartheta_0)$$

and

$$\mathbf{x} \in \mathcal{R}_0^c \text{ if } f_{\mathbf{X}}(\mathbf{x}; \vartheta_1) = g_T(T(\mathbf{x}); \vartheta_1)h(\mathbf{x}) < k g_T(T(\mathbf{x}); \vartheta_0)h(\mathbf{x}) = k f_{\mathbf{X}}(\mathbf{x}; \vartheta_0).$$

Also, by condition (8.9), we have

$$\mathbb{P}(\mathbf{X} \in \mathcal{R}_0; \vartheta_0) = \mathbb{P}(T(\mathbf{X}) \in S; \vartheta_0) = \alpha.$$

So, by the sufficiency part of the Neyman-Pearson Lemma, the test based on  $T$  is a UMP level  $\alpha$  test.  $\square$

*Example 8.16.* Consider

$$f_{\mathbf{X}}(\mathbf{x}; \vartheta) = \exp \{A(\mathbf{x})Q(\vartheta) + C(\mathbf{x}) - K(\vartheta)\} :$$

a sufficient statistic could be  $\sum_{i=1}^n A(x_i)$ . If we wanted to test

$$H_0 : \vartheta = \vartheta_0, \quad H_1 : \vartheta = \vartheta_1 > \vartheta_0,$$

we could use the fact that the shape of the rejection region is

$$\sum_{i=1}^n A(x_i) \geq c.$$

In general, by the previous corollary, the shape of the rejection region defined by the Neyman-Pearson Lemma is a function of the sufficient statistic for the parameter of interest.

## 9

## Week IX

## 9.1 20/11/18

## 9.1.1 Rubin Test

As we have seen, the Neyman-Pearson Lemma comes in handy when dealing with *simple* hypotheses. However, in more realistic problems, the hypotheses of interest specify more than one possible distribution for the sample (*composite hypotheses*). A large class of problems that admit a UMP level  $\alpha$  test involve one-side hypotheses and probability density functions, or probability mass functions, with the Monotone Likelihood Ratio property.

**Definition 9.1.** A family of probability density functions, or probability mass functions, namely

$$\{g_T(t; \vartheta) \mid \vartheta \in \Theta\}$$

for a univariate random variable  $T$  with real-valued parameter  $\vartheta$  has a **Monotone Likelihood Ratio (MLR)** if, for every  $\vartheta_2 > \vartheta_1$ ,

$$\frac{g_T(t; \vartheta_2)}{g_T(t; \vartheta_1)}$$

is a monotone (non-increasing or non-decreasing) function of  $t$  on

$$\{t \mid g_T(t; \vartheta_1) > 0 \vee g_T(t; \vartheta_2) > 0\}.$$

*Example 9.2.* Many common families of distributions have the MLR property.

1. *Gaussian*, with respect to the mean for a fixed variance;
2. *Binomial*;
3. *Poisson*.

In general, any regular exponential family, with density function or probability mass function

$$g(t; \vartheta) = h(t)c(\vartheta)e^{w(\vartheta)t}$$

has the MLR property if  $w(\vartheta)$  is a non-decreasing function.

There exist a generalisation of the Neyman-Pearson Lemma for statistical models with MLR property, which allows us to consider unidirectional (*composite*) hypotheses. This is called **Rubin test**.

**Theorem 9.3 (Rubin Test).** Consider the problem of testing

$$H_0 : \vartheta \leq \vartheta_0, \quad H_1 : \vartheta > \vartheta_0.$$

Suppose that  $T$  is a sufficient statistic for  $\vartheta$ , and that the family of probability density functions, or probability mass functions,

$$\{g(t; \vartheta) \mid \vartheta \in \Theta\}$$

of  $T$  has the MLR property. Then, for any  $t_0$ , the statistical test that rejects  $H_0$  if and only if  $T > t_0$  is a UMP level  $\alpha$  test, where  $\alpha = \mathbb{P}(T > t_0; \vartheta_0)$ .

### 9.1.2 Likelihood Ratio Test

**Definition 9.4.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from a parametric model  $X \sim f_{\mathbf{X}}(\mathbf{x}; \vartheta)$ . Assume we want to test something very general, like

$$H_0 : \vartheta \in \Theta_0, \quad H_1 : \vartheta \notin \Theta_0.$$

We define the **Likelihood Ratio**  $\lambda(\mathbf{x})$  as

$$\lambda(\mathbf{x}) := \frac{\sup_{\vartheta \in \Theta_0} \mathcal{L}(\vartheta; \mathbf{x})}{\sup_{\vartheta \in \Theta} \mathcal{L}(\vartheta; \mathbf{x})}.$$

We consider a rejection region of the form

$$\lambda(\mathbf{x}) \leq c,$$

such that  $\mathbb{P}(\lambda(\mathbf{x}) \leq c \mid H_0) = \alpha$ . Of course,  $H_0$  is not specified.

*Remark 9.5.* The numerator of  $\lambda(\mathbf{x})$  is the maximum probability of the observed sample. Supposing the parameter lives in the null hypothesis, the denominator is the maximum probability of the observed sample over all possible parameters. The ratio is small if there are parameter points in the alternative, for which the observed sample is much more likely than any other parameter in the null.

*Example 9.6.* Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $X \sim \mathcal{N}(\vartheta, \sigma^2)$ , where  $\sigma^2$  is unknown. Consider the following test:

$$H_0 : \vartheta = \vartheta_0, \quad H_1 : \vartheta \neq \vartheta_0.$$

In this case,

$$\Theta = \mathbb{R} \times \mathbb{R}_+ = \{(\vartheta, \sigma^2) \in (\mathbb{R} \times \mathbb{R}_+)\}, \quad \Theta_0 = \{(\vartheta_0, \sigma^2) \in (\vartheta_0, \mathbb{R}_+)\}.$$

Without any restriction on the parametric space  $\Theta$ , the Maximum Likelihood estimators of  $(\vartheta, \sigma^2)$  are

$$\hat{\vartheta} = \bar{X}_n, \quad \hat{\sigma}^2 = S_n^2.$$

Therefore,

$$\begin{aligned}\mathcal{L}(\bar{X}_n, S_n^2; \mathbf{x}) &= \left( \frac{1}{\sqrt{2\pi S_n^2}} \right)^n \exp \left\{ -\frac{1}{2S_n^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \right\} \\ &\propto (S_n^2)^{-\frac{n}{2}} e^{-\frac{n}{2}} \\ &\propto (S_n^2)^{-\frac{n}{2}}\end{aligned}$$

Now, we maximise the Likelihood function with respect to  $\Theta_0$ : if  $\hat{\vartheta} = \vartheta_0$ ,

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \vartheta_0)^2 = S_0^2,$$

therefore

$$\mathcal{L}(\vartheta_0, S_0^2) = \left( \frac{1}{\sqrt{2\pi S_0^2}} \right)^n \exp \left\{ -\frac{1}{2S_0^2} \sum_{i=1}^n (X_i - \vartheta_0)^2 \right\} \propto (S_0^2)^{-\frac{n}{2}},$$

meaning that the Likelihood Ratio test is given by

$$\lambda(\mathbf{x}) = \left( \frac{S_n^2}{S_0^2} \right)^{n/2} \leq c.$$

Since

$$\begin{aligned}S_0^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \vartheta_0)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n + \bar{X}_n - \vartheta_0)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 + (\bar{X}_n - \vartheta_0)^2 + \underbrace{\frac{2}{n} \sum_{i=1}^n (X_i \bar{X}_n - X_i \vartheta_0 - \bar{X}_n^2 + \bar{X}_n \vartheta_0)}_{(*)} \\ &= S_n^2 + (\bar{X}_n - \vartheta_0)^2,\end{aligned}$$

where  $(*) = 0$  because

$$\begin{aligned}&\frac{2}{n} \vartheta_0 \sum_{i=1}^n (\bar{X}_n - X_i) + \frac{2}{n} \bar{X}_n \sum_{i=1}^n (X_i - \bar{X}_n) \\ &= 2\vartheta_0 \left( \bar{X}_n - \frac{1}{n} \sum_{i=1}^n X_i \right) + 2\bar{X}_n \left( \frac{1}{n} \sum_{i=1}^n X_i - \bar{X}_n \right) \\ &= 0.\end{aligned}$$

Hence

$$\lambda(\mathbf{x}) = \left( \frac{S_n^2}{S_n^2 + (\bar{X}_n - \vartheta_0)^2} \right)^{\frac{n}{2}} \leq c$$

$$\begin{aligned} \left(1 + \frac{(\bar{X}_n - \vartheta_0)^2}{S_n^2}\right)^{-\frac{n}{2}} &\leq c \\ \frac{(\bar{X}_n - \vartheta_0)^2}{S_n^2} &\geq c' \\ \frac{(\bar{X}_n - \vartheta_0)^2}{\frac{n-1}{n}S_n^2} &\geq c' \\ \frac{\sqrt{n}|\bar{X}_n - \vartheta_0|}{S_n} &\geq c' \\ |T_n| &\geq c', \end{aligned}$$

where

$$T_n = \frac{\sqrt{n}(\bar{X}_n - \vartheta_0)}{s_n} \sim t_{n-1}$$

under  $H_0$ . Therefore,

$$\mathbb{P}(-t_{\alpha/2} \leq T_n \leq t_{\alpha/2}) = 1 - \alpha$$

and the *rejection region of level  $\alpha$*  is

$$|T_n| > t_{\alpha/2}.$$

*Remark 9.7.* 1. The rejection region defined by the Likelihood Ratio Test is obtained by taking the ratio between the Likelihood function evaluated in the point of maximum, with respect to a parametric space **restricted** by a null hypothesis, and the Likelihood function evaluated in the point of maximum with respect to the **unrestricted** parametric space.

2. The rejection region defined by the Likelihood Ratio test is always a function of the sufficient statistic, for the parameter of interest.
3. If the hypotheses are “simple enough”, then the Likelihood Ratio test becomes the Neyman-Pearson Lemma.

**Exercise 9.8.** The number of breaks of an object every week is a Poisson random variable with parameter  $\lambda$ . Consider the number of breaks in subsequent weeks to be independent. Let  $X_1, \dots, X_{10}$  be the number of breaks in ten weeks, and let

$$\bar{X}_{10} = \frac{1}{10} \sum_{i=1}^{10} X_i.$$

1. Find  $t$  such that  $\{\bar{X}_{10} > t\}$  is a rejection region of level  $\alpha = 0,05$  for testing

$$H_0 : \lambda = 2, \quad H_1 : \lambda = 3.$$

Compute the power of the test.

2. Is the test above a UMP level  $\alpha$  test?

*Proof.* (1) Under  $H_0$ , the random variables  $X_1, \dots, X_{10}$  are Poisson random variables with parameter  $\lambda = 2$ . If we set  $S := \sum_{i=1}^{10} X_i$ , then

$$\mathbb{P}(\bar{X}_{10} > t) = \mathbb{P}(S > 10t).$$

We need to find the point  $k$  such that  $\mathbb{P}(S \leq k) \geq 0,95$ ; then, we set  $t = k/10$  to solve

$$\mathbb{P}(S \leq k) \geq 0,95.$$

Recall that  $S$ , being the sum of Poisson random variables, is itself a Poisson random variable. Therefore,

$$\begin{aligned} \mathbb{P}\left(\sum_{i=1}^{10} X_i \leq 10t; \lambda = 2\right) &= \mathbb{P}(S \leq 10t; \lambda = 2) \\ &= \mathbb{P}\left(\frac{R - 10\lambda}{\sqrt{10\lambda}} \leq \frac{10t - 10\lambda}{\sqrt{10\lambda}}; \lambda = 2\right) \\ &\approx \Phi\left(\frac{10t - 10\lambda}{\sqrt{10\lambda}}; \lambda = 2\right) \\ &= \Phi\left(\frac{10t - 20}{\sqrt{20}}\right) \\ &\geq 0.95 \end{aligned}$$

if and only if

$$\sqrt{5}(t - 2) \geq Z_{0.05}.$$

Therefore,

$$t \geq \frac{1.64}{\sqrt{5}} + 2.$$

In order to find the power of the test with, we have to compute a similar probability as above:

$$\begin{aligned} \gamma &:= 1 - \beta \\ &= 1 - \mathbb{P}\left(\sum_{i=1}^{10} X_i \leq 10t; \lambda = 3\right) \\ &= \mathbb{P}(S \leq 10t; \lambda = 3) \\ &= 1 - \mathbb{P}\left(\frac{R - 10\lambda}{\sqrt{10\lambda}} \leq \frac{10t - 10\lambda}{\sqrt{10\lambda}}; \lambda = 3\right) \\ &\approx 1 - \Phi\left(\frac{10t - 10\lambda}{\sqrt{10\lambda}}; \lambda = 3\right) \\ &= 1 - \Phi\left(\frac{10t - 30}{\sqrt{30}}\right) \end{aligned}$$

(2) In order to find a test with the highest power, we can apply the Neyman-Pearson Lemma.

$$\frac{\mathcal{L}(\lambda_1; \mathbf{x})}{\mathcal{L}(\lambda_0; \mathbf{x})} = \frac{\prod_{i=1}^{10} e^{-\lambda_1} \frac{\lambda_1^{x_i}}{x_i!}}{\prod_{i=1}^{10} e^{-\lambda_0} \frac{\lambda_0^{x_i}}{x_i!}}.$$

Then, the rejection region of point (1) is a UMP level  $\alpha$  rejection region.

□

## 9.2 21/11/18

So far, we have found the following techniques for the problem

$$\alpha = \mathbb{P}(T \geq k; H_0) :$$

1. Neyman-Pearson Lemma: UMP level  $\alpha$  test for simple hypotheses and for an arbitrary statistical model;
2. Rubin Test: UMP level  $\alpha$  test for composite unidirectional hypotheses, for statistical models with Monotone Likelihood Ratio;
3. Likelihood Ratio Test: level  $\alpha$  test for any hypothesis and for arbitrary statistical models.

One of the most useful approach to obtain the shape of the rejection region is the *Likelihood Ratio Test*:

$$\lambda(\mathbf{x}) = \frac{\max_{\vartheta \in \Theta_0} \mathcal{L}(\vartheta; \mathbf{x})}{\max_{\vartheta \in \Theta} \mathcal{L}(\vartheta; \mathbf{x})},$$

with its form is given by

$$\{\mathbf{x} \mid \lambda(\mathbf{x}) \leq c\}.$$

After we gather the data  $\mathbf{X} = \mathbf{x}$ , the Likelihood function  $\mathcal{L}(\vartheta; \mathbf{x})$  is a completely defined function of the variable  $\vartheta$ . Even if we cannot compute analytically the points that maximise the Likelihood function over the sets  $\Theta_0$  and  $\Theta$ , these can be always evaluated *numerically*. Thus, the test statistic  $\lambda(\mathbf{x})$  can be obtained for the observed data point even if no convenient formula defining  $\lambda(\mathbf{x})$  is available.

To define a level  $\alpha$  test, the constant  $c$  must be chosen so that

$$\sup_{\vartheta \in \Theta_0} \mathbb{P}(\lambda(\mathbf{X}) \leq c; \vartheta) \leq \alpha. \quad (9.1)$$

If we cannot derive a simple formula for  $\lambda(\mathbf{x})$ , it might seem that it is hopeless to derive the sampling distribution of  $\lambda(\mathbf{X})$  and thus know how to pick  $c$  to ensure (9.1). This means that we must know the distribution of  $\lambda(\mathbf{x})$ . A solution to this problem consists in determining a general result for the asymptotic behaviour of  $\lambda(\mathbf{x})$ , as  $n \rightarrow \infty$ .

**Theorem 9.9 (Wilks).** Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from a probability density function (or probability mass function)  $X \sim f_X(\mathbf{x}, \vartheta)$ . Under some regularity conditions on the model  $f_X(x; \vartheta)$ , if  $\Theta_0 = \{\vartheta_0\}$ , then the distribution of the statistic

$$-2 \log \lambda(\mathbf{x}) \xrightarrow{n \rightarrow \infty} \chi_1^2,$$

where  $\chi_1^2$  is a Chi-squared random variable with one degree of freedom.

*Example 9.10.* Consider the problem of testing

$$H_0 : \lambda = \lambda_0, \quad H_1 : \lambda \neq \lambda_0.$$

Consider a random sample from  $X \sim \mathcal{P}(\lambda)$ . Therefore, applying Theorem 9.9,

$$-2 \log \lambda(\mathbf{x}) = -2 \log \left( \frac{e^{-n\lambda_0} \lambda_0^{\sum_{i=1}^n x_i}}{e^{-n\hat{\lambda}} \hat{\lambda}^{\sum_{i=1}^n x_i}} \right) = 2n \left[ (\lambda_0 - \hat{\lambda}) - \hat{\lambda} \log \left( \frac{\lambda_0}{\hat{\lambda}} \right) \right],$$

where

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

By a direct application of Wilks' Theorem, we reject  $H_0$  at level  $\alpha$  when

$$-2 \log \lambda(\mathbf{x}) > \chi_{1,\alpha}^2,$$

where the right hand side of the inequality indicates the quantile of order  $\alpha$  of a Chi-squared distribution.

*Remark 9.11.* As the Cramer Theorem, Wilks' Theorem is a "qualitative" result. That is, it does not provide a level of the sample size for which the theorem can be applied.

*Example 9.12.* Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $X \sim f_X(x; \vartheta)$ , where

$$f_X(x, \vartheta) = \exp\{-(x - \vartheta)\} \mathbb{1}_{(\vartheta, \infty)}(x).$$

Consider the problem of testing

$$H_0 : \vartheta = \vartheta_0, \quad H_1 : \vartheta > \vartheta_0.$$

If we apply the Likelihood Ratio Test,

$$\lambda(\mathbf{x}) = \frac{\prod_{i=1}^n \exp\{-(x_i - \vartheta_0)\}}{\prod_{i=1}^n \exp\{-(x_i - t_n)\}} = \exp\{-n(t_n - \vartheta_0)\},$$

where  $t_n = x_{(1)}$ . Therefore,

$$-2 \log \lambda(\mathbf{x}) = 2n(T_n - \vartheta_0),$$

where  $T_n = X_{(1)} = \min\{X_1, \dots, X_n\}$ . If we define  $H = 2n(T_n - \vartheta_0)$ , under  $H_0$ , the distribution of  $H$  is given by

$$\mathbb{P}(H \leq h) = \mathbb{P}\left(T_n - \vartheta_0 \leq \frac{h}{2n}\right),$$

where  $T_n$  is a negative exponential of parameter  $n$  and  $T_n - \vartheta_0$  is a negative exponential of parameter one.

Recall that, if  $X_1, \dots, X_n$  with  $X_i \sim \mathcal{E}_n(1)$ , then  $\min_i X_i \sim \mathcal{E}_n(n)$ . Therefore

$$\mathbb{P}(2n(T_n - \vartheta_0) \leq h) = 1 - \exp\left\{-\frac{h}{2}\right\},$$

meaning that  $2n(T_n - \vartheta_0) \sim \mathcal{E}_n(\frac{1}{2})$ : therefore, the Negative Exponential distribution with parameter  $1/2$  is a Chi-squared distribution with two degrees of freedom.

Wilks' Theorem can be stated also for statistical (regular) models, whose parametric space has dimension greater than one.

**Theorem 9.13 (Kendal-Stuart).** *Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from a probability density function (or probability mass function)  $f_X(x; \vartheta)$ . Under the usual conditions of regularity, the statistic  $-2 \log \lambda(\mathbf{x})$  converges weakly, as  $n \rightarrow \infty$ , to a Chi-squared distribution. The degrees of freedom of the limiting distribution is the difference between the number of free parameters specified by the null hypothesis and the number of free parameters specified of the model.*

*Example 9.14 (Multinomial model).* Let  $\vartheta = (p_1, \dots, p_5)$ , where the  $p_j$ 's are non-negative and sum up to 1. Suppose that we consider a random sample  $(X_1, \dots, X_5)$  such that

$$\mathbb{P}(X_i = j) = p_j, \quad j = 1, \dots, 5.$$

Then, the Likelihood function is proportional to

$$\prod_{i=1}^5 p_i^{y_i},$$

where  $y_i$  is the number of  $x_1, \dots, x_n$  equal to  $i$ . Consider the problem of testing

$$H_0 : p_1 = p_2 = p_3, p_4 = p_5 \quad H_1 : H_0 \text{ is not true.}$$

The parametric space is  $\Theta$  and it is a four-dimensional parametric space: indeed, since  $p_5 = 1 - \sum_{i=1}^4 p_i$ , there are only four free parameters. The parameter set is defined as follows:

$$\left\{ (p_1, \dots, p_5) \mid p_i \geq 0 \wedge \sum_{i=1}^4 p_i \leq 1 \right\},$$

a subset of  $\mathbb{R}^4$  containing an open subset of  $\mathbb{R}^4$ . Thus,  $q = 4$ . How many free parameters are there under the null hypotheses? Only one, because once  $p_1$  is fixed ( $0 \leq p_1 \leq 1/3$ ), then  $p_2 = p_3$  must be equal to  $p_1$  and  $p_4 = p_5$  must be equal to

$$\frac{1 - 3p_1}{2}.$$

Thus  $p_1 = 1$ , and the degrees of freedom is  $4 - 1 = 3$ . According to Wilks' Theorem we have that

$$-2 \log \lambda(\mathbf{x}) \rightarrow \chi_3^2.$$

Then, asymptotically, for  $n$  large, the test rejects  $H_0$  if

$$-2 \log \lambda(\mathbf{x}) \geq \chi_{3,\alpha}^2.$$

Wilks' Theorem is not the only approach for asymptotic test. A common method to construct large sample test statistic is based on estimators that are asymptotically Gaussian. In particular, the Maximum Likelihood estimators. Suppose that we want to test a hypothesis about a real valued parameter, and  $W_n = W(X_1, \dots, X_n)$  is a point estimator for  $\vartheta$ , for instance the Maximum Likelihood estimator. An approximate test, based on the Gaussian approximation, can be justified as follows: if  $\sigma^2$  denoted the variance of  $W_n$  and if we use some form of the Central Limit Theorem (or the Delta method) in order to show that, as  $n \rightarrow \infty$ ,

$$\frac{W_n - \vartheta}{\sigma_n} \rightarrow \mathcal{N}(0, 1),$$

then  $(W_n - \vartheta)/\sigma_n$  can be used as a term of comparison to the standard Gaussian. This implies that we have an approximate test.

## 10

### Week X

#### 10.1 27/11/18

##### 10.1.1 Linear Models

The technique of regression, in particular *linear* regression, probably wins as the most popular statistical tool. The major purpose of regression is to explore the dependence of one variable on another. What we did before was to consider a random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)$  from a parametric model  $Y \sim f_Y(y, \vartheta)$ . Now, we say that

$$Y_i = f(\vartheta, x) \quad \text{with } x \text{ fixed,}$$

and we want to make inference on  $\vartheta$ .

The first restriction that we introduce is the presence of a **response variable**  $y$  and one or more **explanatory variables**  $x_1, \dots, x_p$ . These variables,  $y$  and  $(x_1, \dots, x_p)$ , are treated asymmetrically, since the response variable is considered to vary according to variations of the explanatory variables.

*Example 10.1.* In medical statistics, the variable  $x$  could represent the dose of a drug which is given to individuals, and this is the role of the explanatory variable; the variable  $y$  could be, for instance, a binary response which denotes the effect of the drug.

Now, consider

$$y = f(x_1, \dots, x_p) :$$

for the moment being, we assume  $f$  to be linear (this is the reason why we call it *linear* model).

The second restriction is that the variable  $y$  is the sum of two terms, namely

$$y = f(x_1, \dots, x_p) + \varepsilon. \quad (10.1)$$

The term  $f(x_1, \dots, x_p)$  is called the **systematic component** and it is the contribution of  $x_1, \dots, x_p$  to  $y$ . The term  $\varepsilon$  is called the **erratic component**, and it is the random discrepancy between  $y$  and  $f(x_1, \dots, x_p)$ . This is *free* of any connection with  $x_1, \dots, x_p$  (meaning that they do not interact with each other).

**Definition 10.2.** An equation of the type (10.1) is called a **regression model**.

*Example 10.3.* Consider an electric circuit. Suppose we measure the voltage  $V$  starting when the circuit is closed. It is known that there exist the following relationship

$$V = E \left( 1 - \exp \left\{ -\frac{t}{T} \right\} \right),$$

where

- $E$  denotes the *voltage of the power supply*;
- $t$  is the *time since the closure of the circuit*;
- $T = RC$  is a constant characteristic value of the circuit depending on  $R$  (the *resistor*) and  $C$  (the *capacitor*).

In this framework,  $V$  is the response variable,  $t$  is the explanatory variable, whilst  $E$  and  $T$  are parameters which must be evaluated. Under closed circuit, we take some measurements at different points in time: what we get is a collection of various pairs of the form  $(t, V)$ .

Next, the error enters into the game: if we assume that the instruments that we use to take measurements are error-free, then the measurements are completely accurate, and the pairs  $(t_i, V_i)$ , for  $i = 1, \dots, n$ , are sufficient to determine the parameters. Therefore, we just have to solve a non-linear equation, with some constraints. However, in our case, we will consider the error's contribution, and we want to model its behaviour. In practice, the circuit and the instruments are not expected to be under ideal conditions. In this way we can model this uncertainty with

$$V = E \left( 1 - \exp \left\{ -\frac{t}{T} \right\} \right) + \varepsilon, \quad (10.2)$$

where  $\varepsilon$  includes all the sources of departure from the ideal behaviour. Since we do not know anything about  $\varepsilon$ , we assume it to be random. Therefore, (10.2) can not be solved as a non-linear equation, since  $V$  is random.

**Definition 10.4.** We consider a special case of regression model (10.1), where  $f(\cdot)$  is a linear function of the parameters,

$$y = \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon, \quad (10.3)$$

where  $\beta_i$  are unknown parameters. This model is called **linear regression model**.

This restriction on  $f(\cdot)$  could appear as a limitation of the model, and indeed it is. However, it turns out to be useful in a large number of real examples. Why? First, (10.3) is extremely simple. Second, even if there are cases where  $f(\cdot)$  is non-linear, they can be approximated by a linear relationship.

We consider  $n$  observations from the model (10.3): the objective is to make inference on the values of the parameters  $\beta_1, \dots, \beta_p$ . By inference, we mean *point estimation* and *testing*. Since the error component is considered to be generated at random,  $y$  is also regarded as a random quantity. Therefore, the first step would be to specify the distributional assumption on  $\varepsilon$ .

### 10.1.2 Second Order Assumptions of the Linear Regression Model

**Definition 10.5.** Suppose that  $n$  observations  $(y_1, \dots, y_n)$  are available, where each component is generated by (10.3). These observations are assumed to be a random sample from the random variables  $(Y_1, \dots, Y_n)$  that satisfy

$$Y_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $Y_i$  is the  $i$ -th component of the response variable,  $x_i$  is the  $i$ -th value of the explanatory variable  $X_i$ . Using the matrix notation for convenience, we can write

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where

- $\mathbf{Y} = {}^t(Y_1, \dots, Y_n)$  is the random vector with the  $n$  component of the response variable;
- $X = (x_{ij})$  is an  $n \times p$  matrix (with  $n \gg p$ ) called **regression matrix**, which contains the values of the  $p$  explanatory variables;
- $\boldsymbol{\varepsilon} = {}^t(\varepsilon_1, \dots, \varepsilon_n)$  is the vector containing the  $n$  components of the error term;
- $\boldsymbol{\beta} = {}^t(\beta_1, \dots, \beta_p)$  is the vector of the regression parameters.

Now, the **second order assumptions** are the following:

- $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$ ;
- $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n$ , for some unknown  $\sigma^2$ , implying that

$$\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \quad \text{and} \quad \text{Var}(\varepsilon_i) = \sigma^2, \quad \forall i \neq j;$$

- $X$  is a non-random matrix with full rank  $p$ .

We also refer to  $\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , with the above assumptions, as a **linear regression model**.

*Remark 10.6.* The linear regression model is particularly suited for the mathematical description of problems arising from *controlled experiments*, where the *experimenter* can control the values taken by relevant features to examine the corresponding value of the response variable. In this setting,  $X$  contains the values of the experimental features, which are non-stochastic since they are chosen by the *experimenter*. The error term is determined by the measuring errors and, if the instruments are not biased, then it follows that

$$\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}.$$

Finally, if the various experiments are conducted in such a way as not to influence each other, then the stochastic independence assumption is satisfied, and it implies uncorrelated error.

We aim to make inference on the set of  $(p + 1)$  parameters  $\beta_1, \dots, \beta_p, \sigma^2$ . By the second order assumptions, we have that

- $\mathbb{E}[\mathbf{Y}] = X\boldsymbol{\beta}$ ;
- $\text{Var}(\mathbf{Y}) = \sigma^2 I_n$ .

110 10 Week X

However, we do not have enough ingredients to write down the probability distribution of  $\mathbf{Y}$ . Therefore, we do not have the Likelihood function of the model for the observed  $\mathbf{y}$ . This means that we need more assumptions for  $\boldsymbol{\varepsilon}$  to introduce a probabilistic model for  $\mathbf{Y}$ . We must find a way to estimate  $\boldsymbol{\beta}$  by using only the second order assumptions. A reasonable proposal is to choose  $\boldsymbol{\beta}$  so as to minimise the *euclidean distance*

$$\|\mathbf{y} - \boldsymbol{\mu}\|$$

between the observed vector  $\mathbf{y}$  and its mean value  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{Y}] = X\boldsymbol{\beta}$ , that it to minimise

$$\|\mathbf{y} - X\boldsymbol{\beta}\|.$$

## 10.2 28/11/18

### 10.2.1 Least Square Method

We are looking for some estimate of the parameter  $\beta$  and  $\sigma^2$ . As we noted last time, the problem is to choose  $\beta$  such that the *euclidean distance*  $\|\mathbf{y} - \boldsymbol{\mu}\|$  is minimum. Equivalently, we can minimise the square of the distance; basically

$$Q(\beta) = \|\mathbf{y} - \boldsymbol{\mu}\|^2 = {}^t(\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu}) = {}^t(\mathbf{y} - X\beta)(\mathbf{y} - X\beta).$$

With this choice, we are establishing the **Least Square criterion**, whose name reminds that it operates by minimising the sum of squared differences between the component of  $\mathbf{y}$  and the corresponding component of  $\boldsymbol{\mu}$ .

Consider the vectors  $\mathbf{y}, \mathbf{x}_1, \dots, \mathbf{x}_n$  containing the response and the explanatory variables as elements of the vector space  $\mathbb{R}^n$ . As the parameter  $\beta$  varies in  $\mathbb{R}^p$ , the expression

$$X\beta = \beta_1\mathbf{x}_1 + \dots + \beta_p\mathbf{x}_p$$

is a linear combination of the column vectors  $\mathbf{x}_1, \dots, \mathbf{x}_p$  with coefficients  $\beta$ . In other terms, we have a parametric equation (in  $\beta$ ) of the linear subspace of  $\mathbb{R}^n$  generated by the columns of the design matrix. We denote this linear subspace by  $\mathcal{C}(X)$  (with dimension  $p$ ). In the model  $\mathbf{Y} = X\beta + \varepsilon$  we stated that  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{Y}]$  lies in  $\mathcal{C}(X)$  and the Least Squares criterion selects an element of  $\mathcal{C}(X)$  which minimises the *euclidean distance* between  $\mathbf{y}$  and the space  $\mathcal{C}(X)$ . We denote this element by  $\hat{\boldsymbol{\mu}} := X\hat{\beta}$ , which is uniquely identified by the coefficient  $\hat{\beta} \in \mathbb{R}^p$ .

In order to identify  $\hat{\beta}$ , we minimise the function  $Q(\beta)$  by making use of differentiation.

*Remark 10.7.* Remember that, given two vectors

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, \quad \mathbf{a} = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix}$$

such that  $\mathbf{y} = {}^t\mathbf{x}\mathbf{a}$ ,

$$\frac{d}{d\mathbf{x}}({}^t\mathbf{x}\mathbf{a}) = {}^t\left(\frac{d}{dx_1}({}^t\mathbf{x}\mathbf{a}), \dots, \frac{d}{dx_n}({}^t\mathbf{x}\mathbf{a})\right) = \mathbf{a}, \quad (10.4)$$

$$\frac{d}{d{}^t\mathbf{x}}({}^t\mathbf{x}\mathbf{a}) = \left(\frac{d}{dx_1}({}^t\mathbf{x}\mathbf{a}), \dots, \frac{d}{dx_n}({}^t\mathbf{x}\mathbf{a})\right) = \mathbf{a}, \quad (10.5)$$

$$\frac{d}{d\mathbf{x}}A\mathbf{x} = {}^tA,$$

$$\frac{d}{d\mathbf{x}}{}^t\mathbf{x}B\mathbf{x} = 2B\mathbf{x},$$

where  $B$  has to be symmetric (note that in (10.4) we considered the transpose since  $\mathbf{x}$  is a column vector, and in (10.5) we did not since  ${}^t\mathbf{x}$  is a row vector).

By using these two rules, we finally differentiate  $Q(\beta)$  with respect to  $\beta$ :

$$\begin{aligned}\frac{d}{d\beta}Q(\beta) &= \frac{d}{d\beta} {}^t(\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu}) \\ &= \frac{d}{d\beta} ({}^t\mathbf{y}\mathbf{y} - 2{}^t\mathbf{y}X\boldsymbol{\beta} + {}^t\boldsymbol{\beta}{}^tXX\boldsymbol{\beta}) \\ &= \frac{d}{d\beta} (-2{}^t\mathbf{y}X\boldsymbol{\beta}) + \frac{d}{d\beta} ({}^t\boldsymbol{\beta}{}^tXX\boldsymbol{\beta}) \\ &= -2{}^tX\mathbf{y} + 2{}^tXX\boldsymbol{\beta} \\ &= 2({}^tXX\boldsymbol{\beta} - {}^tX\mathbf{y}).\end{aligned}$$

If we equate the previous expression to zero, then we deduce that the required minimum point satisfies

$${}^tXX\boldsymbol{\beta} = {}^tX\mathbf{y},$$

also known as **normal equation**. Observe that the inversion of the matrix  ${}^tXX$  is legitimate, since we are assuming the condition on the rank of  $X$ : recall that if  $\text{rk } X = p$ , then  $\text{rk } {}^tXX = p$ . This implies that a point which minimises  $Q(\beta)$  is

$$\hat{\boldsymbol{\beta}} = ({}^tXX)^{-1} {}^tX\mathbf{y}. \quad (10.6)$$

The estimates  $\hat{\boldsymbol{\beta}}$  are called **regression coefficients**. To check that (10.6) gives a minimum point for  $Q(\beta)$  we simply consider the matrix of the second derivatives

$$\frac{d}{d\beta} \frac{d}{d\beta} Q(\beta) = 2{}^tXX,$$

which is a positive definite matrix (therefore  $Q(\beta)$  is convex).

### 10.2.2 The Projection Matrix Operator

**Definition 10.8.** Associated to  $\hat{\boldsymbol{\beta}}$  there is the projection vector  $\mathbf{y}$  onto  $\mathcal{C}(X)$ , which is called the **vector of fitted values**, that is

$$\hat{\boldsymbol{\mu}} = X\hat{\boldsymbol{\beta}} = X({}^tXX)^{-1} {}^tX\mathbf{y} = P\mathbf{y},$$

where  $P := X({}^tXX)^{-1} {}^tX$  is the **projection matrix operator** (or simply **projection matrix**) on  $\mathcal{C}(X)$ . The projection matrix defines another operator associated with the matrix  $X$  whose role is to project any vector  $\mathbf{y} \in \mathbb{R}^n$  and transform it into another vector, that we call  $P\mathbf{y} \in \mathcal{C}(X)$ : this minimises the distance from  $\mathbf{y}$ .

*Remark 10.9.*  $P$  has some interesting properties, that we will often use.

- $P$  is symmetric.
- $P$  is idempotent, that is  $P^2\mathbf{y} = P\mathbf{y}$ .
- $\text{rk } P = \text{tr } P = \text{tr}({}^tXX)^{-1} {}^tXX = p$ .

We can decompose  $\mathbf{y}$  into two terms:

1. the projection  $\hat{\boldsymbol{\mu}}$  onto  $\mathcal{C}(X)$ ;

2. the residual component

$$\mathbf{y} - \hat{\boldsymbol{\mu}} = \mathbf{y} - X({}^tXX)^{-1} {}^tX\mathbf{y} = (I_n - P)\mathbf{y}.$$

*Remark 10.10.* The equation  ${}^tXX\boldsymbol{\beta} = {}^tX\mathbf{y}$  can be also obtained by a simple geometric argument. Indeed, the vector  $\hat{\boldsymbol{\mu}}$  which minimises the distance from  $\mathbf{y}$  is such that

$$(\hat{\mathbf{y}} - X\hat{\boldsymbol{\beta}}) \perp \mathcal{C}(X),$$

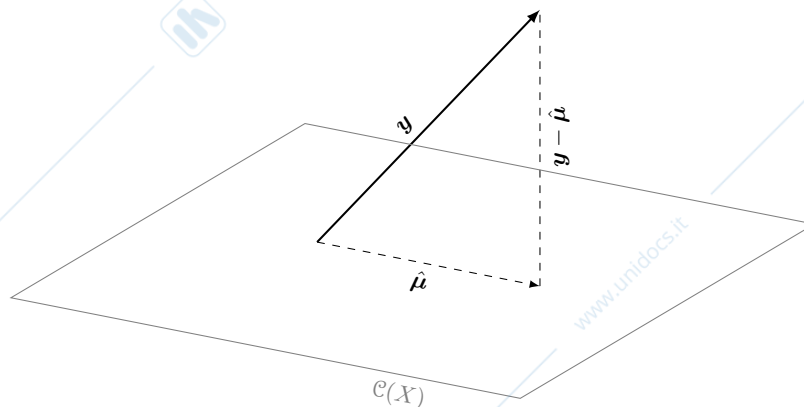
In particular, we have that  $(\mathbf{y} - \hat{\boldsymbol{\mu}})$  is orthogonal to all the elements in the subspace  $\mathcal{C}(X)$ . Indeed, for every  $X\mathbf{a} \in \mathcal{C}(X)$ , we can write

$$\begin{aligned} {}^t(X\mathbf{a})(\mathbf{y} - \hat{\boldsymbol{\mu}}) &= {}^t(X\mathbf{a})(\mathbf{y} - P\mathbf{y}) \\ &= {}^t\mathbf{a} {}^tX(\mathbf{y} - X({}^tXX)^{-1} {}^tX\mathbf{y}) \\ &= {}^t\mathbf{a} {}^tX\mathbf{y} - {}^t\mathbf{a} {}^tXX({}^tXX)^{-1} {}^tX\mathbf{y} \\ &= {}^t\mathbf{a} {}^tX\mathbf{y} - {}^t\mathbf{a} {}^tX\mathbf{y} \\ &= 0. \end{aligned}$$

In other terms, we need that

$${}^t(\mathbf{y} - X\hat{\boldsymbol{\beta}})X = 0,$$

which is equivalent to the normal equation.



**Fig. 10.1.** Graphical representation of the projection onto  $\mathcal{C}(X)$ .

### 10.2.3 Estimation of $\hat{\boldsymbol{\beta}}$

Consider the estimates  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$ . Can we say something about  $\hat{\boldsymbol{\beta}}$ ? For instance, we have that

$$\begin{aligned} \mathbb{E}[\hat{\boldsymbol{\beta}}] &= \mathbb{E}[({}^tXX)^{-1} {}^tX\mathbf{y}] \\ &= ({}^tXX)^{-1} {}^tX\mathbb{E}[\mathbf{y}] \\ &= ({}^tXX)^{-1} {}^tXX\boldsymbol{\beta} \end{aligned}$$

114 10 Week X

$$= \beta.$$

The Least Squares estimator  $\hat{\beta}$  is **unbiased** for  $\beta$ . Moreover, this implies that

$$\mathbb{E}[\hat{\mu}] = \mathbb{E}[X\hat{\beta}] = X\mathbb{E}[\hat{\beta}] = X\beta = \mu.$$

The dispersion matrix of the estimates is

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \text{Var}(({}^tXX)^{-1}{}^tX\mathbf{y}) \\ &= ({}^tXX)^{-1}{}^tX \text{Var}(\mathbf{y}) ({}^tXX)^{-1}{}^tX \\ &= ({}^tXX)^{-1}{}^tX \sigma^2 I_n ({}^tXX)^{-1}{}^tX \\ &= \sigma^2 I_n ({}^tXX)^{-1}{}^tXX ({}^tXX)^{-1} \\ &= \sigma^2 ({}^tXX)^{-1} \end{aligned}$$

Also,

$$\begin{aligned} \text{Var}(\hat{\mu}) &= X \text{Var}(\hat{\beta}) X \\ &= \sigma^2 X ({}^tXX)^{-1} X \\ &= \sigma^2 P \end{aligned}$$

We can also check consistency. By the above variances, consistency holds if and only if the diagonal element of  $V = ({}^tXX)^{-1}$  converges to 0 as  $n \rightarrow \infty$ . This is determined by the design matrix  $X$ , whose number of rows must increase as  $n \rightarrow \infty$ . In general, it is impossible to assess consistency of  $\hat{\beta}$  in general terms without making some assumptions over the designed matrix.

*Remark 10.11.* So far, we have regarded the matrix  $X$  as a non-stochastic matrix. However, it is very frequently the case that the response and explanatory variables are determined simultaneously. Therefore, we deal with stochastic explanatory variables. The most typical example occurs when two or more variables are observed on the same statistical unit, such as weight ( $\mathbf{y}$ ) and height ( $\mathbf{x}$ ).

The assumption of a non-stochastic matrix  $X$  is supported by the following arguments. In most cases, the distribution of the explanatory variables does not contain any information on the relationship with the response variable. With a regression model, we are interested in making inference on these relationships, and not on the distribution of the explanatory variables. In some way, we examine the variables conditionally on the values taken by  $X$ .

### 10.3 29/11/18

#### 10.3.1 Estimation of $\sigma^2$

So far, we considered the estimation of  $\beta$ . However, we are also interested in estimating  $\sigma^2$ , but the Least Square criterion does not say how to estimate  $\sigma^2$ . Since the generic term  $\varepsilon_i$  is such that  $\text{Var}(\varepsilon_i) = \sigma^2$ , it is reasonable to estimate  $\sigma^2$  by the arithmetic mean of  $\hat{\varepsilon}_i$ , where  $\hat{\varepsilon}_i$  is the generic component of the residual vector

$$\hat{\varepsilon} = \mathbf{y} - \hat{\boldsymbol{\mu}}.$$

Therefore, we can consider

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{\|\hat{\varepsilon}\|^2}{n} \quad (10.7)$$

as an estimate of  $\sigma^2$ . Now,

$$\begin{aligned} \|\hat{\varepsilon}\|^2 &= Q(\hat{\boldsymbol{\beta}}) \\ &= {}^t(\mathbf{y} - \hat{\boldsymbol{\mu}})(\mathbf{y} - \hat{\boldsymbol{\mu}}) \\ &= {}^t\mathbf{y} {}^t(I_n - P)(I_n - P)\mathbf{y} \\ &= {}^t\mathbf{y}(I_n - P)\mathbf{y}, \end{aligned}$$

since  $(I_n - P)$  is idempotent. We can compute the expectation of  $\hat{\sigma}^2$  as follows, using (10.7):

$$\mathbb{E}[n\hat{\sigma}^2] = \mathbb{E}[{}^t\mathbf{y}(I_n - P)\mathbf{y}].$$

In order to do the computation, recall Lemma 1.22 in our case,

$$\mathbb{E}[{}^t\mathbf{y}(I_n - P)\mathbf{y}] = {}^t\boldsymbol{\mu}(I_n - P)\boldsymbol{\mu} + \text{tr}((I_n - P)\sigma^2 I_n) :$$

now, the first addend in the right-hand side vanishes, since  $(I_n - P)$  projects onto the orthogonal space  $\mathcal{C}(X)$  which contains  $\boldsymbol{\mu}$ , and hence

$$(I_n - P)\boldsymbol{\mu} = (I_n - X({}^tXX)^{-1}{}^tX)X\boldsymbol{\beta} = 0.$$

This implies that

$$\mathbb{E}[{}^t\mathbf{y}(I_n - P)\mathbf{y}] = \sigma^2(n - p) :$$

we proved that  $\mathbb{E}[n\hat{\sigma}^2] = \sigma^2(n - p)$ , which implies that the estimator  $\hat{\sigma}^2$  is *biased* for  $\sigma^2$ , but it is *asymptotically unbiased* for  $\sigma^2$ .

*Remark 10.12.* We could easily find an unbiased estimator for  $\sigma^2$  by using the relation

$$\mathbb{E}[n\hat{\sigma}^2] = (n - p)\sigma^2.$$

This estimator is

$$\hat{\sigma}^2 \frac{n}{n - p} = \frac{Q(\hat{\boldsymbol{\beta}})}{n - p}.$$

*Example 10.13.* The simplest case that we can consider has  $p = 1$  and  $\mathbf{X} = \mathbf{1}_n$ . Then, the linear model is the following:

$$\mathbf{Y} = \mathbf{1}_n \beta + \boldsymbol{\varepsilon},$$

where  $\beta$  is a scalar. This is just an unusual way of introducing a model where all components of  $\mathbf{Y}$  have the same mean and the same variance, and are uncorrelated. In this scenario,

$$({}^t\mathbf{X}\mathbf{X})^{-1} = \frac{1}{n}, \quad {}^t\mathbf{X}\mathbf{y} = \sum_{i=1}^n y_i, \quad \hat{\beta} = \frac{1}{n} \sum_{i=1}^n y_i.$$

*Example 10.14.* We consider  $p = 2$ , with the first column of  $\mathbf{X}$  equal to  $\mathbf{1}_n$ , and the second column of  $\mathbf{X}$  equal to  $\mathbf{x} = (x_1, \dots, x_n)$ , meaning

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}.$$

In this situation, the linear regression model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \beta_1 + \beta_2 \mathbf{x} + \boldsymbol{\varepsilon}$$

is called **simple linear regression model**. Here,

$$\begin{aligned} {}^t\mathbf{X}\mathbf{X} &= \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix} \\ ({}^t\mathbf{X}\mathbf{X})^{-1} &= \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix} \\ {}^t\mathbf{X}\mathbf{y} &= \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix} \\ \hat{\boldsymbol{\beta}} &= \begin{pmatrix} \bar{y} - \frac{s_{xy}}{s_{xx}} \bar{x} \\ \frac{s_{xy}}{s_{xx}} \end{pmatrix} \end{aligned} \quad (10.8)$$

where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad s_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2, \quad s_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

How did we obtain 10.8? Indeed,

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= ({}^t\mathbf{X}\mathbf{X})^{-1} ({}^t\mathbf{X}\mathbf{y}) \\ &= \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix} \cdot \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix} \\ &= \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i \\ n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i \end{pmatrix}. \end{aligned}$$

In order to show that this is equal to

$$\left( \bar{y} - \frac{s_{xy}\bar{x}}{s_{xx}} \right)$$

first notice that

$$\begin{aligned} s_{xx} &= \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \sum_{i=1}^n (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \\ &= \sum_{i=1}^n x_i^2 - 2 \left( \sum_{i=1}^n x_i \right) \left( \frac{1}{n} \sum_{i=1}^n x_i \right) + n \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 \\ &= \sum_{i=1}^n x_i^2 - \frac{2}{n} \left( \sum_{i=1}^n x_i \right)^2 + \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \\ &= \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \\ &= \frac{1}{n}, \end{aligned}$$

while

$$\begin{aligned} s_{xy} &= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \sum_{i=1}^n (x_i y_i - x_i \bar{y} - \bar{x} y_i + \bar{x} \bar{y}) \\ &= \sum_{i=1}^n x_i y_i - \left( \sum_{i=1}^n x_i \right) \left( \frac{1}{n} \sum_{i=1}^n y_i \right) - \left( \frac{1}{n} \sum_{i=1}^n x_i \right) \left( \sum_{i=1}^n y_i \right) + \left( \frac{1}{n} \sum_{i=1}^n x_i \right) \left( \frac{1}{n} \sum_{i=1}^n y_i \right) \\ &= \frac{1}{n}. \end{aligned}$$

Finally,

$$\begin{aligned} \bar{y} - \frac{s_{xy}\bar{x}}{s_{xx}} &= \frac{1}{n} \sum_{i=1}^n y_i - \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \left( \frac{1}{n} \sum_{i=1}^n x_i \right) \\ &= \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n y_i \left( \sum_{i=1}^n x_i \right)^2 - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i + \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \\ &= \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}. \end{aligned}$$

The expression for  $\hat{\beta}$  simplifies if  $\bar{x} = 0$ . This condition can always be satisfied by changing the parametrisation from

$$\mathbf{y} = \mathbf{1}_n \beta_1 + x \beta_2 + \boldsymbol{\varepsilon}$$

118 10 Week X

to

$$\mathbf{y} = \mathbf{1}_n(\beta_1 + \beta_2\bar{x}) + (\mathbf{x} - \mathbf{1}_n\bar{x})\beta_2 + \boldsymbol{\varepsilon}$$

which is an equation that has

- *intercept*  $\alpha = \beta_1 + \beta_2\bar{x}$ ;
- *slope*  $\beta_2$ .

Under the new parametrisation, the new explanatory variables  $\mathbf{z} = \mathbf{x} - \mathbf{1}_n\bar{x}$  has  $\sum_{i=1}^n z_i = 0$ .

### 10.3.2 Sample Correlation and Coefficient of Determination

The *sample correlation* is defined as

$$r := \frac{s_{xy}}{\sqrt{s_{xx}s_{yy}}},$$

where

$$s_{yy} := \sum_{i=1}^n (y_i - \bar{y})^2.$$

This quantity is typically used as a descriptive statistic to quantify the degree of closeness of the observed data points to the fitted line. In fact,  $r \in [-1, 1]$ ; if  $r = \pm 1$ , this means that the fitted point lie exactly on the regression line. On the other hand,  $r = 0$  denotes a poor fit of the regression line to the points.

A related quantity is the *coefficient of determination*, which is defined as a function of the sample correlation

$$r^2 := \frac{s_{xy}^2}{s_{xx}s_{yy}},$$

which represents the fraction of variability of the  $y_i$ 's that can be imputed to the linear component of the relationship.

*Remark 10.15.* One of the second order hypotheses is that the error component  $\boldsymbol{\varepsilon}$  must have  $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$ . In a large number of examples, this assumption does not work, but at the same time it does not represent an issue. In particular, suppose that one column of the design matrix  $X$ , say the first, is  $\mathbf{1}_n$ . Denote by  $\mu_{\varepsilon_i}$  the mean value of  $\varepsilon_i$ , and by  $\bar{X}$  and  $\bar{\boldsymbol{\beta}}$  the remaining column of  $X$  and the remaining rows of  $\boldsymbol{\beta}$  after deleting the first. Then, we can write our linear model

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

as

$$\mathbf{Y} = \mathbf{1}_n\beta_1 + \bar{X}\bar{\boldsymbol{\beta}} + \boldsymbol{\varepsilon} = \mathbf{1}_n(\beta_1 + \mu_{\varepsilon_i}) + \bar{X}\bar{\boldsymbol{\beta}} + (\boldsymbol{\varepsilon} - \mathbf{1}_n\mu_{\varepsilon_i}),$$

where the error component is  $\boldsymbol{\varepsilon} - \mathbf{1}_n\mu_{\varepsilon_i}$  which has expected value equal to zero. Therefore, **the assumption on the mean value of  $\boldsymbol{\varepsilon}$  is not a relevant assumption.**

**Theorem 10.16 (Gauss-Markov).** *Under the second order assumptions for the linear model  $\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , the Least Squares estimator  $\hat{\boldsymbol{\beta}}$  is such that*

$$\text{Var}(\hat{\boldsymbol{\beta}}) \leq \text{Var}(\mathbf{T}),$$

where  $T$  is any estimator of the form

$$\mathbf{T} = {}^t C \mathbf{Y}$$

and  $C$  is an  $n \times p$  matrix such that  $\mathbb{E}[\mathbf{T}] = \boldsymbol{\beta}$ , for every  $\boldsymbol{\beta}$ .

*Proof.* Since the equality

$$\mathbb{E}[{}^t C \mathbf{Y}] = {}^t C X \boldsymbol{\beta} = \boldsymbol{\beta}$$

for every  $\boldsymbol{\beta}$ , it follows that

$${}^t C X = I_p = {}^t X C.$$

Now, recall that the variance of  $\hat{\boldsymbol{\beta}}$  is

$$\text{Var}(\hat{\boldsymbol{\beta}}) = \sigma^2 ({}^t X X)^{-1}$$

and since

$$\text{Var}(\mathbf{T}) = \text{Var}({}^t C \mathbf{Y}) = C \text{Var}(\mathbf{Y}) {}^t C = \sigma^2 {}^t C C,$$

we have that the statement of the Theorem is equivalent to prove that

$${}^t C C - ({}^t X X)^{-1} \geq 0.$$

Since  ${}^t C X = I_p = {}^t X C$ , we must show

$$\begin{aligned} {}^t C C - {}^t C X ({}^t X X)^{-1} {}^t X C &= {}^t C [C - X ({}^t X X)^{-1} {}^t X C] \\ &= {}^t C (I_n - P) C \\ &\geq 0. \end{aligned}$$

Indeed, we have

$$\begin{aligned} {}^t \mathbf{a} {}^t C (I_n - P) C \mathbf{a} &= {}^t \mathbf{a} {}^t C (I_n - P) (I_n - P) C \mathbf{a} \\ &= \|(I_n - P) C \mathbf{a}\|^2 \\ &\geq 0, \quad \forall \mathbf{a} \in \mathbb{R}^p. \end{aligned}$$

We proved that in the class of *linear unbiased estimators* for the parameters of a linear model  $\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , the Least Squares estimator  $({}^t X X)^{-1} {}^t X \mathbf{y}$  is the one with the minimum variance.

□



## 11

## Week XI

## 11.1 04/12/18

*Remark 11.1.* An implication of the Gauss-Markov Theorem is the following: consider a linear combination of the parameter  $\vartheta$ , usually known as *contrast*, of the form

$$\psi = {}^t\mathbf{a}\boldsymbol{\beta},$$

where  $\mathbf{a} \in \mathbb{R}^p$  is a vector of given constants. It follows that the estimator  $\hat{\psi} = {}^t\mathbf{a}\hat{\boldsymbol{\beta}}$  is unbiased for  $\psi$  and its variance is the minimum achievable variance among linear estimators.

In particular, if we choose a vector  $\mathbf{a}$  with null components except the component of order  $j$ -th, we have that  $\hat{\psi} = \hat{\beta}_j$  is the best linear estimator for the component  $j$ -th.

*Remark 11.2.* The Gauss-Markov Theorem states that, in the class of unbiased linear estimators of  $\boldsymbol{\beta}$ , the Least Squares estimator has the minimum variance. Therefore,

$$\hat{\boldsymbol{\beta}} = ({}^tXX)^{-1} {}^tX\mathbf{y}$$

is the best linear estimator. This provides with the (unique) theoretical guarantee for the Least Squares approach.

## 11.1.1 Constrained Estimates

## Sums-of-squares Decomposition

We discussed about the orthogonality of  $P\mathbf{y}$  and of the residual vector. If we write the squared norm of  $\hat{\boldsymbol{\mu}} + (\mathbf{y} - \hat{\boldsymbol{\mu}})$ , then by orthogonality

$$\|\mathbf{y}\|^2 = \|\hat{\boldsymbol{\mu}}\|^2 + \|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2.$$

Now, if we replace  $\mathbf{y}$  with the corresponding random variable, we have that

$$\|\mathbf{Y}\|^2 = \|P\mathbf{Y}\|^2 + \|(I_n - P)\mathbf{Y}\|^2,$$

or alternatively

$${}^t\mathbf{Y}I_n\mathbf{Y} = {}^t\mathbf{Y}P\mathbf{Y} + {}^t\mathbf{Y}(I_n - P)\mathbf{Y}.$$

These are quadratic forms with associated (idempotent) matrices

$$I_n, \quad P, \quad (I_n - P),$$

and such that

$$\text{rk } I_n = n, \quad \text{rk } P = p, \quad \text{rk}(I_n - P) = n - p.$$

The above quadratic form (on the above square norms) are called

- $\|\mathbf{Y}\|^2$  *total sum of squares*;
- $\|P\mathbf{Y}\|^2$  *regression sum of squares*;
- $\|(I_n - P)\mathbf{Y}\|^2$  *residual sum of squares*.

It is important to study the distribution of these quantities. However, this study is possible only if we specify a distributional assumption on  $\boldsymbol{\varepsilon}$ . For the moment, with the second order assumptions, we can only compute the expectation of this quadratic form by applying Lemma [1.22](#):

$$\begin{aligned} \mathbb{E} \left[ \|P\mathbf{Y}\|^2 \right] &= p\sigma^2 + {}^t\boldsymbol{\mu}P\boldsymbol{\mu}, \\ \mathbb{E} \left[ \|(I_n - P)\mathbf{Y}\|^2 \right] &= (n - p)\sigma^2, \quad \text{since } {}^t\boldsymbol{\mu}(I_n - P)\boldsymbol{\mu} = 0 \\ \mathbb{E} \left[ \|\mathbf{Y}\|^2 \right] &= n\sigma^2 + {}^t\boldsymbol{\mu}\boldsymbol{\mu}. \end{aligned}$$

### 11.1.2 Constrained Optimisation: Lagrange Multipliers

Consider the following problem: we want to estimate  $\boldsymbol{\beta}$  in the linear model when there exist some linear constraints among the components of  $\boldsymbol{\beta}$ . For instance, imagine that the vector  $\boldsymbol{\beta}$  is such that

$$H\boldsymbol{\beta} = \mathbf{0},$$

where  $H$  is a  $q \times p$  ( $q \ll p$ ) matrix of rank  $q$ , whose elements are given constants. The solution of this constrained estimation problem is particularly useful in the contest of hypotheses testing of the parameter  $\boldsymbol{\beta}$ . With respect to the Least Squares approach, we are dealing with the additional condition

$$H\boldsymbol{\beta} = \mathbf{0}.$$

This condition constraints  $\boldsymbol{\beta}$  to stay in a certain subset of the linear space  $\mathcal{C}(X)$ , say  $\mathcal{C}_0(X)$ . It is easy to show that  $\mathcal{C}_0(X)$  is a linear subspace of  $\mathcal{C}(X)$ , with dimension  $p - q$ .

We can solve the problem of constrained optimisation by using *Lagrange Multipliers*. In other terms, we want to minimise the function  $Q(\boldsymbol{\beta})$  under  $H\boldsymbol{\beta} = \mathbf{0}$ . Basically, we must minimise the function

$$f(\boldsymbol{\alpha}, \boldsymbol{\beta}) := {}^t(\mathbf{y} - X\boldsymbol{\beta})(\mathbf{y} - X\boldsymbol{\beta}) + 2{}^t(H\boldsymbol{\beta})\boldsymbol{\alpha},$$

where  $\boldsymbol{\alpha}$  is a vector of Lagrange multipliers. By differentiating the function  $f$  with respect to  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$ , and equating these expressions to zero, we obtain the following system of equations:

$$\begin{cases} \frac{\partial}{\partial \beta} f(\alpha, \beta) = {}^t X X \beta + {}^t H \alpha = {}^t X \mathbf{y} \\ \frac{\partial}{\partial \alpha} f(\alpha, \beta) = H \beta = \mathbf{0} \end{cases}$$

Consider the first equation and multiply it by  $({}^t X X)^{-1}$ . Then, we can isolate  $\beta$ :

$$\beta = ({}^t X X)^{-1} {}^t X \mathbf{y} - ({}^t X X)^{-1} {}^t H \alpha = \hat{\beta} - ({}^t X X)^{-1} {}^t H \alpha \quad (11.1)$$

and (11.1) leads, multiplying by  $H$ , to

$$0 = H \beta = H \hat{\beta} - H ({}^t X X)^{-1} {}^t H \alpha,$$

thus isolating  $\alpha$  we obtain

$$\alpha = [H ({}^t X X)^{-1} {}^t H]^{-1} H \hat{\beta}.$$

Therefore, the minimum of the function  $f$  is obtained for  $\beta$  equal to

$$\hat{\beta}_0 = \hat{\beta} - ({}^t X X)^{-1} {}^t H K H \hat{\beta},$$

where  $K := [H ({}^t X X)^{-1} {}^t H]^{-1}$ . The vector  $\hat{\beta}_0$  is the estimator of  $\beta$  under  $H \beta = \mathbf{0}$ . This implies that we have a fitted value with respect to  $\hat{\beta}_0$  which lives in  $\mathcal{C}_0(X)$ . The fitted value is obtained as follows:

$$\begin{aligned} \hat{\mu}_0 &= X \hat{\beta}_0 \\ &= X [\hat{\beta} - ({}^t X X)^{-1} {}^t H K H \hat{\beta}] \\ &= \hat{\mu} - X ({}^t X X)^{-1} {}^t H K H \hat{\mu} \\ &= P \mathbf{y} - X ({}^t X X)^{-1} {}^t H K H ({}^t X X)^{-1} {}^t X \mathbf{y} \\ &= (P - P_H) \mathbf{y} \\ &= P_0 \mathbf{y}, \end{aligned}$$

where

$$P_0 := P - P_H, \quad P_H := X ({}^t X X)^{-1} {}^t H K H ({}^t X X)^{-1} {}^t X.$$

We obtained a new projection matrix  $P_0$  that applies to  $\mathbf{y}$  in order to obtain the fitted value  $\hat{\mu}_0$  that minimises the distance with respect to  $\mathbf{y}$  and lives in  $\mathcal{C}_0(X)$ .

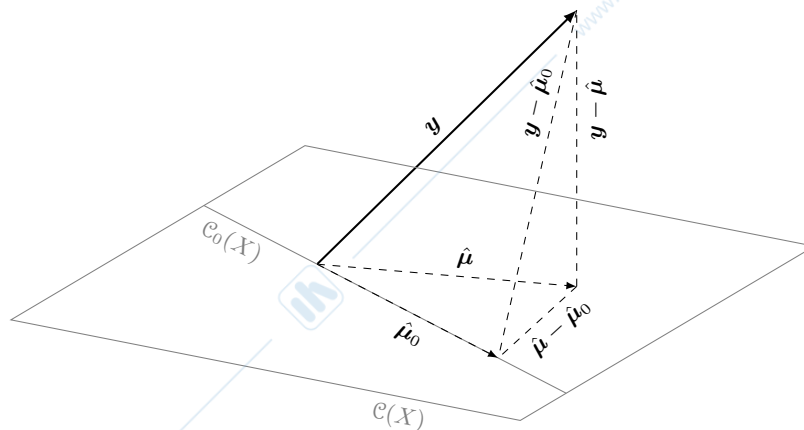


Fig. 11.1. Graphical representation of the projection onto  $\mathcal{C}(X)$ , taking  $\mathcal{C}_0(X)$  into account.

### Properties

We obtained a new projection matrix  $P_0$  that projects any vector of  $\mathbb{R}^n$  into the linear subspace  $\mathcal{C}_0(X)$ . The following properties of the projection hold.

1. The vector  $\mathbf{y} - \hat{\boldsymbol{\mu}}_0$  is orthogonal to all the elements of  $\mathcal{C}_0(X)$ . Indeed, consider  $\mathbf{c}$  to be any element of  $\mathcal{C}_0(X)$  such that  $H\mathbf{c} = \mathbf{0}$ ; then, we have

$${}^t(\mathbf{y} - \hat{\boldsymbol{\mu}}_0)X\mathbf{c} = 0,$$

and in particular

$$(\mathbf{y} - \hat{\boldsymbol{\mu}}) \perp \hat{\boldsymbol{\mu}}_0.$$

2. The projection of  $\mathbf{y} - \hat{\boldsymbol{\mu}}_0$  onto  $\mathcal{C}(X)$  is

$$P(\mathbf{y} - \hat{\boldsymbol{\mu}}_0) = \hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}_0,$$

such that

$$\hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}_0 \perp \hat{\boldsymbol{\mu}}_0.$$

We can have a decomposition of  $\mathbf{y}$  which is similar to the one obtained without constraint on  $\boldsymbol{\beta}$ :

$$\mathbf{y} = \hat{\boldsymbol{\mu}}_0 + (\hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}_0) + (\mathbf{y} - \hat{\boldsymbol{\mu}}),$$

where the three elements on the right-hand side are orthogonal and then we have an even more general version of Pitagora's Theorem:

$$\|\mathbf{y}\|^2 = \|\hat{\boldsymbol{\mu}}_0\|^2 + \|\hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}_0\|^2 + \|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2,$$

which extends the one used above.

#### 11.1.3 Gaussian Theory

So far, we considered the linear model under the second order assumptions on  $\boldsymbol{\varepsilon}$ . However, it is natural to specify a distributional assumption on  $\boldsymbol{\varepsilon}$ , and then make use of the theory developed in the first part of this course. The most natural (and simplest) assumption is to consider  $\boldsymbol{\varepsilon}$  to be Gaussian. In particular, we assume that

$$\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 I_n),$$

therefore

$$\mathbf{Y} \sim \mathcal{N}_n(X\boldsymbol{\beta}, \sigma^2 I_n).$$

We have a distributional assumption for  $\mathbf{Y}$ , that is, we have a parametric model for the observable values. Let  $\boldsymbol{\vartheta} = ({}^t\boldsymbol{\beta}, \sigma^2)$ . Then, the Log-Likelihood function is

$$\begin{aligned} \log \mathcal{L}(\boldsymbol{\vartheta}, \mathbf{y}) &= -\frac{1}{2}n \log \sigma^2 - \frac{1}{2}\sigma^{-2} \|\mathbf{y} - X\boldsymbol{\beta}\|^2 \\ &= -\frac{1}{2}n \log \sigma^2 - \frac{1}{2}\sigma^{-2} ({}^t\mathbf{y}\mathbf{y} - 2{}^t\mathbf{y}X\boldsymbol{\beta} + {}^t\boldsymbol{\beta}{}^tXX\boldsymbol{\beta}). \end{aligned}$$

This Log-Likelihood function has a clear exponential structure, such that we can identify the sufficient statistic:

$${}^t(\mathbf{y}X, \mathbf{y}\mathbf{y}),$$

which is a vector of dimension  $p+1$ . In the Log-Likelihood function, we can identify the function  $Q(\boldsymbol{\beta})$ , which is the same quantity used by the Least Squares approach. This implies that, for  $\sigma^2$  fixed, the maximum of the Log-Likelihood function coincides with the minimum of  $Q(\boldsymbol{\beta})$ . We conclude that **the Maximum Likelihood estimator of  $\boldsymbol{\beta}$  coincides with the Least Squares estimator of  $\boldsymbol{\beta}$** ,

$$\hat{\boldsymbol{\beta}} = ({}^tXX)^{-1} {}^tX\mathbf{y}.$$

This happens because we assumed that  $\boldsymbol{\varepsilon}$  is a multi-dimensional Gaussian, and then  $\mathbf{Y}$  is a multi-dimensional Gaussian. We can also maximise the Log-Likelihood function with respect to  $\sigma^2$ : if we did that, we obtain that the Maximum Likelihood estimator for  $\sigma^2$  is

$$\hat{\sigma}^2 = \frac{\|\hat{\boldsymbol{\varepsilon}}\|^2}{n}.$$

Recall that  $\hat{\sigma}^2$  is a biased estimator for  $\sigma^2$ . We can fix this issue and consider the unbiased estimator

$$\hat{\sigma}^2 \frac{n}{n-p}.$$

## 11.2 05/12/18

## 11.2.1 The Fisher-Cochran Theorem

**Theorem 11.3.** Consider  $\mathbf{Y} = (Y_1, \dots, Y_k) \sim \mathcal{N}_k(\boldsymbol{\mu}, I_k)$  and let  $A_1, \dots, A_m$  be positive semidefinite matrices of ranks  $r_1, \dots, r_m$ , respectively. Assume that

$$\sum_{i=1}^m A_i = I_k.$$

Then, the following are equivalent:

1. the quadratic form

$$Q_j = {}^t\mathbf{Y}A_j\mathbf{Y}, \quad j = 1, \dots, m$$

are distributed as non-centred Chi-squared distributions, with degrees of freedom  $r_j$  and parameter of non-centrality  ${}^t\boldsymbol{\mu}A_j\boldsymbol{\mu}$ . Moreover,  $Q_i$  is independent of  $Q_j$ , for  $j \neq i$ ;

2. it holds that

$$r_1 + \dots + r_m = k.$$

*Proof.* (1)  $\Rightarrow$  (2) We have that

$${}^t\mathbf{Y}\mathbf{Y} = {}^t\mathbf{Y}(A_1 + \dots + A_m)\mathbf{Y} = Q_1 + \dots + Q_m \sim \chi_r^2(\delta),$$

where

$$r = \sum_{j=1}^m r_j \quad \text{and} \quad \delta = {}^t\boldsymbol{\mu} \left( \sum_{j=1}^m A_j \right) \boldsymbol{\mu} = {}^t\boldsymbol{\mu}\boldsymbol{\mu},$$

$\delta$  also known as the *non-centrality parameter*. Therefore,  ${}^t\mathbf{Y}\mathbf{Y} \sim \chi_k^2({}^t\boldsymbol{\mu}\boldsymbol{\mu})$ , meaning that  $r = k$ .

(2)  $\Rightarrow$  (1) We can write  $A_j = B_j {}^tB_j$ , where  $B_j$  is a matrix of order  $k \times r_j$ . We can do that by taking a square root of  $A_j$  and then removing the columns of the null eigenvalues. Define

$$B = (B_1, \dots, B_m),$$

which is by hypothesis a square matrix, so that

$$B {}^tB = B_1 {}^tB_1 + \dots + B_m {}^tB_m = A_1 + \dots + A_m = I_k.$$

Therefore, the matrix  $B$  is orthogonal. Now, define a random vector  $\mathbf{Z}$  as

$$\mathbf{Z} = {}^tB\mathbf{Y}.$$

This means that  $\mathbf{Z} \sim \mathcal{N}_k({}^tB\boldsymbol{\mu}, I_k)$ . But then we have

$$Q_j = {}^t\mathbf{Y}A_j\mathbf{Y} = {}^t\mathbf{Y}B_j {}^tB_j\mathbf{Y} = {}^t\mathbf{Z}_j\mathbf{Z}_j,$$

where  $\mathbf{Z} = ({}^tZ_1, \dots, {}^tZ_m)$ . The  $\mathbf{Z}_j$ 's are independent and each  $\mathbf{Z}_j$  is distributed as  $N_{r_j} \sim ({}^tB_j\boldsymbol{\mu}, I_{r_j})$ . This implies that

$$Q_j = {}^t\mathbf{Z}_j\mathbf{Z}_j \sim \chi_{r_j}^2({}^t\boldsymbol{\mu}A_j\boldsymbol{\mu})$$

and  $Q_j$  independent of  $Q_i$ , for  $i \neq j$ . □

**Corollary 11.4.** If  $\mathbf{Y} \sim N_k(\boldsymbol{\mu}, I_k)$  and  $A$  is a symmetric idempotent matrix of order  $k$ , then

$${}^t\mathbf{Y}\mathbf{A}\mathbf{Y} \sim \chi_r^2({}^t\boldsymbol{\mu}\mathbf{A}\boldsymbol{\mu}),$$

where  $r = \text{rk } A = \text{tr } A$ .

*Proof.* The proof is an application of Theorem 11.3 with  $m = 2$  and

$$A_1 = A, \quad A_2 = I_k - A.$$

□

We introduced the decomposition

$${}^t\mathbf{Y}I_n\mathbf{Y} = {}^t\mathbf{Y}P\mathbf{Y} + {}^t\mathbf{Y}(I_n - P)\mathbf{Y}.$$

Assuming the second order hypotheses, we only found the expectation of the above three terms. Now, since we have distributional assumptions on  $\boldsymbol{\varepsilon}$ , we can determine the distributions of  ${}^t\mathbf{Y}\mathbf{Y}$ ,  ${}^t\mathbf{Y}P\mathbf{Y}$  and  ${}^t\mathbf{Y}(I_n - P)\mathbf{Y}$ . These distributions follow from Theorem 11.3

- the total residual

$${}^t\mathbf{Y}\mathbf{Y} \sim \sigma^2 \chi_n^2(\delta),$$

where  $\delta = \frac{{}^t\boldsymbol{\mu}\boldsymbol{\mu}}{\sigma^2}$ ;

- the regression residual

$${}^t\mathbf{Y}P\mathbf{Y} \sim \sigma^2 \chi_p^2(\delta)$$

- the residual

$${}^t\mathbf{Y}(I_n - P)\mathbf{Y} \sim \sigma^2 \chi_{n-p}^2(0),$$

since  ${}^t\boldsymbol{\mu}(I_n - P)\boldsymbol{\mu} = \mathbf{0}$ .

Finally, we have that  ${}^t\mathbf{Y}P\mathbf{Y}$  is independent of  ${}^t\mathbf{Y}(I_n - P)\mathbf{Y}$ .

*Remark 11.5.* From the perspective of point estimation, the use of the Maximum Likelihood approach does not make any difference, therefore we have the same estimator. The major improvement in introducing distributional hypotheses on  $\boldsymbol{\varepsilon}$  is related to hypothesis testing on  $\boldsymbol{\beta}$ .

The Least Squares approach does not provide any insight on testing hypotheses on the parameter  $\boldsymbol{\beta}$  and/or  $\sigma^2$ .

*Example 11.6.* Consider the problem of testing if the parameter  $\boldsymbol{\beta}$  is equal or different from zero:

$$H_0 : \boldsymbol{\beta} = \mathbf{0}, \quad H_1 : \boldsymbol{\beta} \neq \mathbf{0}.$$

To do that, we can apply the Likelihood Ratio test. Note that the parameter  $\sigma^2$  is assumed to be unknown.

- Under  $H_0$ , the maximum of the Log-Likelihood function is achieved at the point

$$\boldsymbol{\vartheta} = \begin{pmatrix} 0 \\ \hat{\sigma}_0^2 \end{pmatrix} \quad \hat{\sigma}_0^2 = \frac{\|\mathbf{Y}\|^2}{n},$$

as we have seen before (notice that  $\|\mathbf{y}\|^2 = \|\boldsymbol{\varepsilon}\|^2$  when  $\boldsymbol{\beta} = \mathbf{0}$ ).

128 11 Week XI

- The unconstrained maximum value of the Log-Likelihood function is achieved at the point

$$\boldsymbol{\vartheta} = \begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\sigma}^2 \end{pmatrix},$$

with  $\hat{\boldsymbol{\beta}}$  obtained by the Least Squares approach.

The Log-Likelihood ratio criterion leads to the following test statistic

$$\lambda(\mathbf{y}) = \frac{(\hat{\sigma}_0^2)^{-n/2} \exp\left\{-\frac{1}{2}\hat{\sigma}_0^2 \|\mathbf{y}\|^2\right\}}{(\hat{\sigma}^2)^{-n/2} \exp\left\{-\frac{1}{2}\hat{\sigma}^2 \|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2\right\}} \propto \left(\frac{\|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2}{\|\mathbf{y}\|^2}\right)^{n/2}.$$

The Likelihood ratio  $\lambda$  can be replaced by the monotone function

$$\lambda^*(\mathbf{y}) = \lambda^{-2/n}(\mathbf{y}) - 1 = \frac{\|\hat{\boldsymbol{\mu}}\|^2}{\|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2} = \frac{\frac{\|\mathbf{PY}\|^2}{p}}{\frac{\|(I_n - P)\mathbf{Y}\|^2}{n-p}}.$$

According to the Likelihood ratio test, we reject the null hypothesis if  $\lambda$  is sufficiently small or, in other terms, if  $\lambda^*$  is sufficiently large. In order to do that, we have to determine the distribution of  $\lambda^*$ :

$$\lambda^*(\mathbf{y}) = \frac{\|\hat{\boldsymbol{\mu}}\|^2}{\|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2} = \frac{{}^t\mathbf{Y}\mathbf{P}\mathbf{Y}}{{}^t\mathbf{Y}(I_n - P)\mathbf{Y}}.$$

The terms on the right-hand side have a Chi-squared distribution up to a multiplicative constant  $\sigma^2$ . Moreover, the terms  $\|\hat{\boldsymbol{\mu}}\|^2$  and  $\|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2$  are independent. This means that

$$F = \lambda^* \left( \frac{n-p}{p} \right)$$

is a statistic distributed as a Fisher distribution with  $(p, n-p)$  degrees of freedom and non-centrality parameter

$$\delta = \frac{{}^t\boldsymbol{\beta} {}^t\mathbf{X}\mathbf{X}\boldsymbol{\beta}}{\sigma^2}.$$

Note that the non-centrality parameter vanishes under  $H_0$ .

## 12

### Week XII

#### 12.1 11/12/18

##### 12.1.1 Generalised Linear Models

What are the limitations of linear models? There are some kind of problems that cannot be reduced to the form required by the linear model. This might happen in different scenarios.

1. The relationship between the response and the explanatory variables is of the type

$$y = r(x_1, \dots, x_p) + \varepsilon,$$

where the function  $r(\cdot)$  is non-linear in the parameters. In this case, linearisation does not work.

2. In addition, even the mathematical form of  $r(\cdot)$  might not be known exactly. We often know enough about the nature of the problem. The most common example is related to the range of  $y$ . For instance, suppose that  $y$  is the maximum fraction of skin which is affected by a certain disease and the explanatory variables are some factor for these particular diseases. In this case,  $y \in [0, 1]$  but a fixed regression hyperplane cannot satisfy this constraint.
3. The variance of the error term, and then the response variable, should be constant while in practical situations the observed data does not satisfy this requirement.
4. In linear models, we assumed that the response variable is distributed according to a Gaussian distribution. There are distributions in which this is not the case: *binary response, counting response, positive response*.

We consider a new class of models to deal with the above situations. This class is called **Generalised Linear Models**. They rely on the exponential family.

For this purpose, we introduce a specific notation for the exponential family, different from the one we used before.

**Definition 12.1.** A random variable  $Y$  is in the **exponential family**, meaning

$$Y \sim EF \left( b(\vartheta), \frac{\psi}{w} \right),$$

if the probability density function (or the probability mass function) has the following form:

$$f_Y(y) = \exp \left\{ \frac{w}{\psi} [y\vartheta - b(\vartheta)] + c(y, \psi) \right\}, \quad (12.1)$$

where  $\vartheta$  and  $\psi$  are scalar parameters,  $w$  is a known constant, whilst  $b(\cdot)$  and  $c(\cdot, \cdot)$  are fixed functions that determine the specific parametric family of distributions.

For any given choice of the **dispersion parameter**  $\psi$ , (12.1) forms an exponential family with parameter  $\vartheta$ . However, (12.1) is not an exponential family when  $\vartheta$  and  $\psi$  vary simultaneously. The variability of the random variable  $Y$  distributed as (12.1) is not controlled by a single parameter, but by the ratio  $\psi/w$  instead. The parameter  $w$  is known as the **weight**.

**Proposition 12.2.** Given a random variable  $Y$  such that

$$Y \sim EF \left( b(\vartheta), \frac{\psi}{w} \right),$$

then

$$\mathbb{E}[Y] = b'(\vartheta) \quad \text{Var}(Y) = b''(\vartheta).$$

*Proof.* Take

$$\log(f_Y(y)) = \frac{w}{\psi} [y\vartheta - b(\vartheta)] + c(y, \psi)$$

and consider

$$\begin{aligned} \frac{\partial}{\partial \vartheta} \log(f_Y(y)) &= \frac{\partial}{\partial \vartheta} \left( \frac{(y\vartheta - b(\vartheta))w}{\psi} + c(y, \psi) \right) \\ &= \frac{y - b'(\vartheta)}{\psi} w. \end{aligned}$$

Now, by Bartlett's Identities, we have

$$\begin{aligned} 0 &= \mathbb{E} \left[ \frac{\partial}{\partial \vartheta} \log f_Y(y) \right] \\ &= \mathbb{E} \left[ \frac{y - b'(\vartheta)}{\psi} w \right] \\ &= \frac{\mathbb{E}[Y] - b'(\vartheta)}{\psi} w, \end{aligned}$$

and

$$\begin{aligned} 0 &= \text{Var} \left( \frac{\partial}{\partial \vartheta} \left( \frac{y\vartheta - b(\vartheta)}{\psi} w + c(y, \psi) \right) \right) + \mathbb{E} \left[ \frac{\partial^2}{\partial \vartheta^2} \left( \frac{y\vartheta - b(\vartheta)}{\psi} w + c(y, \psi) \right) \right] \\ &= \frac{w^2}{\psi^2} \text{Var}(y - b'(\vartheta)) + \mathbb{E} \left[ -\frac{b''(\vartheta)}{\psi} w \right]. \end{aligned}$$

Therefore,

$$\mathbb{E}[Y] = b'(\vartheta) \quad \text{and} \quad \text{Var}(Y) = b''(\vartheta) \frac{\psi}{w}.$$

□

Usually, for simplicity's sake, we write

$$\mu := b'(\vartheta) \quad \text{and} \quad V(\mu) := b''(\vartheta),$$

noting that  $V$  is also known as the **variance function**.

*Example 12.3.* Consider  $Y \sim \mathcal{P}(\mu)$ . We can rewrite as usual, for  $y = 0, 1, \dots$ ,

$$\begin{aligned} f_Y(y) &= \frac{e^{-\mu} \mu^y}{y!} \\ &= \exp \{y \log \mu - \mu - \log y!\} \\ &= \exp \{y\vartheta - e^\vartheta - \log y!\} \end{aligned}$$

where  $\vartheta = \log \mu$  and, consequently,  $\mu = e^\vartheta$ . Therefore, by 12.1, we have

$$\begin{aligned} b(\vartheta) &= b'(\vartheta) = b''(\vartheta) = e^\vartheta, \\ \psi &= w = 1, \\ c(\psi, y) &= \log y!, \end{aligned}$$

with variance function equal to  $V(\mu) = \mu$ . In practice, here  $\psi$  is not present.

*Example 12.4.* Consider  $Y \sim \mathcal{N}(\mu, \sigma^2)$ , where

$$\begin{aligned} f_Y(y) &= \exp \left\{ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y - \mu)^2 \right\} \\ &= \exp \left\{ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y^2 - 2y\mu + \mu^2) - \frac{1}{2} \right\} \\ &= \exp \left\{ \frac{1}{\sigma^2} \left( y\mu - \frac{\mu^2}{2} \right) - \frac{1}{2} \left[ \frac{y^2}{\sigma^2} + \log(2\pi\sigma^2) \right] \right\}. \end{aligned}$$

The density function is inside the exponential family, from where we conclude that

$$\begin{aligned} \vartheta &= \mu, \\ b(\vartheta) &= \frac{\mu^2}{2}, \quad b'(\vartheta) = \mu, \quad b''(\vartheta) = 1, \\ \psi &= \sigma^2, \\ w &= 1, \\ \frac{\psi}{w} &= \sigma^2. \end{aligned}$$

*Remark 12.5.* The class of Generalised Linear Models has the main advantage of providing (by means of the exponential family) a unified treatment for a set of relevant models. Consider the linear model. For each observation  $Y_i$ , we defined a linear predictor

$$\eta_i = \mathbf{x}_i \boldsymbol{\beta},$$

where  $\mathbf{x}_i$  is the row of order  $i$  of the matrix  $X$ , for  $i = 1, \dots, n$ . Suppose that each observed  $y_i$  is sampled from  $Y_i \sim \mathcal{N}(\mu_i, \sigma^2)$ , where the relation between the mean value  $\mu_i$  and the linear predictor  $\eta_i$  is the identity. In other terms, we assumed that

$$\eta_i = \mu_i, \quad \eta_i = \mathbf{x}_i^t \boldsymbol{\beta}.$$

The Generalised Linear Models are obtained by expanding the above formulation of the linear model in two directions:

1. the distribution of  $Y_i$  is not restricted to being Gaussian, but it can be any distribution of the form  $EF(b(\vartheta_i), \psi/w_i)$ , such that  $b'(\vartheta_i) = \mu_i$  is provided;
2. other forms of relationship between the linear predictor  $\eta_i$  and the mean value  $\mu_i$  are possible. We consider

$$g(\mu_i) = \eta_i,$$

where  $g(\cdot)$  is a differentiable monotonic function called **link function**.

In a schematic way, we can specify a Generalised Linear Model as follows:

1.  $Y_i \sim EF(b(\vartheta_i), \psi/w_i)$  known as **error structure**;
2.  $g(\mu_i) = \eta_i$  known as **link function**;
3.  $\eta_i = \mathbf{x}_i^t \boldsymbol{\beta}$  known as **linear predictor**.

**Definition 12.6.** A statistical model is a Generalised Linear Model when it satisfies the following requirements:

1. the observations  $y_1, \dots, y_n$  are realisations of independent random variables  $Y_1, \dots, Y_n$ ;
2. each  $Y_i \sim EF(b(\vartheta_i), \psi/w_i)$  with  $\mathbb{E}[Y_i] = \mu_i = b'(\vartheta_i)$ , for  $i = 1, \dots, n$ ;
3. there exist a function  $g(\cdot)$  such that

$$g(\mu_i) = \mathbf{x}_i^t \boldsymbol{\beta},$$

where  $\mathbf{x}_i$  is a vector of constants and  $\boldsymbol{\beta}$  is a vector of parameters;

4. the function  $b(\cdot)$  and the parameter  $\psi$  are common to all the random variables  $Y_i$ , and the parameter  $w_i$  can change;
5. the functions  $b(\cdot)$ ,  $c(\cdot, \cdot)$  and  $g(\cdot)$  are known, as also the weight  $w_i$ .

Note that, when the link function  $g(\cdot)$  is the identity function and the error structure is Gaussian, then we go back to the linear model. The terminology “error structure” is preserved because of this connection with the linear model. However, in the generalised linear model, there is not an explicit error term  $\varepsilon$ .

In fact, in the Gaussian case, we can equivalently write

$$Y_i \sim \mathcal{N}(\mu, \sigma^2)$$

or

$$Y_i = \mu_i + \varepsilon_i,$$

where  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ . In this case, there is a complete separation between the systematic component  $\mu_i$  (which depends on the explanatory variable  $\mathbf{x}_i$  and  $\boldsymbol{\beta}$ ) and the purely erratic term  $\varepsilon_i$  (which does not). In the case of generalised linear models, this precise separation of the response variable in two components is no longer possible.

*Example 12.7 (Poisson Regression).* Consider  $Y_1, \dots, Y_n$  independent Poisson random variables, with mean values  $\mu_1, \dots, \mu_n$ , respectively. Assume in addition that

$$\log \mu_i = \mathbf{x}_i^t \boldsymbol{\beta} = \eta_i,$$

where the meaning of  $\mathbf{x}_i$  and  $\beta$  is the usual. In this particular situation,  $Y_i$  cannot be decomposed as

$$Y_i = \mu_i + \varepsilon_i,$$

where  $\varepsilon_i$  represents a purely erratic component, because in order to satisfy an identity as the above, the random variable  $\varepsilon_i$  should have a distribution depending on  $\mu_i$ , and this is not appropriate.

*Example 12.8 (Logistic Regression).* An important case of Generalised Linear Models is  $Y_i \sim \text{Ber}(\mu_i)$ . These are distributed as

$$\begin{aligned} f_{Y_i}(y_i) &= \exp \left\{ \log \left( \mu_i^{y_i} (1 - \mu_i)^{1-y_i} \right) \right\} \\ &= \exp \left\{ y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i) \right\} \\ &= \exp \left\{ y_i \log \left( \frac{\mu_i}{1 - \mu_i} \right) - \log(1 - \mu_i) \right\} \\ &= \exp \left\{ y_i \vartheta_i - \log \left( 1 + \frac{e^{\vartheta_i}}{1 - \mu_i} \right) \right\} \\ &= \exp \left\{ y_i \vartheta_i - \log(1 + e^{\vartheta_i}) \right\}, \end{aligned}$$

where

$$\vartheta_i = \log \left( \frac{\mu_i}{1 - \mu_i} \right) = \text{LOGIT}(\mu_i) \implies \mu_i = \frac{e^{\vartheta_i}}{1 + e^{\vartheta_i}}$$

with the assumption that  $\vartheta_i$  is a linear function of an explanatory variable. We can assume that  $\vartheta_i$  is a linear function of an explanatory variable, therefore

$$\begin{aligned} b(\vartheta_i) &= \log(1 + e^{\vartheta_i}) \\ b'(\vartheta_i) &= \frac{e^{\vartheta_i}}{1 + e^{\vartheta_i}} = \mu_i \\ b''(\vartheta_i) &= \frac{e^{\vartheta_i}}{(1 + e^{\vartheta_i})^2} = \mu(1 - \mu) = V(\mu_i) \end{aligned}$$

This is the classical form of the logistic regression model.

**12.2 12/12/18**

We can extend the basic logistic regression in two directions. The first extension consists in allowing several explanatory variables, forming the vector  $\mathbf{x}_i$ , and assuming

$$\vartheta_i = \mathbf{x}_i \boldsymbol{\beta}.$$

The second extension consists in the possibility of replicating the experiment for any given combination of experimental factors.

If we denote by  $m_i$  the number of replicates at a given combination  $\mathbf{x}_i$  of the explanatory variables, then the corresponding number of successes is

$$\bar{Y}_i \sim \mathcal{B}(m_i, \mu_i).$$

Since our goal is to investigate the relation between the explanatory variables and the probability of success, we do not use  $\bar{Y}_i$  as a response variable, but rather the proportion of successes

$$Y_i = \frac{\bar{Y}_i}{m_i}.$$

In this way,  $\mathbb{E}[Y_i] = \mu_i$ , which is the probability under consideration. In this case, the probability function is

$$\begin{aligned} f_Y(y_i) &= \binom{m_i}{m_i y_i} \mu_i^{m_i y_i} (1 - \mu_i)^{m_i(1-y_i)} \\ &= \exp \left\{ \log \binom{m_i}{m_i y_i} + \log \mu_i^{m_i y_i} + \log (1 - \mu_i)^{m_i(1-y_i)} \right\} \\ &= \exp \left\{ m_i y_i \log \mu_i + m_i(1 - y_i) \log(1 - \mu_i) + \log \binom{m_i}{m_i y_i} \right\} \\ &= \exp \left\{ m_i y_i [\log \mu_i - \log(1 - \mu_i)] + m_i \log(1 - \mu_i) + \log \binom{m_i}{m_i y_i} \right\} \\ &= \exp \left\{ m_i y_i \log \left( \frac{\mu_i}{1 - \mu_i} \right) - m_i \log \left( \frac{1}{1 - \mu_i} \right) + \log \binom{m_i}{m_i y_i} \right\} \\ &= \exp \left\{ m_i \left[ y_i \log \left( \frac{\mu_i}{1 - \mu_i} \right) - \log \left( \frac{1}{1 - \mu_i} \right) \right] + \log \binom{m_i}{m_i y_i} \right\} \\ &= \exp \left\{ m_i [y_i \vartheta_i - \log(1 + e^{\vartheta_i})] + \log \binom{m_i}{m_i y_i} \right\}, \end{aligned}$$

where

$$y_i \in \left\{ 0, \frac{1}{m_i}, \frac{2}{m_i}, \dots, 1 \right\}.$$

**How to make inference on the parameters  $\boldsymbol{\beta}$  and  $\psi$** 

Denote by  $p$  the size of  $\boldsymbol{\beta}$  and let  $X = (x_{ij})$  be the  $n \times p$  matrix. Let  $\mathbf{x}_i$  be the  $i$ -th row of  $X$ . Considering the observations  $Y_1, \dots, Y_n$  with mean values  $\mu_1, \dots, \mu_n$ , respectively, such that

$$g(\mu_i) = \mathbf{x}_i \boldsymbol{\beta},$$

we wish to estimate the parameters  $\beta$  and  $\psi$ .

Let us first focus on the parameter  $\beta$ , considering  $\psi$  as a noise. Since the  $Y_i$ 's are independent and we have a precise distributional assumption on their nature, we can write the Log-Likelihood function (which will be a sum):

$$\begin{aligned}\log \mathcal{L}(\beta) &= \sum_{i=1}^n \left( \frac{w_i [y_i \vartheta_i - b(\vartheta_i)]}{\psi} + c_i(y_i, \psi) \right) \\ &= \sum_{i=1}^n \mathcal{L}_i(\beta).\end{aligned}$$

To obtain the Likelihood equations, we can write

$$\frac{d\mathcal{L}_i}{d\beta_j} = \frac{d\mathcal{L}_i}{d\vartheta_i} \frac{d\vartheta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{d\eta_i}{d\beta_j},$$

whose components are

$$\begin{aligned}\frac{d\mathcal{L}_i}{d\vartheta_i} &= \frac{w_i [y_i - b'(\vartheta_i)]}{\psi} = \frac{w_i (y_i - \mu_i)}{\psi} \\ \frac{d\mu_i}{d\vartheta_i} &= b''(\vartheta_i) = \frac{w_i \text{Var}(Y_i)}{\psi} \\ \frac{d\eta_i}{d\beta_j} &= x_{ij}.\end{aligned}$$

Now,

$$\frac{d\mathcal{L}_i}{d\beta_j} = \frac{w_i (y_i - \mu_i)}{\psi} \frac{\psi}{w_i \text{Var}(Y_i)} \frac{d\mu_i}{d\eta_i} x_{ij}$$

and we obtain the Likelihood equation

$$\sum_{i=1}^n \frac{d\mathcal{L}_i}{d\beta_j} = \sum_{i=1}^n \frac{(y_i - \mu_i) x_{ij}}{\text{Var}(Y_i)} \frac{d\mu_i}{d\eta_i} = 0, \quad j = 1, \dots, p.$$

We can also compute the Fisher information. Consider the second derivatives of  $\mathcal{L}_i(\beta)$  and write

$$\begin{aligned}I_n(\beta) &= -\mathbb{E} \left[ \frac{d^2 \mathcal{L}_i}{d\beta_j d\beta_k} \right] \\ &= -\mathbb{E} \left[ \frac{d\mathcal{L}_i}{d\beta_j} \frac{d\mathcal{L}_i}{d\beta_k} \right] \\ &= -\mathbb{E} \left[ \left( \frac{(Y_i - \mu_i) x_{ij}}{\text{Var}(Y_i)} \frac{d\mu_i}{d\eta_i} \right) \left( \frac{(Y_i - \mu_i) x_{ik}}{\text{Var}(Y_i)} \frac{d\mu_i}{d\eta_i} \right) \right] \\ &= -\frac{x_{ij} x_{ik}}{\text{Var}(Y_i)} \left( \frac{d\mu_i}{d\eta_i} \right)^2, \quad i = 1, \dots, n \quad j, k = 1, \dots, p.\end{aligned}$$

This provides the entry of order  $(j, k)$  of the Fisher information matrix. In general, we have

136 12 Week XII

$$-\sum_{i=1}^n \mathbb{E} \left[ \frac{d^2 \mathcal{L}_i}{d\beta_j d\beta_k} \right] = -\sum_{i=1}^n \frac{x_{ij} x_{ik}}{\text{Var}(Y_i)} \left( \frac{d\mu_i}{d\eta_i} \right)^2,$$

which can be re-written as  $I(\beta) = {}^t X \tilde{W} X$ , where

$$\tilde{W} = \begin{pmatrix} \tilde{w}_1 & 0 & \dots & 0 \\ 0 & \tilde{w}_2 & \dots & 0 \\ \vdots & & \ddots & 0 \\ 0 & \dots & \dots & \tilde{w}_p \end{pmatrix}$$

and

$$\tilde{w}_i = \frac{1}{\text{Var}(Y_i)} \left( \frac{d\mu_i}{d\eta_i} \right)^2.$$

This is all we need to make inference.

*Example 12.9.* Consider a random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)$  from  $Y \sim f_Y(y)$ , where

$$f_Y(y) = \exp\{-\rho y + \log \rho\} = \exp\{\vartheta y + \log(-\vartheta)\},$$

where

$$\vartheta = -\rho = e^{-\eta}.$$

Here

$$\begin{aligned} b(\vartheta) &= -\log(-\vartheta), \\ b'(\vartheta) &= -\frac{1}{\vartheta} = \mu, \\ b''(\vartheta) &= \frac{1}{\vartheta^2} = \mu^2 = V(\mu), \\ \frac{d\mu}{d\eta} &= e^\eta = -\frac{1}{\vartheta}. \end{aligned}$$

We can apply the Likelihood equations. We have

$$\sum_{i=1}^n \frac{(y_i + \vartheta_i^{-1})}{\vartheta_i^{-2}} x_{ij} \frac{1}{\vartheta_i}.$$

Suppose for instance that  $p = 2$ . Set  $x_{i1} = 1$ ,  $x_{i2} = x_i$ . The Likelihood equations are

$$\begin{cases} \sum_{i=1}^n y_i \rho_i = 0 \\ \sum_{i=1}^n x_i y_i \rho_i - \sum_{i=1}^n x_i = 0 \end{cases}.$$

Try to compute the Fisher information

$$I(\beta) = {}^t X \hat{W} X.$$

### Canonical Link and Sufficient Statistics

The link function  $g(\mu)$  does not need to satisfy strong conditions, and it can be specified in different ways. There exist a specific choice of the limit function, also known as **canonical link**, which enjoys special properties:  $g(\mu_i) = \vartheta_i$ , such that  $\eta_i = \vartheta_i$ . The reason behind the name choice is evident, since the linear predictor  $\eta_i = \mathbf{x}_i \boldsymbol{\beta}$  coincides with the parameter  $\vartheta_i$  of the exponential family. In the case of the canonical link, we have that

$$\begin{aligned} \log \mathcal{L}(\boldsymbol{\beta}) &= \sum_{i=1}^n \frac{w_i}{\psi} [y_i \vartheta_i - b(\vartheta_i)] + \sum_{i=1}^n c_i(y_i, \eta_i) \\ &= \sum_{i=1}^n \frac{1}{\psi} [w_i y_i \mathbf{x}_i \boldsymbol{\beta} - w_i b(\mathbf{x}_i \boldsymbol{\beta})] + \sum_{i=1}^n c_i(y_i, \psi) \\ &= \frac{1}{\psi} \left[ \left( \sum_{i=1}^n w_i y_i \mathbf{x}_i \right) \boldsymbol{\beta} - \sum_{i=1}^n w_i b(\mathbf{x}_i \boldsymbol{\beta}) \right] + \sum_{i=1}^n c_i(y_i, \psi). \end{aligned}$$

This implies that, by the factorisation result,

$$\sum_{i=1}^n w_i y_i \mathbf{x}_i$$

is a sufficient statistic for  $\boldsymbol{\beta}$ , when  $\psi$  is a fixed quantity.

*Example 12.10.* Suppose that  $Y_i \sim \Gamma(w_i, w_i/\mu_i)$ , where the index  $w$  is constant on  $i$ :

$$\begin{aligned} f_{Y_i}(y_i) &= \frac{\left(\frac{w}{\mu_i}\right)^w}{\Gamma(w)} y_i^{w-1} \exp\left\{-\frac{w y_i}{\mu_i}\right\} \\ &= \exp\left\{w \left(-\frac{y_i}{\mu_i} - \log \mu_i\right) + (w-1) \log y_i - \log \Gamma(w) + w \log w_i\right\} \\ &= \exp\{w(\vartheta y_i + \log(-\vartheta)) + c_i(y_i, w)\}, \end{aligned}$$

where we set  $\vartheta = -1/\mu_i$  and

$$c_i(y_i, m_i) = (w-1) \log y_i - \log \Gamma(w) + w \log w.$$

The canonical link can be obtained when

$$-\mu_i^{-1} = \vartheta_i = \mathbf{x}_i \boldsymbol{\beta}.$$

Therefore,

$$\mathcal{L}(\boldsymbol{\beta}) = \exp\left\{w \left[ \left( \sum_{i=1}^n \mathbf{x}_i y_i \right) \boldsymbol{\beta} + \sum_{i=1}^n \log(-\mathbf{x}_i \boldsymbol{\beta}) \right] + \sum_{i=1}^n c(y_i, w) \right\},$$

confirming that

$$\sum_{i=1}^n \mathbf{x}_i y_i$$

138 12 Week XII

is the sufficient statistic for  $\beta$  if  $w$  is a fixed quantity.

The main advantage of using the canonical link is in terms of the Likelihood equations: indeed, they simplify quite a bit, since

$$\frac{d\mu_i}{d\eta_i} = \frac{d\mu_i}{d\vartheta_i} = \frac{d}{d\vartheta_i} b'(\vartheta_i) = b''(\vartheta_i).$$

This allows to write, using the simplification, more Likelihood equations:

$$\frac{d\mathcal{L}_i}{d\beta_j} = \left( \frac{y_i - \mu_i}{\text{Var}(Y_i)} \right) x_{ij} b''(\vartheta_i) = \frac{w_i (y_i - \mu_i) x_{ij}}{\psi} = 0, \quad j = 1, \dots, p,$$

which become independent from the specification of a particular distribution of  $Y_i$ . Moreover, with the canonical link, the Fisher information can also be simplified. Indeed, we have

$$\frac{d^2 \mathcal{L}_i}{d\beta_j d\beta_k} = - \frac{w_i x_{ij}}{\psi} \frac{d\mu_i}{d\beta_k},$$

which does not depend on the observations. Therefore,

$$\frac{d^2 \mathcal{L}}{d\beta_j d\beta_k} = -\mathbb{E} \left[ \frac{d^2 \mathcal{L}_i}{d\beta_j d\beta_k} \right],$$

so that the expected Fisher information and the observed information coincide.

## 12.3 13/12/18

### 12.3.1 Solving the Likelihood Equations

In most cases, the Likelihood equations

$$\sum_{i=1}^n \frac{(y_i - \mu_i)x_{ij}}{\text{Var}(Y_i)} \frac{d\mu_i}{d\eta_i} = 0$$

cannot be solved explicitly. Therefore, we resort to numerical methods. We have all the ingredients to apply Newton-Rhapsody method.

One of the reasons of the success of Generalised Linear Models, is indeed the possibility of using a single algorithm with a small adjustment to be made for the choice of the link function and the distribution of  $Y$ . The algorithm operates by solving a sequence of Least Squares problems.

The method works by adjusting the  $k$ -th approximate solution  $\beta^{(k)}$  to obtain the next approximation  $\beta^{(k+1)}$ . We will repeat this process until convergence.

Denote by  $\mathbf{u}^{(k)}$  the Fisher score computed at the solution  $\beta^{(k)}$ , and use the up-script notation for denoting iterations. In order to improve the approximate solution at the step  $k$ , we will use the Fisher scoring algorithm which can be expressed as follows:

$$\beta^{(k+1)} = \beta^{(k)} + \left[ I(\beta^{(k)}) \right]^{-1} \mathbf{u}^{(k)}, \quad (12.2)$$

which is similar to the Newton-Rapson method but with the negative Hessian matrix replaced by the expected Fisher information. Expression (12.2) can be re-written as follows

$$I(\beta^{(k)}) \beta^{(k+1)} = I(\beta^{(k)}) \beta^{(k)} + \mathbf{u}^{(k)}.$$

The right-hand side is a vector with the following  $h$ -th component:

$$\sum_{j=1}^n \left[ \sum_{i=1}^n \frac{x_{ih}x_{ij}}{\text{Var}(Y_i)} \left( \frac{d\mu_i}{d\eta_i} \right)^2 \right] \beta_j^{(k)} + \sum_{i=1}^n \frac{(y_i - \mu_i^{(k)})x_{ih}}{\text{Var}(Y_i)} \frac{d\mu_i}{d\eta_i},$$

therefore we can write

$$I(\beta^{(k)})\beta^{(k)} + \mathbf{u}^{(k)} = {}^t X \tilde{W}^{(k)} \mathbf{z}^{(k)},$$

where

$$\begin{aligned} z_i^{(k)} &= \sum_{j=1}^n x_{ij} \beta_j^{(k)} + (y_i - \mu_i^{(k)}) \left( \frac{\beta_i^{(k)}}{d\mu_i^{(k)}} \right) \\ &= \eta_i^{(k)} + (y_i - \mu_i^{(k)}) \left( \frac{d\eta_i^{(k)}}{d\mu_i^{(k)}} \right) \end{aligned}$$

and the entries of  $\tilde{W}$  have been defined yesterday. By combining

140 12 Week XII

$$I(\boldsymbol{\beta}^{(k)})\boldsymbol{\beta}^{(k+1)} = I(\boldsymbol{\beta}^{(k)})\boldsymbol{\beta}^{(k)} + \mathbf{u}^{(k)}$$

and

$$I(\boldsymbol{\beta}) = {}^tX\tilde{W}X$$

we get

$$\begin{aligned}({}^tX\tilde{W}X)\boldsymbol{\beta}^{(k+1)} &= ({}^tX\tilde{W}X)\boldsymbol{\beta}^{(k)} + \mathbf{u}^{(k)} \\ &= {}^tX\tilde{W}^{(k)}\mathbf{z}^{(k)}\end{aligned}$$

which is equivalent to

$$\boldsymbol{\beta}^{(k+1)} = ({}^tX\tilde{W}^{(k)}X)^{-1} {}^tX\tilde{W}^{(k)}\mathbf{z}^{(k)}.$$

This is the recursion to set the approximate solution at level  $(k+1)$  given the approximate solution at level  $k$ .

*Remark 12.11.* Using this approach, the dispersion parameter is not relevant for estimating  $\boldsymbol{\beta}$ .

### The Algorithm

The algorithm proceeds in two main steps.

1. When  $\boldsymbol{\beta}^{(k)}$  is given, the vector  $\mathbf{z}^{(k)}$  and the matrix  $\tilde{W}^{(k)}$  can be computed.
2. Using

$$\boldsymbol{\beta}^{(k+1)} = ({}^tX\tilde{W}^{(k)}X)^{-1} {}^tX\tilde{W}^{(k)}\mathbf{z}^{(k)}$$

we can compute  $\boldsymbol{\beta}^{(k+1)}$ .

These two steps are iterated until the sequence  $(\boldsymbol{\beta}^{(k)})_{k \geq 1}$  converges. This algorithm is called **Iteratively Reweighed Least Squares**.

*Remark 12.12.* If we consider the following Taylor expansion of  $g(y_i)$  around  $\mu_i$ ,

$$\begin{aligned}g(y_i) &\approx g(\mu_i) + (y_i - \mu_i)g'(\mu_i) \\ &= \mu_i(y_i - \mu_i)\frac{d\eta_i}{d\mu_i} \\ &= z_i,\end{aligned}$$

which means that the quantity  $z_i$  (component of the vector  $\mathbf{z}$ ) is a local approximation of  $g(y_i)$ .

This Remark also provides us with a hint for selecting the starting point of the algorithm. In particular, we can set

$$z_i^{(0)} = g(y_i), \quad \tilde{W}^{(0)} = I.$$

Given these two, we obtain  $\boldsymbol{\beta}^{(1)}$ , and therefore a starting point for the iteration.

One application of this algorithm is the Logistic Regression (repeated experiment):

$$\begin{aligned}\text{Var}(Y_i) &= \frac{1}{m_i} \mu_i(1 - \mu_i); \\ g(\mu_i) &= \log\left(\frac{\mu}{1 - \mu}\right) = \eta; \\ \mu &= \frac{e^\eta}{1 + e^\eta}; \\ \eta_i &= \mathbf{x}_i \boldsymbol{\beta} = \vartheta_i; \\ \frac{d\mu}{d\eta} &= \frac{e^\eta}{(1 + e^\eta)^2} = \mu(1 - \mu).\end{aligned}$$

By combining these, we have

$$\begin{aligned}z_i^{(k)} &= \eta_i^{(k)} + \frac{y_i - \mu_i^{(k)}}{\mu_i^{(k)}(1 - \mu_i^{(k)})} \\ \tilde{W}_i^{(k)} &= m_i \mu_i^{(k)}(1 - \mu_i^{(k)}).\end{aligned}$$

The equality  $\eta_i = \vartheta_i$  means that the LOGIT function is the canonical link for the logistic regression. To start, we can set

$$z_i^{(0)} = g(y_i) = \log\left(\frac{y_i}{1 - y_i}\right) = \log\left(\frac{\tilde{y}_i}{m_i - \tilde{y}_i}\right).$$

To avoid issues with  $\tilde{y}_i = 0$  and  $\tilde{y}_i = m_i$ , typically we start with

$$z_i^{(0)} = \log\left(\frac{\tilde{y}_i + \frac{1}{2}}{m_i - \tilde{y}_i + \frac{1}{2}}\right).$$

*Initial Settings:* we need a value for  $z_i$  we assign

$$z_i^{(0)} \leftarrow \log\left(\frac{\tilde{y}_i + 2^{-1}}{m_i - \tilde{y}_i + 2^{-1}}\right)$$

and the identity matrix to  $\tilde{W}$ .

*Repeat these steps until convergence:*

$$\begin{aligned}R &\leftarrow (\mathbf{X}^t \tilde{W} \mathbf{X})^{-1}, \\ \boldsymbol{\beta} &\leftarrow R \mathbf{X}^t \tilde{W} \mathbf{X}, \\ \boldsymbol{\eta} &\leftarrow \mathbf{X} \boldsymbol{\beta}, \\ \mu_i &\leftarrow \frac{e^{\eta_i}}{(1 + e^{\eta_i})}, \\ \Delta_i &\leftarrow \mu_i(1 - \mu_i), \\ z_i &\leftarrow \frac{\mu_i + (\tilde{y}_i - \mu_i)}{\Delta_i}, \quad i = 1, \dots, n \\ \tilde{W} &\leftarrow \text{diag}(\Delta_1 m_1, \dots, \Delta_n m_n)\end{aligned}$$

142 12 Week XII

*Final setting:*

$\hat{\beta} \leftarrow$  last value of  $\beta$ ,

$\text{Var}(\hat{\beta}) \leftarrow$  last value of R.