

# ulm university universität **UUI**

# **Mathematical Statistics**

Lecture Notes

Prof. Dr. Evgeny Spodarev

Ulm

2024

# Preface

The present lecture notes aim to give an introduction to different aspects of modern statistics. They are, in their current state, a result of holding lectures on statistics at Ulm University in the years 2010-2023 for students of mathematical bachelor's and master's programs.

The goal of the lectures is to provide an overview of typical problem settings and approaches to statistical inference. Additionally it aims to present a middle ground between practically orientated applied statistical monographs (which are usually mathematically sparse) and arid books on mathematical statistics. Whether I actually succeeded in finding said middle ground, shall be decided by the reader.

I would like to thank my colleagues at the Institute of Stochastics for their support and exhilarating discussions during the making of these notes. A special thanks goes to Linus Lach for the English translation of the German version and the creation of figures which accompany the text. I am also indebted to Tobias Brosch for the initial creation of the German LATEX—version and to Viet Hoang for the many corrections.

Ulm, July 11, 2025 Evgeny Spodarev

# Contents

Ta	Table of Contents								
1	Point Estimation								
	1.1	Paran	netric families of reference distributions	2					
		1.1.1	Gamma distribution	2					
		1.1.2	Student's t distribution	6					
		1.1.3	Fisher-Snedecor distribution (F distribution)	8					
	1.2	Metho	ods for obtaining point estimators	11					
		1.2.1	Plug-In estimator	12					
		1.2.2	Method of moments estimator	13					
		1.2.3	Maximum-likelihood estimator	15					
		1.2.4	Bayesian estimation	28					
		1.2.5	Resampling methods for obtaining point estimators	31					
	1.3	Furth	er quality properties of point estimators	36					
		1.3.1	Cramér-Rao inequality	36					
		1.3.2	Sufficiency	42					
		1.3.3	Completeness	47					
		1.3.4	Best unbiased estimator	49					
		1.3.5	$\delta ext{-Method}$	52					
2	Cor	Confidence Intervals 58							
	2.1	Introd	luction	58					
	2.2	One-s	ample problems	60					
		2.2.1	Normal distribution	60					
		2.2.2	Confidence intervals and stochastic inequalities	63					
		2.2.3	Asymptotic confidence intervals	64					
	2.3	Two-s	sample problems	68					
		2.3.1	Normally distributed samples	68					
		2.3.2	Poisson distributed random samples	70					
3	Testing Statistical Hypotheses 74								
	3.1	Gener	ral philosophy of testing	74					
	3.2	Non-r	andomized tests	84					

CONTENTS ii

		3.2.1	Parametric significance tests 84		
	3.3	3 Randomized test			
		3.3.1	Fundamentals		
		3.3.2	Neyman-Pearson test for simple hypotheses 90		
		3.3.3	One-sided Neyman-Pearson tests 96		
		3.3.4	Unbiased two-sided tests		
	3.4	Goodn	ess-of-fit tests		
		3.4.1	$\chi^2$ -goodness-of-fit test		
		3.4.2	$\chi^2$ -goodness-of-fit test of Pearson-Fisher 115		
		3.4.3	Shapiros goodness-of-fit test		
	3.5	More r	nonparametric tests		
		3.5.1	Binomial test		
		3.5.2	Randomness iteration tests		
4	Line	ear Reg	gression 128		
	4.1	Multiv	ariate normal distribution		
		4.1.1	Properties of the multivariate normal distribution 132		
		4.1.2	Linear and quadratic forms of normally distributed		
			random variables		
	4.2	Multiv	ariate linear regression models with full rank 141		
		4.2.1	Method of least squares		
		4.2.2	Estimator of the variance $\sigma^2$		
		4.2.3	Maximum likelihood estimator for $\beta$ and $\sigma^2$ 148		
		4.2.4	Tests for regression parameters		
		4.2.5	Confidence region		
	4.3	Multiv	ariate linear regression with $rank(X) < m \dots 158$		
		4.3.1	Generalized inverse		
		4.3.2	OLS estimator for $\beta$		
		4.3.3	Functions that can be estimated without bias 163		
		4.3.4	Normally distributed error terms 166		
		4.3.5	Hypothesis testing		
		4.3.6	Confidence regions		
		4.3.7	Introduction to variance analysis		
5	Ger	ieralize	ed linear models 177		
	5.1	Expon	ential family of distributions		
	5.2	Link fu	nctions		
	5.3	Maxim	num likelihood estimator for $\beta$		
	5.4	Asymptotic tests for $\beta$			
	5.5	Criteri	a for model selection or model adjustment 196		

CONTENTS	iii

6	Pri	ncipal Component Analysis	199
	6.1	Introduction	. 199
	6.2	PCA on model level	. 200
	6.3	PCA on data level	. 209
	6.4	Asymptotic distributions of principal components for normal	
		distributed random samples	. 212
	6.5	Outlier detection	. 214
	6.6	PCA and regression	. 217
	6.7	Numeric calculation of principal components	. 222
Literature			225
In	dex		228

# Chapter 1

# Point Estimation

Let  $(x_1, \ldots, x_n)$  be a given sample. Assume that it is a realization of a random sample  $(X_1, \ldots, X_n)$ , where  $X_1, \ldots, X_n$  are independent identically distributed (i.i.d.) random variables with unknown distribution F. Further assume that F is an element of a parametric family of distributions given by  $\{F_{\theta} : \theta \in \Theta\}$ . Here  $\theta = (\theta_1, \ldots, \theta_m) \in \Theta$  denotes the m-dimensional parameter vector of the distribution  $F_{\theta}$ , and  $\Theta \subset \mathbb{R}^m$  is the so called parameter space (a Borel subset of  $\mathbb{R}^m$ , which is composed of all valid parameter values). The parametrization  $\theta \mapsto F_{\theta}$  is set to be identifiable, under the assumption that  $F_{\theta_1} \neq F_{\theta_2}$  for  $\theta_1 \neq \theta_2$ .

An important task in statistics, discussed in this chapter, is the estimation of the parameter vector  $\theta$  (or a part of  $\theta$ ) on the basis of a given sample  $(x_1, \ldots, x_n)$ . In this context, the described procedure is called *point estimation* with respect to a *point estimator*  $\hat{\theta} : \mathbb{R}^n \to \mathbb{R}^m$ , which is a valid sample function. Usually one assumes that

$$P\left(\hat{\theta}(X_1,\ldots,X_n)\in\Theta\right)=1\,,$$

even though exceptions exist. The probability space  $(\Omega, \mathcal{F}, P)$  on which the random sample is defined has yet to be specified thoroughly. Here, the so called *canonical probability space* comes into play, which is defined by

$$\Omega = \mathbb{R}^{\infty}, \qquad \mathcal{F} = \mathcal{B}_{\mathbb{R}}^{\infty} = \mathcal{B}_{\mathbb{R}} \otimes \mathcal{B}_{\mathbb{R}} \otimes \dots$$

with probability measure P given by

$$P\left(\left\{\omega = (\omega_1, \dots, \omega_n, \dots\right) \in \mathbb{R}^\infty : \omega_{i_1} \leq x_{i_1}, \dots, \omega_{i_k} \leq x_{i_k}\right\}\right) = \prod_{j=1}^k F_\theta(x_{i_j})$$

for all  $k \in \mathbb{N}$  and  $1 \le i_1 < \cdots < i_k$ . In order to emphasize that P depends on  $\theta$ , the notation  $P_{\theta}$ ,  $E_{\theta}$  and  $Var_{\theta}$  is introduced for the measure P as well as the expectation E and variance Var with respect to P.

On the canonical probability space  $(\Omega, \mathcal{F}, P_{\theta})$  it holds that  $X_i(\omega) = \omega_i$  (projection on the *i*'th coordinate),  $i = 1, \ldots, n$ , with

$$P_{\theta}(X_i \leq x_i) = P_{\theta}(\{\omega \in \Omega : \omega_i \leq x_i\}) = F_{\theta}(x_i), \quad i = 1, \dots, n, \quad x_i \in \mathbb{R}.$$

# 1.1 Parametric families of reference distributions

In the lecture "Elementary Probability Theory" some parametric families have already been introduced. In this section, more parametric families of distributions that play a special role (e.g. as reference distributions in estimation theory, statistical tests and confidence intervals) will be presented.

## 1.1.1 Gamma distribution

First, consider the following special functions:

1. The Gamma function:

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

The following properties hold:

- $\Gamma(1) = 1$
- $\Gamma(1/2) = \sqrt{\pi}$ ,
- $\Gamma(p+1) = p\Gamma(p)$  for all p > 0,
- $\Gamma(n+1) = n!$  for all  $n \in \mathbb{N}$ .
- 2. The Beta function:

$$B(p,q) = \int_0^1 t^{p-1} (1-t)^{q-1} dt \,, \quad p,q > 0 \,.$$

The following properties hold:

- B(p,q) = B(q,p),
- $B(p,q) = \frac{\Gamma(p)\Gamma(q)}{\Gamma(p+q)}$  for all p,q > 0,

**Definition 1.1.1.** The Gamma distribution with parameters  $\lambda > 0$  and p > 0 is an absolutely continuous distribution with probability density function

$$f_X(x) = \begin{cases} \frac{\lambda^p x^{p-1}}{\Gamma(p)} e^{-\lambda x}, & x \ge 0, \\ 0, & x < 0. \end{cases}$$
 (1.1)

Denote by  $X \sim \Gamma(\lambda, p)$  a random variable X which is Gamma distributed with parameters  $\lambda$  and p. Obviously  $X \geq 0$  almost surely (a.s.).

**Exercise 1.1.2.** Show that (1.1) is indeed a probability density function.

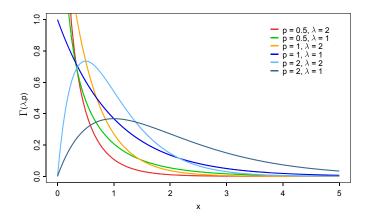


Figure 1.1: Probability density function of the Gamma distribution with various choices for the parameters  $\lambda > 0$  and p > 0.

## Example 1.1.3.

- 1. The Gamma distribution is often used for modeling small and medium sized insurance claims.
- 2. If p = 1, then  $\Gamma(\lambda, 1) = Exp(\lambda)$ , i.e. the Exponential distribution with parameter  $\lambda > 0$ .

# **Theorem 1.1.4.** Let $X \sim \Gamma(\lambda, p)$ .

1. The moment generating function  $\Psi_X(s)$  of X is given by

$$\Psi_X(s) = \mathbf{E}e^{sX} = \frac{1}{(1 - s/\lambda)^p}, \quad s < \lambda.$$

The characteristic function  $\varphi_X(s)$  of X is given by

$$\varphi_X(s) = \mathbb{E}e^{isX} = \frac{1}{(1 - is/\lambda)^p}, \quad s \in \mathbb{R}.$$

2. The k-th moments of X are given by

$$\mathrm{E}X^k = \frac{p(p+1)\cdot\dots\cdot(p+k-1)}{\lambda^k}\,,\quad k\in\mathbb{N}\,.$$

## Proof

#### 1. Consider

$$\Psi_X(s) = \int_0^\infty e^{sx} f_X(x) \, dx = \frac{\lambda^p}{\Gamma(p)} \int_0^\infty x^{p-1} e^{\frac{s^2}{(s-\lambda)}x} \, dx$$

$$= \frac{\lambda^p}{\Gamma(s-\lambda)x = y} \int_0^\infty \frac{y^{p-1}}{(-(s-\lambda))^p} e^{-y} \, dy = \frac{\lambda^p \Gamma(p)}{\Gamma(p)(\lambda - s)^p}$$

$$= \left(\frac{\lambda}{\lambda - s}\right)^p = \frac{1}{(1 - s/\lambda)^p}, \quad \lambda > s.$$

If  $s \in \mathbb{C}$  and  $\operatorname{Re}(s) < \lambda$ , then  $\Psi_X(s)$  is holomorphic on D, where  $D := \{s = x + iy \in \mathbb{C} : x < \lambda\}$ . It holds that

$$\Psi_X(s) = \varphi_X(-is) \,,$$

for s = it,  $0 < \lambda$ , which implies that

$$\varphi_X(t) = \Psi_X(it), \qquad t \in \mathbb{R}$$

Ultimately, this yields

$$\varphi_X(t) = \frac{1}{(1 - it/\lambda)^p}, \quad t \in \mathbb{R}.$$

2.

$$EX^k = \Psi^{(k)}(0) \Longrightarrow EX^k = \frac{p \cdot (p+1) \cdot \dots \cdot (p+k-1)}{\lambda^k}, \quad k \in \mathbb{N}.$$

**Corollary 1.1.5** (Stability of the Gamma distribution). If  $X \sim \Gamma(\lambda, p_1)$ ,  $Y \sim \Gamma(\lambda, p_2)$  and X and Y are independent, then  $X + Y \sim \Gamma(\lambda, p_1 + p_2)$ .

**Proof** It holds that

$$\varphi_{X+Y}(s) = \varphi_X(s) \cdot \varphi_Y(s)$$

$$= \frac{1}{(1 - is/\lambda)^{p_1}} \cdot \frac{1}{(1 - is/\lambda)^{p_2}}$$

$$= \left(\frac{1}{1 - is/\lambda}\right)^{p_1 + p_2}$$

$$= \varphi_{\Gamma(\lambda, p_1 + p_2)}(s).$$

Since the characteristic function uniquely determines the distribution of a random variable,  $X + Y \sim \Gamma(\lambda, p_1 + p_2)$  holds.

**Example 1.1.6.** Let  $X_1, \ldots, X_n \sim Exp(\lambda)$  be independent. By Corollary 1.1.5 it holds that  $X = X_1 + \ldots + X_n \sim \Gamma(\lambda, \underbrace{1 + \ldots + 1}) = \Gamma(\lambda, n)$ , since

 $Exp(\lambda) = \Gamma(\lambda, 1)$ . This special case of the Gamma distribution is also called *Erlang distribution* with parameters  $\lambda > 0$  and  $n \in \mathbb{N}$ . Notation:  $X \sim Erl(\lambda, n)$ .

In summary: 
$$Erl(\lambda, n) = \Gamma(\lambda, n), \quad \lambda > 0, n \in \mathbb{N}.$$

Interpretation: In risk theory the random variables  $X_i$  represent interarrival times for the individual damages. Here  $X = \sum_{i=1}^{n} X_i$  represents the occurrence time of the *n*-th loss with  $X \sim Erl(\lambda, n)$ .

**Definition 1.1.7** ( $\chi^2$  distribution). X is a  $\chi^2$  distributed random variable with k degrees of freedom (Notation:  $X \sim \chi_k^2$ ), if  $X \stackrel{d}{=} X_1^2 + \ldots + X_k^2$ , where  $X_1, \ldots, X_k \sim N(0, 1)$  are i.i.d. random variables.

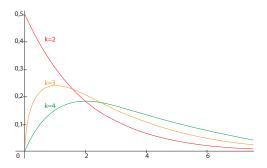


Figure 1.2: Probability density function of the  $\chi_k^2$  distribution with k=2,3,4 degrees of freedom.

**Theorem 1.1.8** ( $\chi^2$  distribution: Special case of the Gamma distribution with  $\lambda=1/2,\ p=k/2$ ). If  $X\sim\chi^2_k$ , then

1.  $X \sim \Gamma(1/2, k/2)$ , i.e.

$$f_X(x) = \begin{cases} \frac{x^{k/2 - 1} e^{-x/2}}{2^{k/2} \Gamma(k/2)}, & x \ge 0\\ 0, & x < 0 \end{cases}$$
 (1.2)

2. In particular EX = k, Var X = 2k.

#### Proof

1. Let  $X=X_1^2+\ldots+X_k^2$  with  $X_i\sim N(0,1)$  i.i.d. random variables. Calculating the distribution function of  $X_i^2$  by [33, Satz 3.6.4] and using

$$f_{X_1^2}(x) = \frac{1}{2\sqrt{x}} \left( f_{X_1}(\sqrt{x}) + f_{X_1}(-\sqrt{x}) \right)$$

yields

$$P(X_1^2 \le x) = \int_0^x \left( \frac{1}{\sqrt{2\pi}} e^{\frac{-t}{2}} \frac{1}{2\sqrt{t}} dt + \frac{1}{\sqrt{2\pi}} e^{-t/2} \frac{1}{2\sqrt{t}} \right) dt$$
$$= \int_0^x \frac{(1/2)^{-1/2} t^{1/2 - 1}}{\Gamma(1/2)} e^{-t/2} dt, \qquad x \ge 0.$$

Thus  $X_1^2 \sim \Gamma(1/2, 1/2) \Longrightarrow X \sim \Gamma(1/2, \underbrace{1/2 + \ldots + 1/2}_{k \text{ times}}) = \Gamma(1/2, k/2)$  and therefore (1.2) with respect to the density holds.

Because of the additivity of the expected value and the indepen

2. Because of the additivity of the expected value and the independence of the  $X_i$ , it holds that

 $EX = k \cdot EX_1^2$ ,  $Var X = k Var X_1^2$ ,  $E(X_1^2) = E(\Gamma(1/2, 1/2)) = 1$  by Theorem 1.1.4, 2. Indeed,

$$E(X_1^2) = \frac{1/2}{1/2} = 1$$
,  $E(X_1^4) = \frac{1/2(1/2+1)}{(1/2)^2} = \frac{3/4}{1/4} = 3$ ,  
 $Var X_1^2 = E(X_1^4) - (E(X_1^2))^2 = 3 - 1 = 2$ .

#### 1.1.2 Student's t distribution

**Definition 1.1.9.** Let X and Y be independent random variables, where  $X \sim N(0,1)$  and  $Y \sim \chi_r^2$ . The random variable

$$U \stackrel{d}{=} \frac{X}{\sqrt{Y/r}}$$

is called Student or t distributed with r degrees of freedom. Notation:  $U \sim t_r$ .

**Theorem 1.1.10** (Probability density function of the t distribution). If  $X \sim t_r$ , then

$$f_X(x) = \frac{1}{\sqrt{r}B\left(\frac{r}{2}, \frac{1}{2}\right)} \cdot \frac{1}{\left(1 + \frac{x^2}{r}\right)^{\frac{r+1}{2}}}, \qquad x \in \mathbb{R}.$$

<sup>&</sup>lt;sup>1</sup>Named after the mathematician William Sealy Gosset, who signed his work under the pseudonym "Student".

7

2. 
$$EX = 0$$
,  $Var X = \frac{r}{r-2}$ ,  $r \ge 3$ .

#### Remark 1.1.11.

1. Student's t distribution is symmetric. In particular,

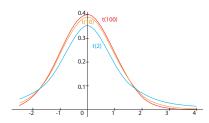


Figure 1.3: Probability density function f of the t distribution for r=2,10,100

$$t_{r,\alpha} = -t_{r,1-\alpha}, \quad \alpha \in (0,1),$$

where  $t_{r,\alpha}$  is the  $\alpha$  quantile of the Student's distribution with r degrees of freedom.

- 2. For  $r \to \infty$  it holds that  $f_r(x) \to \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ ,  $x \in \mathbb{R}$ . (Proof: Exercise)
- 3. For r=1 the t distribution coincides with the standard Cauchy distribution, i.e. it holds that  $t_1 = Cauchy(0,1)$  with probability density function  $f(x) = \frac{1}{\pi(1+x^2)}$ . The expected value of  $t_1$  doesn't exist.

# Proof of Theorem 1.1.10:

1. It holds that  $X:=\varphi(Y,Z)$ , where  $\varphi(x,y)=\frac{x}{\sqrt{y/r}}$  and V=(Y,Z) is a two dimensional random vector with  $Y\sim N(0,1)$  and  $Z\sim\chi^2_r$  independent of Y. The density transformation theorem [33, Theorem 3.6.6] states that

$$f_{\varphi(V)}(x) = f_V(\varphi^{-1}(x))|J|,$$

where  $|J| = |\det J|$ , where  $J = \left(\frac{\partial \varphi_i^{-1}(x)}{\partial x_j}\right)_{i,j=1}^n$  denotes the Jacobi matrix of the function  $\varphi = (\varphi_1, \dots, \varphi_n) : \mathbb{R}^n \to \mathbb{R}^n$ . Computing  $\varphi^{-1}$ , where  $\varphi : (x,y) \mapsto (v,w)$  as above, with  $v = \frac{x}{\sqrt{y/r}}$ , w = y yields

$$v = \frac{x}{\sqrt{\frac{y}{r}}} \Rightarrow x = v\sqrt{\frac{y}{r}} = v\sqrt{\frac{w}{r}}$$
.

Thus,

$$\varphi^{-1}:(v,w)\mapsto \left(v\sqrt{\frac{w}{r}},w\right),$$

and the Jacobi matrix is given by

$$J = \begin{pmatrix} \frac{\partial \varphi_1^{-1}}{\partial v} & \frac{\partial \varphi_1^{-1}}{\partial w} \\ \frac{\partial \varphi_2^{-1}}{\partial v} & \frac{\partial \varphi_2^{-1}}{\partial w} \end{pmatrix} = \begin{pmatrix} \sqrt{\frac{w}{r}} & \frac{v}{2\sqrt{wr}} \\ 0 & 1 \end{pmatrix}.$$

For V = (Y, Z) with Y and Z independent it follows that

$$f_V(x,y) = f_Y(x) \cdot f_Z(y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \cdot \frac{y^{r/2 - 1} e^{-y/2}}{\Gamma(r/2) 2^{r/2}} = \frac{y^{r/2 - 1} e^{-\frac{y + x^2}{2}}}{2^{\frac{r+1}{2}} \Gamma(1/2) \Gamma(r/2)}$$

for all  $x \in \mathbb{R}$  and y > 0. The density transformation theorem ultimately yields

$$\begin{split} f_X(v) &= \int_0^\infty f_{\varphi(V)}(u,w) dw = \int_0^\infty f_V(\varphi^{-1}(v,w)) |J| \, dw \\ &= \int_0^\infty \frac{e^{-(v^2 \frac{w}{r} + w)/2} w^{r/2 - 1}}{2^{\frac{r+1}{2}} \Gamma(1/2) \Gamma(r/2)} \sqrt{w/r} \, dw \\ &= \frac{1}{\sqrt{r} 2^{\frac{r+1}{2}} \Gamma(1/2) \Gamma(r/2)} \cdot \int_0^\infty w^{\frac{r-1}{2}} e^{-\frac{v^2}{r} + 1} \frac{1}{2} \cdot w \, dw \\ &= \frac{1}{w = \frac{2t}{v^2/r + 1}} \frac{1}{\sqrt{r} 2^{\frac{r+1}{2}} \Gamma(1/2) \Gamma(r/2)} \cdot \int_0^\infty \frac{2^{\frac{r-1}{2} + 1} t^{\frac{r-1}{2}}}{(v^2/r + 1)^{\frac{r-1}{2} + 1}} e^{-t} dt \\ &= \frac{2^{\frac{r+1}{2}} \Gamma(\frac{r+1}{2})}{(\frac{v^2}{r} + 1)^{\frac{r+1}{2}} \sqrt{r} 2^{\frac{r+1}{2}} \Gamma(1/2) \Gamma(r/2)} \\ &= \frac{1}{\sqrt{r} B(r/2, 1/2) (1 + v^2/r)^{\frac{r+1}{2}}} \end{split}$$

#### 2. Exercise.

# 1.1.3 Fisher-Snedecor distribution (F distribution)

**Definition 1.1.12.** Let  $X \stackrel{d}{=} \frac{U_r/r}{U_s/s}$ , where  $U_r \sim \chi_r^2$ ,  $U_s \sim \chi_s^2$ ,  $r,s \in \mathbb{N}$ ,  $U_r, U_s$  are independent. Then, X is Fisher or F distributed with  $r,s \in \mathbb{N}$  degrees of freedom. Notation:  $X \sim F_{r,s}$ .

**Lemma 1.1.13.** Let  $X \sim F_{r,s}$ ,  $r, s \in \mathbb{N}$ . Then, X is absolutely continuously distributed with probability density function

$$f_X(x) = \frac{x^{r/2-1}}{B(r/2, s/2)(r/s)^{-r/2}(1 + (r/s) \cdot x)^{\frac{r+s}{2}}} \cdot I(x > 0),$$

where  $I_B(x)$  denotes the indicator function of the set B.

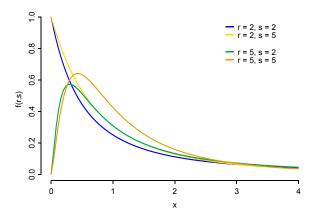


Figure 1.4: Probability density function of the  $F_{r,s}$  distribution for various choices of r and s.

**Proof** Note that for  $U_r \sim \chi_r^2$  the probability density function is given by

$$f_{U_r}(x) = \frac{x^{r/2-1}e^{-x/2}}{\Gamma(r/2)2^{r/2}}, \quad x > 0, \quad r \in \mathbb{N}.$$

Consequently,

$$P(U_r/r < x) = P(U_r < rx) = F_{U_r}(rx),$$

and therefore

$$f_{U_{r/r}}(x) = (F_{U_r}(rx))' = r \cdot f_{U_r}(rx) = \frac{r(rx)^{r/2 - 1} e^{\frac{-rx}{2}}}{\Gamma(r/2) 2^{r/2}} \cdot I(x > 0)$$
$$= \frac{r^{r/2} x^{r/2 - 1} e^{-r/2 \cdot x}}{\Gamma(r/2) 2^{r/2}} \cdot I(x > 0).$$

By the density transformation theorem for the ratio of two random variables [33, Theorem 3.6.9., 2] it holds that

$$f_{\frac{U_{r/r}}{U_{s/s}}}(x) = \int_0^\infty t f_{U_{r/r}}(xt) \cdot f_{U_{s/s}}(t) dt \cdot I(x>0).$$

Hence,

$$\begin{split} f_X(x) &= \int_0^\infty t \frac{r^{r/2}(tx)^{r/2-1}e^{-\frac{rtx}{2}}}{\Gamma(r/2)2^{r/2}} \cdot \frac{s^{s/2}t^{s/2-1}e^{-st/2}}{\Gamma(s/2)2^{s/2}} \, dt \\ &= \frac{r^{r/2}s^{s/2}x^{r/2-1}}{\Gamma(r/2)\Gamma(s/2)2^{\frac{r+s}{2}}} \cdot \int_0^\infty t^{r/2+s/2-1}e^{-\frac{st/2}{2}} \, dt \\ &= \frac{r^{r/2}s^{s/2}x^{r/2-1}}{\Gamma(r/2)\Gamma(s/2)} \cdot \int_0^\infty \frac{y^{\frac{r+s}{2}-1}}{(rx+s)^{\frac{r+s}{2}}} \cdot e^{-y} \, dy \\ &= \frac{r^{r/2}s^{s/2}x^{r/2-1}}{\Gamma(r/2)\Gamma(s/2)s^{\frac{r+s}{2}}(1+\frac{r}{s}\cdot x)^{\frac{r+s}{2}}} \cdot \Gamma\left(\frac{r+s}{2}\right) \\ &= \frac{(r/s)^{r/2}x^{r/2-1}}{B(r/2,s/2)(1+\frac{r}{s}x)^{\frac{r+s}{2}}} \cdot I(x>0) \, . \end{split}$$

**Remark 1.1.14.** Let  $X \sim F_{r,s}, r, s \in \mathbb{N}$ , with probability density function  $f_X$ .

- 1. Some graphs of the F distribution are shown in Figure 1.4.
- 2. Some properties of the F distribution:

**Lemma 1.1.15.** Let  $X \sim F_{r,s}, r, s \in \mathbb{N}$ . Then,

(a) 
$$EX = \frac{s}{s-2}, \qquad s \ge 3.$$

(b) 
$$\operatorname{Var} X = \frac{2s^2(r+s-2)}{r(s-4)(s-2)^2}, \qquad s \ge 5.$$

(c) Denote by  $F_{r,s,\alpha}$  the  $\alpha$ -quantile of the  $F_{r,s}$  distribution. Then,

$$F_{r,s,\alpha} = \frac{1}{F_{s,r,1-\alpha}}, \quad \alpha \in (0,1).$$

Exercise 1.1.16. Prove Lemma 1.1.15.

3. The following approximation formula holds for quantiles  $F_{r,s,\alpha}$  (cf.

Abramowitz, Stegun (1972)):  $F_{r,s,\alpha} \approx e^{\omega}$ , where

$$\omega = 2\left(\frac{\alpha(h+a)^{1/2}}{h} - \left(\frac{1}{r-1} - \frac{1}{s-1}\right) \cdot \left(a + \frac{5}{6} - \frac{2}{3h}\right)\right),$$

$$h = 2\left(\frac{1}{r-1} + \frac{1}{s-1}\right)^{-1},$$

$$a = \frac{z_{\alpha}^2 - 3}{6}$$

and  $z_{\alpha}$  is the  $\alpha$  quantile of the N(0,1) distribution.

# 1.2 Methods for obtaining point estimators

The following introductory examples were given in the lecture "Elementary Probability Theory and Statistics".

#### Definition 1.2.1.

- 1. The function  $\hat{F}_n(x) = \#\{x_i : x_i \leq x, i = 1, \dots, n\}/n$  for all  $x \in \mathbb{R}$  is called *empirical distribution function of a realized sample*  $(x_1, \dots, x_n)$ . Here  $\hat{F}_n : \mathbb{R}^{n+1} \to [0,1]$  holds, since  $\hat{F}_n(x) = \varphi(x_1, \dots, x_n, x)$ .
- 2. The random variable  $\hat{F}_n: \Omega \times \mathbb{R} \to [0,1]$  which is indexed by  $x \in \mathbb{R}$  is called *empirical distribution function of the random sample* given by  $(X_1, \ldots, X_n)$ , if

$$\hat{F}_n(x,\omega) = \hat{F}_n(x) = \frac{1}{n} \# \{X_i, i = 1, \dots, n : X_i(\omega) \le x\}, \quad x \in \mathbb{R}.$$

Equivalently to Definition 1.2.1 it can be shown that

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \le x), \quad x \in \mathbb{R},$$

where

$$I(x \in A) = \begin{cases} 1, & x \in A \\ 0, & \text{otherwise.} \end{cases}$$

It holds that

$$\hat{F}_n(x) = \begin{cases} 1, & x \ge x_{(n)}, \\ \frac{i}{n}, & x_{(i)} \le x < x_{(i+1)}, & i = 1, \dots, n-1, \\ 0, & x < x_{(1)}. \end{cases}$$

for  $x_{(1)} < x_{(2)} < \ldots < x_{(n)}$ .

The height of the jump at  $x_{(i)}$  is equal to the relative frequency  $f_i$  of  $x_{(i)}$ . If  $x_{(i)} = x_{(i+1)}$  for a  $i \in \{1, \ldots, n\}$ , the value i/n does not occur (cf. [33, Section 6.3.2]).

# 1.2.1 Plug-In estimator

Based on the empirical distribution function  $\hat{F}_n$ , the *Plug-in method* yields the class of *Plug-in estimators*. Let  $M := \{F : F \text{ is a distribution function}\}$ .

**Definition 1.2.2.** Let the parameter  $\theta$  of the distribution F be given as a functional  $T: M \to \mathbb{R}$  of F, i.e  $\theta = T(F)$ . Then,  $\hat{\theta} = T(\hat{F}_n)$  is called the *Pluq-in estimator* for  $\theta$ .

**Definition 1.2.3.** Let F be an arbitrary distribution function. The functional  $T: M \to \mathbb{R}$  is called *linear*, if

$$T(aF_1+bF_2) = aT(F_1)+bT(F_2)$$
 for all  $a, b \in \mathbb{R}_+, a+b=1, F_1, F_2 \in M$ .

Consider a special class of linear functionals given by

$$T(F) = \int_{\mathbb{R}} r(x) dF(x),$$

where r(x) is an arbitrary continuous function with  $\mathbb{E}(r(X)) < \infty$ . An example for such T is given by

$$EX^k = \int_{\mathbb{R}} x^k dF(x), \qquad k \in \mathbb{N}.$$

**Lemma 1.2.4.** The Plug-in estimator for  $\theta = \int_{\mathbb{R}} r(x) dF(x)$  is given by

$$\hat{\theta} = \int_{\mathbb{R}} r(x) d\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n r(x_i).$$

Exercise 1.2.5. Prove Lemma 1.2.4!

Example 1.2.6 (Plug-in estimator).

- 1.  $\bar{X}_n$  is a Plug-in estimator for the expected value  $\mu$ .
- 2. Plug-in estimator for  $\sigma^2 = \text{Var } X$ : It holds that  $\text{Var } X = \mathbf{E} X^2 (\mathbf{E} X)^2$  and therefore

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

3. Estimator for skewness and kurtosis  $\hat{\gamma}_1$  and  $\hat{\gamma}_2$  (cf. [33, Section 6.4.4]) are Plug-in estimators, since the coefficient of skewness is defined as

$$\gamma_1 = \mathrm{E}\left(\frac{X-\mu}{\sigma}\right)^3$$

where  $\mu = EX$ ,  $\sigma^2 = Var X$ , implies

$$\hat{\gamma}_1 \underset{\sigma^2 \mapsto \hat{\sigma}^2}{\overset{\mu \to \bar{X}_n}{=}} \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^3}{(\hat{\sigma}_n^2)^{3/2}} = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^3}{\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2\right)^{3/2}}.$$

The construction of  $\hat{\gamma}_2$  can be done in the same spirit.

4. The empirical coefficient of correlation  $\varrho_{XY}$  is a Plug-in estimator,

$$\hat{\varrho}_{XY} = \frac{S_{XY}^2}{\sqrt{S_{XX}^2} \sqrt{S_{YY}^2}} = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n)}{\sqrt{\sum_{i=1}^n (X_i - \bar{X}_n)^2 \sum_{i=1}^n (Y_i - \bar{Y}_n)^2}}.$$

Indeed

$$\varrho_{XY} = \frac{\mathrm{E}(X - \mathrm{E}X)(Y - \mathrm{E}Y)}{\sqrt{\mathrm{Var}\,X \cdot \mathrm{Var}\,Y}} = \frac{\mathrm{E}(XY) - \mathrm{E}X \cdot \mathrm{E}Y}{\sqrt{(\mathrm{E}X^2 - (\mathrm{E}X)^2)(\mathrm{E}Y^2 - (\mathrm{E}Y)^2)}}$$

and therefore, considering the linear functionals

$$T_1(F) = \int x \, dF(x), \ T_2(F) = \int x^2 \, dF(x), \ T_{12}(G) = \int xy \, dG(x,y)$$
$$\varrho_{XY} = \frac{T_{12}(F_{XY}) - T_1(F_X) \cdot T_1(F_Y)}{\sqrt{(T_2(F_X) - (T_1(F_X))^2) (T_2(F_Y) - (T_1(F_Y))^2)}}.$$

 $\hat{\varrho}_{XY}$  is obtained by replacing  $F_X$ ,  $F_Y$  and  $F_{XY}$  in  $T_1$ ,  $T_2$  and  $T_{12}$  with  $\hat{F}_{n,X}$ ,  $\hat{F}_{n,Y}$  and  $\hat{F}_{n,XY}$ :

$$\hat{\varrho}_{XY} = \frac{T_{12}(\hat{F}_{n,XY}) - T_{1}(\hat{F}_{n,X}) \cdot T_{1}(\hat{F}_{n,Y})}{\sqrt{\left(T_{2}(\hat{F}_{n,X}) - \left(T_{1}(\hat{F}_{n,X})\right)^{2}\right)\left(T_{2}(\hat{F}_{n,Y}) - \left(T_{1}(\hat{F}_{n,Y})\right)^{2}\right)}}.$$

# 1.2.2 Method of moments estimator

In the following let  $(X_1, ..., X_n)$  be a sample of i.i.d. random variables  $X_i$  with distribution function  $F \in \{F_{\theta} : \theta \in \Theta\}$ ,  $\Theta \subset \mathbb{R}^m$  (parametric model). Assume that the parametrisation  $\theta \mapsto F_{\theta}$  is distinguishable, i.e.  $F_{\theta} \neq F_{\theta'} \iff \theta \neq \theta'$ .

**Goal:** Construction of an estimator  $\hat{\theta}(X_1, \ldots, X_n)$  for  $\theta = (\theta_1, \ldots, \theta_m)$ . [33, Theorem 4.5.6<sup>2</sup>] implies that under certain conditions on F (e.g. uniform distribution on a compact interval) the underlying distribution can be determined, if the moments  $EX^k$ ,  $k \in \mathbb{N}$  are known. The *method of moments estimation* is based on the idea of estimating F by using the moments and was introduced by Karl Pearson in the end of the 19th century.

**Assumptions:** There exists  $r \geq m$  such that  $E_{\theta}|X_i|^r < \infty$ . Assume that the moments  $E_{\theta}X_i^k = g_k(\theta), k = 1, ..., r$  are given as functions of the parameter vector  $\theta = (\theta_1, ..., \theta_m) \in \Theta$ .

$$P(X \in C) = 1.$$

If  $C \subset [a, b]$ , a < b, then  $\{\mu_k\}_{k \in \mathbb{N}}$  defines  $P_X$  uniquely.

**Theorem 4.5.6.** Let X be a random variable with values in  $C \subset \mathbb{R}$ , i.e.

**Moment equation system:**  $\hat{\mu}_k = g_k(\theta), k = 1, \dots, r$ , where  $\hat{\mu}_k$  are the k-th empirical moments defined by  $\hat{\mu}_k = \frac{1}{n} \sum_{i=1}^n X_i^k$ .

**Definition 1.2.7.** If the system above is uniquely solvable for  $\theta$ , then its solution  $\hat{\theta}(X_1, \ldots, X_n)$  is called *moment estimator* (*M-estimator*) of  $\theta$ .

**Lemma 1.2.8.** Let  $g = (g_1, \ldots, g_r) : \Theta \to C \subset \mathbb{R}^r$  be a bijective function, and let its inverse function  $g^{-1} : C \to \Theta$  be continuous. Then the moment estimator  $\hat{\theta}(X_1, \ldots, X_n)$  of  $\theta$  is strongly consistent.

**Proof** It holds that  $\hat{\theta}(X_1, \dots, X_n) = g^{-1}(\hat{\mu}_1, \dots, \hat{\mu}_r) \xrightarrow[n \to \infty]{\text{a.s.}} \theta$ , since  $\hat{\mu}_k \xrightarrow[n \to \infty]{\text{a.s.}} g_k(\theta)$ ,  $k = 1, \dots, r$  (strong consistency of the empirical moments) and  $g^{-1}$  is continuous.

#### Remark 1.2.9.

1. Under certain conditions with respect to the regularity of  $F_{\theta}$  the moment estimator  $\hat{\theta}(X_1, \ldots, X_n)$  for  $\theta$  is asymptotically normally distributed:

$$\sqrt{n}\left(\hat{\theta}(X_1,\ldots,X_n)-\theta\right) \xrightarrow[n\to\infty]{d} N(0,\Sigma),$$

where  $N(0, \Sigma)$  is the multivariate normal distribution with covariance matrix

$$\Sigma = G^T \mathbf{E}(YY^T)G$$

with

$$Y = (X, X^2, \dots, X^r)^T, \quad X \stackrel{d}{=} X_i$$

and

$$G = \left(\frac{\partial g_i^{-1}}{\partial \theta_j}\right)_{\substack{i=1\dots r,\\j=1\dots m}},$$

- 2. Other properties for the moment estimator do not hold in general (e.g. not all moment estimators are unbiased (cf. Example 1.2.10, 1)).
- 3. Sometimes r > m equations in the moment equation system are necessary in order to obtain a moment estimator. It can occur for example, if some  $g_i = const$ , i.e they do not provide additional information about  $\theta$  (cf. Example 1.2.10, 2)).

## Example 1.2.10.

1. Normal distribution:  $X_i \stackrel{d}{=} X$ , i = 1, ..., n,  $X \sim N(\mu, \sigma^2)$ ; The goal is to obtain a moment estimator for  $\mu$  and  $\sigma^2$ , so  $\theta = (\mu, \sigma^2)$ . It holds that

$$g_1(\mu, \sigma^2) = \mathcal{E}_{\theta} X = \mu,$$
  

$$g_2(\mu, \sigma^2) = \mathcal{E}_{\theta} X^2 = \operatorname{Var}_{\theta} X + (\mathcal{E}_{\theta} X)^2 = \sigma^2 + \mu^2.$$

Consider the system of equations

$$\begin{cases} \frac{1}{n} \sum_{i=1}^{n} X_i = \mu, \\ \frac{1}{n} \sum_{i=1}^{n} X_i^2 = \mu^2 + \sigma^2. \end{cases}$$

It follows that

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

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

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

hence, the moment estimators are given by  $\hat{\mu} = \bar{X}_n$ ,  $\hat{\sigma}^2 = \frac{n-1}{n} S_n^2$ . Note that  $\hat{\sigma}^2$  is not unbiased, since

$$E_{\theta}\hat{\sigma}^2 = \frac{n-1}{n} \cdot E_{\theta} S_n^2 = \frac{n-1}{n} \sigma^2.$$

2. Uniform distribution:  $X_i \stackrel{d}{=} X$ , i = 1, ..., n,  $X \sim U[-\theta, \theta]$ ,  $\theta > 0$ . The goal is to obtain a moment estimator for  $\theta$ . It holds that

$$g_1(\theta) = \mathcal{E}_{\theta} X = 0,$$
  
 $g_2(\theta) = \mathcal{E}_{\theta} X^2 = \operatorname{Var}_{\theta} X = \frac{(\theta - (-\theta))^2}{12} = \frac{(2\theta)^2}{12} = \frac{\theta^2}{3}.$ 

Thus, the following system of equations can be set up:

$$\begin{cases} \frac{1}{n} \sum_{i=1}^{n} X_i = 0 & \text{(useless)}, \\ \frac{1}{n} \sum_{i=1}^{n} X_i^2 = \frac{\theta^2}{3}. \end{cases}$$

Solving the above system of equations for  $\theta$  yields the moment estimator  $\hat{\theta} = \sqrt{\frac{3}{n} \sum_{i=1}^{n} X_i^2}$ . Here, two equations for the estimation of one parameter  $\theta$  were necessary, i.e. r = 2 > m = 1.

#### 1.2.3 Maximum-likelihood estimator

Maximum-likelihood estimators were discovered by Carl Friedrich Gauss (beginning of the 19th century) and Sir Ronald Fisher (1922). Assume that all distributions in the parametric family  $\{F_{\theta}: \theta \in \Theta\}$  are either discrete or continuous.

**Definition 1.2.11.** Consider the random sample  $X = (X_1, \ldots, X_n)$ .

1. Let  $X_i$ , i = 1, ..., n, be absolutely continuous random variables with probability density function  $f_{\theta}(x)$ . Then,

$$L(x_1, \dots, x_n, \theta) = \prod_{i=1}^n f_{\theta}(x_i), \qquad (x_1, \dots, x_n) \in \mathbb{R}^n, \quad \theta \in \Theta$$

is called *likelihood funktion* of the sample  $(x_1, \ldots, x_n)$ .

2. Let  $X_i$ , i = 1, ..., n, be discrete random variables with probability mass function  $p_{\theta}(x) = P_{\theta}(X_i = x)$ ,  $x \in C$ , where C is the range of X. Then,

$$L(x_1, \dots, x_n, \theta) = \prod_{i=1}^n p_{\theta}(x_i), \qquad (x_1, \dots, x_n) \in \mathbb{C}^n, \quad \theta \in \Theta$$

is called *Likelihood function* of the sample  $(x_1, \ldots, x_n)$ .

By this definition

- the discrete case yields  $L(x_1,\ldots,x_n,\theta)=P_{\theta}(X_1=x_1,\ldots,X_n=x_n)$
- the continuous case yields

$$L(x_1, \dots, x_n, \theta) \prod_{i=1}^n \Delta x_i$$

$$= f_{(X_1, \dots, X_n), \theta}(x_1, \dots, x_n) \Delta x_1 \cdot \dots \cdot \Delta x_n$$

$$\approx P_{\theta}(X_1 \in [x_1, x_1 + \Delta x_1], \dots, X_n \in [x_n, x_n + \Delta x_n])$$

for 
$$\Delta x_i \to 0$$
,  $i = 1, \ldots, n$ .

The goal is to construct an estimator  $\theta$  such that the probability

$$P_{\theta}(X_1 = x_1, \dots, X_n = x_n)$$
 resp.  $P_{\theta}(X_i \in [x_i, x_i + \Delta x_i], i = 1, \dots, n)$ 

is maximized. This procedure is called Maximum-likelihood method.

**Definition 1.2.12.** Assume that the maximization problem given by  $L(x_1, \ldots, x_n, \theta) \mapsto \max_{\theta \in \Theta}$  is uniquely solvable. Then,

$$\hat{\theta}(x_1,\ldots,x_n) = \operatorname*{argmax}_{\theta \in \Theta} L(x_1,\ldots,x_n,\theta)$$

is called Maximum-Likelihood estimator of  $\theta$  (ML estimator).

## Remark 1.2.13.

1. There are only very few cases in which the ML estimator  $\hat{\theta}$  for  $\theta$  is explicitly expressible. In most cases, the constant factor of the likelihood function is omitted. By taking the logarithm of the remaining function

$$\log L(x_1,\ldots,x_n,\theta)$$

the so called *log-likelihood function* is obtained. Consequently

$$\prod_{i=1}^{n} f_{\theta}(x_i) \quad \text{resp.} \quad \prod_{i=1}^{n} p_{\theta}(x_i)$$

turn into sums

$$\sum_{i=1}^{n} \log f_{\theta}(x_i) \quad \text{resp.} \quad \sum_{i=1}^{n} \log p_{\theta}(x_i) \,,$$

which are easier to differentiate with respect to  $\theta$ . To compute the maximum of the log-likelihood function, one considers the first order conditions

$$\frac{\partial \log L(x_1, \dots, x_n, \theta)}{\partial \theta_j} = 0, \qquad j = 1, \dots, m,$$

which are a necessary condition for an extremum of  $\log L$  (and thus of L, since the logarithm is monotonically increasing). If this system is uniquely solvable and the obtained solution  $\hat{\theta}(X_1, \ldots, X_n)$  is a maximum, it is also the ML estimator.

2. In most applied cases, the ML estimators need to be calculated by numerical methods.

# Example 1.2.14.

1. Bernoulli distribution: Let  $X_i \sim Bernoulli(p)$  i.i.d., i = 1, ..., n, with  $p \in [0, 1]$ . Since

$$X_i = \begin{cases} 1, & \text{with probability } p, \\ 0, & \text{else,} \end{cases}$$

where the respective probability mass function is given by

$$p_{\theta}(x) = p^{x}(1-p)^{1-x}, \quad x \in \{0, 1\}.$$

The likelihood function of the random sample  $(X_1, \ldots, X_n)$  is given by

$$L(x_1, \dots, x_n, \theta) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$
$$= p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i} \stackrel{\text{def.}}{=} h(p).$$

- (a) If  $\sum_{i=1}^{n} x_i = 0$  ( $\iff x_1 = x_2 = \dots = x_n = 0$ ), then  $h(p) = (1-p)^n$  is maximized at p = 0. The ML estimator is then given by  $\hat{p}(0, \dots, 0) = 0$ .
- (b) If  $\sum_{i=1}^{n} x_i = n \iff x_1 = x_2 = \dots = x_n = 1$ , then  $h(p) = p^n$  is maximized at p = 1. The ML estimator is then given by  $\hat{p}(1, 1, \dots, 1) = 1$ .
- (c) If  $0 < \sum_{i=1}^{n} x_i < n$ , then

$$\log L(x_1, ..., x_n, p) = n\bar{x}_n \log p + n(1 - \bar{x}_n) \log(1 - p) = n \cdot g(p).$$

Since  $g(p) \xrightarrow[p \to 0.1]{} -\infty$  and

$$\frac{\partial g(p)}{\partial p} = \frac{\bar{x}_n}{p} + \frac{1 - \bar{x}_n}{1 - p} \cdot (-1) = \frac{\bar{x}_n}{p} + \frac{\bar{x}_n - 1}{1 - p} = 0$$

 $\iff$   $(1-p)\bar{x}_n + (\bar{x}_n - 1)p = 0 \iff p = \bar{x}_n$ , the continuity of g implies that g attains exactly one  $\operatorname{argmax}_p g(p) = \bar{x}_n$ .

Thus, the ML estimator is given by  $\hat{p}(X_1, \ldots, X_n) = \bar{X}_n$ .

2. Uniform distribution: Let  $X_i \sim U[0, \theta]$ , i = 1, ..., n, i.i.d. with  $\theta > 0$ . The goal is to obtain a ML estimator for  $\theta$ . It holds that

$$f_{X_i}(x) = 1/\theta \cdot I(x \in [0, \theta]), \qquad i = 1, \dots, n.$$

Thus, the likelihood function is given by

$$L(x_1, \dots, x_n, \theta) = \begin{cases} (1/\theta)^n, & 0 \le x_1, \dots, x_n \le \theta \\ 0, & \text{else} \end{cases}$$
$$= \begin{cases} (1/\theta)^n, & \text{if } \min\{x_1, \dots, x_n\} \ge 0 \\ & \text{and } \max\{x_1, \dots, x_n\} \le \theta \\ 0, & \text{else} \end{cases}$$
$$= g(\theta), \quad \theta > 0.$$

Therefore,  $\hat{\theta} = \operatorname{argmax}_{\theta>0} g(\theta) = \max\{x_1, \dots, x_n\} = x_{(n)}$ . So the ML estimator is given by  $\hat{\theta}(X_1, \dots, X_n) = X_{(n)}$ .

It can be shown that under certain conditions, the ML estimator is weakly consistent and asymptotically normal distributed.

### **Definition 1.2.15.** Let

$$L(x,\theta) = \begin{cases} f_{\theta}(x), & \text{if continuous,} \\ p_{\theta}(x), & \text{if discrete} \end{cases}$$

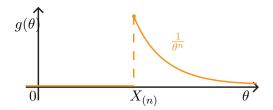


Figure 1.5: Illustration of the function g.

be the likelihood function of X. For  $\theta, \theta' \in \Theta$  and  $X \stackrel{d}{=} X_i$ ,  $P_{\theta}(L(X, \theta') = 0) = 0$  define the Kullback-Leibler information (distance)  $H(P_{\theta}, P_{\theta'})$  in the continuous case as

$$H(P_{\theta}, P_{\theta'}) = \mathcal{E}_{\theta} \log L(X, \theta) - \mathcal{E}_{\theta} \log L(X, \theta') = \int_{\mathbb{R}} \log \frac{L(x, \theta)}{L(x, \theta')} \cdot L(x, \theta) dx.$$

If  $P_{\theta}(L(X, \theta') = 0) > 0$ , then define  $H(P_{\theta}, P_{\theta'}) = \infty$ . In the discrete case take the sum over all non-trivial  $p_{\theta}(x)$  instead of the integral.

The following lemma will show that  $H(\cdot, \cdot)$  has the properties  $H(P_{\theta}, P_{\theta'}) = 0 \iff \theta = \theta'$  and  $H(P_{\theta}, P_{\theta'}) \geq 0 \quad \forall \theta, \theta' \in \Theta$ . It is, on the other hand, easy to prove that  $H(P_{\theta}, P_{\theta'})$  is not symmetric with respect to  $\theta$  and  $\theta'$ . Thus,  $H(\cdot, \cdot)$  is not a metric.

#### Lemma 1.2.16. It holds that

- 1.  $H(P_{\theta}, P_{\theta'})$  is well-defined and  $\geq 0$ .
- 2. If  $H(P_{\theta}, P_{\theta'}) = 0$ , then  $\theta = \theta'$ .

**Proof** Consider the continuous case  $P_{\theta}$ ,  $\theta \in \Theta$  (the discrete case can be shown in the same spirit).

1. Define

$$f(x) = \begin{cases} \frac{L(x,\theta)}{L(x,\theta')}, & \text{if } L(x,\theta') > 0, \\ 1, & \text{else.} \end{cases}$$

If  $P_{\theta}(L(X, \theta') = 0) = 0$ , then  $P_{\theta}(L(X, \theta') > 0) = 1$ . On the other hand, if  $H(P_{\theta}, P_{\theta'}) = \infty > 0$ , then H is well-defined and positive. With probability 1 it holds that  $L(x, \theta) = f(x) \cdot L(x, \theta')$ .

Let  $g(x) = 1 - x + x \log x$ , x > 0. It can be shown that g is convex with  $g(x) \ge 0$ . Indeed, it holds that

$$q'(x) = -1 + \log x + 1 = \log x$$
,  $q''(x) = 1/x > 0$ .

Thus, g admits exactly one zero at x = 1, which is a minimum. Consider g(f(X)),  $X \sim L(x, \theta')$ . Then,

$$0 \leq \mathcal{E}_{\theta'}g(f(X)) = 1 - \mathcal{E}_{\theta'}f(X) + \mathcal{E}_{\theta'}(f(X)\log f(X))$$

$$= 1 - \int \frac{L(x,\theta)}{L(x,\theta')} \cdot L(x,\theta') dx + \int \frac{L(x,\theta)}{L(x,\theta')} \cdot \log \frac{L(x,\theta)}{L(x,\theta')} \cdot L(x,\theta') dx$$

$$= H(P_{\theta}, P_{\theta'}).$$

Therefore,  $H(P_{\theta}, P_{\theta'}) \geq 0$ , which was to be shown.

2. If  $H(P_{\theta}, P_{\theta'}) = 0 \implies E_{\theta'}g(f(X)) = 0$ ,  $g(f(X)) \geq 0$ . Thus,  $L(x, \theta')$ -almost surely  $g(f(X)) = 0 \implies f(X) \stackrel{\theta'\text{-a.s.}}{=} 1$ , which implies either  $L(X, \theta') = 0$  or  $L(x, \theta) = L(x, \theta')$  for  $L(x, \theta')$ -almost all x and therefore  $P_{\theta} = P_{\theta'}$ .

Example 1.2.17.

1. Let  $\Theta = \mathbb{R}_+$  and  $\{P_{\lambda}, \lambda > 0\}$  be the family of exponential distributions with parameter  $\lambda > 0$  and probability density functions  $L(x,\lambda) = \lambda e^{-\lambda x} I(x \geq 0)$ . Computing the Kullback-Leibler information  $H(P_{\lambda}, P'_{\lambda})$  for any  $\lambda, \lambda' > 0$  yields

$$H(P_{\lambda}, P'_{\lambda}) = \int_{0}^{\infty} \log \left( \frac{\lambda e^{-\lambda x}}{\lambda' e^{-\lambda' x}} \right) \lambda e^{-\lambda x} dx$$

$$= \log \left( \frac{\lambda}{\lambda'} \right) \cdot \underbrace{\int_{0}^{\infty} \lambda e^{-\lambda x} dx}_{=1} - (\lambda - \lambda') \underbrace{\int_{0}^{\infty} x \lambda e^{-\lambda x} dx}_{=\frac{1}{\lambda}}$$

$$= \log \left( \frac{\lambda}{\lambda'} \right) - \frac{\lambda - \lambda'}{\lambda}$$

$$= \frac{\lambda'}{\lambda} - 1 - \log \left( \frac{\lambda'}{\lambda} \right).$$

For  $\lambda = \lambda'$  we get  $H(P_{\lambda}, P_{\lambda}) = 1 - 1 - \log(1) = 0$ .

2. It may also happen that  $H(P_{\theta}, P_{\theta'}) = +\infty$  for absolutely continuous distributions  $P_{\theta}$ . As an example, consider the family  $\{U[0, \theta], \theta > 0\}$  of uniform distributions on  $[0, \theta]$  with the likelihood  $L(x, \theta) = \frac{I(x \in [0, \theta])}{\theta}$ .

Then,

$$H(P_{\theta}, P_{\theta'}) = \frac{1}{\theta} \int_0^{\theta} \log \left( \frac{\frac{1}{\theta}}{\frac{1}{\theta'} I(x \in [0, \theta'])} \right) dx$$
$$= \frac{1}{\theta} \int_0^{\theta} \log \left( \frac{\theta'}{\theta I(x \in [0, \theta'])} \right) dx$$
$$= \begin{cases} \log \left( \frac{\theta'}{\theta} \right), & \text{if } \theta' \ge \theta, \\ +\infty, & \text{if } \theta' < \theta. \end{cases}$$

**Theorem 1.2.18** (Weak consistency of ML estimators). Let m = 1 and  $\Theta$  be an open interval in  $\mathbb{R}$ . Furthermore, let  $L(x_1, \ldots, x_n, \theta)$  be unimodal, i.e. for the ML estimator  $\hat{\theta}$  of  $\theta$  it holds that

$$\begin{cases} \forall \, \theta < \hat{\theta}(x_1, \dots, x_n) & \Longrightarrow L(x_1, \dots, x_n, \theta) \text{ is increasing} \\ \forall \, \theta > \hat{\theta}(x_1, \dots, x_n) & \Longrightarrow L(x_1, \dots, x_n, \theta) \text{ is decreasing} \end{cases}$$

(i.e.  $\max_{\theta \in \Theta} L(x_1, \dots, x_n, \theta)$  exists and is unique). Then,

$$\hat{\theta}(X_1,\ldots,X_n) \xrightarrow[n\to\infty]{P} \theta.$$

**Proof** For the weak consistency (the convergence in probability) of  $\hat{\theta}$  to hold, the following needs to be shown:

$$P_{\theta}\left(\left|\hat{\theta}(X_1,\dots,X_n) - \theta\right| > \varepsilon\right) \underset{n \to \infty}{\longrightarrow} 0, \qquad \varepsilon > 0.$$
 (1.3)

Let  $\varepsilon > 0$ :  $\theta \pm \varepsilon \in \Theta$  be arbitrary. Then, the Kullback-Leibler information satisfies  $H(P_{\theta}, P_{\theta \pm \varepsilon}) > \sigma > 0$ , because of the distinguishability of the parametrization of  $P_{\theta}$  and Lemma 1.2.16. Consider  $\{|\hat{\theta} - \theta| \le \varepsilon\}$ . In order to show (1.3), it is sufficient to find a lower bound for  $P_{\theta}(|\hat{\theta} - \theta| \le \varepsilon)$ , which converges to 1 for  $n \to \infty$ . By unimodality it holds that

$$\begin{cases}
|\hat{\theta} - \theta| < \varepsilon \\
& \stackrel{\text{unimod}}{\supseteq} \left\{ L(X_1, \dots, X_n, \theta) \in \left( L(X_1, \dots, X_n, \theta - \varepsilon), L(X_1, \dots, X_n, \theta + \varepsilon) \right)^3 \right\} \\
&= \cup \left\{ \frac{L(X_1, \dots, X_n, \theta)}{L(X_1, \dots, X_n, \theta \pm \varepsilon)} > 1 \right\} \stackrel{\sigma > 0 \Rightarrow e^{n\delta} > 1}{\supseteq} \cup \left\{ \frac{L(X_1, \dots, X_n, \theta)}{L(X_1, \dots, X_n, \theta \pm \varepsilon)} > e^{n\delta} \right\} \\
&= \cup \left\{ \frac{1}{n} \log \frac{L(X_1, \dots, X_n, \theta)}{L(X_1, \dots, X_n, \theta \pm \varepsilon)} > \sigma \right\} = A_+ \cup A_-,$$

where

$$A_{\pm} = \left\{ \frac{1}{n} \log \frac{L(X_1, \dots, X_n, \theta)}{L(X_1, \dots, X_n, \theta \pm \varepsilon)} > \sigma \right\}.$$

<sup>&</sup>lt;sup>3</sup>This means an interval with these endpoints, even though we don't immediately know which one is larger.

Hence,

$$P_{\theta}(|\hat{\theta} - \theta| < \varepsilon) \ge P_{\theta}(A_{+} \cup A_{-}).$$

Showing that

$$\lim_{n \to \infty} P_{\theta}(A_{\pm}) = 1 \tag{1.4}$$

then implies

$$1 \ge \lim_{n \to \infty} P_{\theta}(A_+ \cup A_-) \ge \lim_{n \to \infty} P_{\theta}(A_{\pm}) = 1,$$

in particular this yields

$$\lim_{n \to \infty} P_{\theta}(A_+ \cup A_-) = 1$$

and

$$1 \ge \lim_{n \to \infty} P_{\theta} \left( |\hat{\theta} - \theta| < \varepsilon \right) \ge 1.$$

which implies that

$$\lim_{n \to \infty} P_{\theta} \left( |\hat{\theta} - \theta| > \varepsilon \right) \le 1 - \underbrace{\lim_{n \to \infty} P_{\theta} \left( |\hat{\theta} - \theta| < \varepsilon \right)}_{=1} = 0,$$

i.e. 
$$\hat{\theta} \xrightarrow[n \to \infty]{P} \theta$$
.

In the following it will now be shown that  $P_{\theta}(A_{+}) \xrightarrow[n \to \infty]{} 1$  (similar for  $P_{\theta}(A_{-}) \xrightarrow[n \to \infty]{} 1$ ):

1. Let  $H(P_{\theta}, P_{\theta+\varepsilon}) < \infty$  and

$$f(x) = \begin{cases} \frac{L(x,\theta)}{L(x,\theta+\varepsilon)}, & \text{if } L(x,\theta+\varepsilon) > 0, \\ 1, & \text{else.} \end{cases}$$

By Definition 1.2.15, it holds that  $P_{\theta}(L(X_1, \theta + \varepsilon) > 0) = 1$ . Furthermore, the strong law of large numbers implies

$$\frac{1}{n}\log\frac{L(X_1,\ldots,X_n,\theta)}{L(X_1,\ldots,X_n,\theta+\varepsilon)} = \frac{1}{n}\sum_{i=1}^n\log\frac{L(X_i,\theta)}{L(X_i,\theta+\varepsilon)} = \frac{1}{n}\sum_{i=1}^n\log f(X_i)$$

$$\xrightarrow[n\to\infty]{\text{f.s.}} \text{E}_{\theta}\log f(X_1) = \int L(x,\theta)\cdot\log\frac{L(x,\theta)}{L(x,\theta+\varepsilon)}\,dx = H(P_{\theta},P_{\theta+\varepsilon}) > \sigma > 0,$$

since  $\log f(X_1) \in L^1(\Omega, \mathcal{F}, P)$  and

$$E_{\theta} \log f(X_1) = H(P_{\theta}, P_{\theta+\varepsilon}) < \infty \implies P(A_+) \underset{n \to \infty}{\longrightarrow} 1.$$

2. Let  $H(P_{\theta}, P_{\theta+\varepsilon}) = \infty$  and  $P_{\theta}(L(X_1, \theta + \varepsilon) = 0) = 0$ , then

$$f(x) \stackrel{\text{a.s.}}{=} \frac{L(x,\theta)}{L(x,\theta+\varepsilon)}$$

with respect to the distribution  $P_{X_1}$ . Now,  $\log \min\{f(X_1), c\} \in L^1(\Omega, \mathcal{F}, P)$  for all c > 0. Thus, similarly to 1 it holds that

$$\frac{1}{n} \sum_{i=1}^{n} \log \min\{f(X_i), c\} \xrightarrow[n \to \infty]{\text{a.s.}} \operatorname{E}_{\theta} \log \min\{f(X_1), c\} \in (0, \infty)$$

$$\xrightarrow[c \to \infty]{} H(P_{\theta}, P_{\theta + \varepsilon}) = \infty$$

and therefore

$$A_{+} \supset \left\{ \frac{1}{n} \sum_{i=1}^{n} \log \min\{f(X_{i}), c\} > \sigma \right\}$$

$$\Longrightarrow P(A_{+}) \ge P\left(\frac{1}{n} \sum_{i=1}^{n} \log \min\{f(X_{i}), c\} > \sigma\right) \xrightarrow[n \to \infty]{} 1.$$

3. Let  $H(P_{\theta}, P_{\theta+\varepsilon}) = \infty$  and  $P_{\theta}(L(X_1, \theta+\varepsilon) = 0) = a > 0$ . Then,

$$P_{\theta} \left( \frac{1}{n} \log \frac{L(X_1, \dots, X_n, \theta)}{L(X_1, \dots, X_n, \theta + \varepsilon)} = \infty \right)$$

$$= 1 - P \left( \frac{1}{n} \log \frac{L(X_1, \dots, X_n, \theta)}{L(X_1, \dots, X_n, \theta + \varepsilon)} < \infty \right)$$

$$= 1 - P \left( \bigcap_{i=1}^n \{ L(X_i, \theta + \varepsilon) > 0 \} \right)$$

$$\stackrel{X_i \text{ i.i.d.}}{=} 1 - (1 - a)^n \underset{n \to \infty}{\longrightarrow} 1.$$

In summary,  $P(A_+) \xrightarrow[n \to \infty]{} 1$ .

**Definition 1.2.19.** Let  $X = (X_1, ..., X_n)$  be a random sample of i.i.d. random variables  $X_i \sim F_\theta$ ,  $\theta \in \Theta$ . Let  $L(x, \theta)$  be the likelihood function of  $X_i$ . Then,

$$I(\theta) = \mathcal{E}_{\theta} \left( \frac{\partial}{\partial \theta} \log L(X_1, \theta) \right)^2, \quad \theta \in \Theta$$
 (1.5)

is called the *Fisher information* of the sample  $(X_1, \ldots, X_n)$ .

From now on it will be assumed that  $0 < I(\theta) < \infty$ . In the following some necessary conditions with respect to the asymptotically normal distribution of the ML estimator will be presented.

- 1.  $\Theta \subset \mathbb{R}$  is an open interval (m=1).
- 2. It holds that  $P_{\theta} \neq P_{\theta'}$  if and only if  $\theta \neq \theta'$ .
- 3. The family  $\{P_{\theta}, \theta \in \Theta\}$ ,  $\theta \in \Theta$  consists only of discrete or continuous distributions and no mixtures.
- 4.  $B = \text{supp } L(x, \theta) = \{x \in \mathbb{R} : L(x, \theta) > 0\}$  does not depend on  $\theta \in \Theta$ . Here supp f denotes the support of f, which is defined as

$$\operatorname{supp} f = \{ x \in \mathbb{R} : f(x) \neq 0 \},\$$

and the likelihood function  $L(x, \theta)$  is given by

$$L(x,\theta) = \begin{cases} p(x,\theta), & \text{in the discrete case,} \\ f(x,\theta), & \text{in the continuous case,} \end{cases}$$
 (1.6)

where  $p(x, \theta)$  resp.  $f(x, \theta)$  denotes the probability mass or density function of  $P_{\theta}$ .

5. The mapping  $L(x,\theta)$  is three times continuously differentiable and

$$0 = \frac{d^k}{d\theta^k} \int_B L(x,\theta) \, dx = \int_B \frac{\partial^k}{\partial \theta^k} L(x,\theta) \, dx \,, \quad k = 1, 2, \, \theta \in \Theta \,.$$

Since the integral of  $L(x, \theta)$  is equal to 1, the above derivative is equal to 0. In the discrete case, the integral is replaced by a sum over all values  $x \in \mathbb{R}$  with positive probability mass  $p(x, \theta) > 0$ .

6. For all  $\theta_0 \in \Theta$  there exists a constant  $\sigma_{\theta_0} > 0$  and a measurable function  $g_{\theta_0}: B \to [0, \infty)$ , such that

$$\left| \frac{\partial^3 \log L(x, \theta)}{\partial \theta^3} \right| \le g_{\theta_0}(x), \quad \forall x \in B, \quad |\theta - \theta_0| < \sigma_{\theta_0},$$

where  $E_{\theta_0} g_{\theta_0}(X_1) < \infty$ .

# Remark 1.2.20. It holds that

$$n \cdot I(\theta) = \operatorname{Var}_{\theta} \left( \frac{\partial}{\partial \theta} \log L(X_1, \dots, X_n, \theta) \right),$$

where

$$L(X_1, \dots, X_n, \theta) = \prod_{i=1}^n L(X_i, \theta)$$
(1.7)

is the likelihood function of the sample  $(X_1, \ldots, X_n)$  with  $L(X_i, \theta)$  given in (1.6).

**Proof** Note that

$$\frac{\partial}{\partial \theta} \log L(X_1, \dots, X_n, \theta) = \frac{\partial}{\partial \theta} \sum_{i=1}^n \log L(X_i, \theta) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log L(X_i, \theta)$$
$$= \sum_{i=1}^n \frac{L'(X_i, \theta)}{L(X_i, \theta)}.$$

Furthermore

$$E_{\theta} \left( \frac{\partial}{\partial \theta} \log L(X_1, \dots, X_n, \theta) \right) = \sum_{i=1}^n E_{\theta} \frac{L'(X_i, \theta)}{L(X_i, \theta)}$$
$$= \sum_{i=1}^n \int_B \frac{L'(X, \theta)}{L(X, \theta)} \cdot L(X, \theta) dx$$
$$\stackrel{5)}{=} 0.$$

In summary

$$\begin{aligned} \operatorname{Var}_{\theta} \left( \frac{\partial}{\partial \theta} \log L(X_{1}, \dots, X_{n}, \theta) \right) &= \operatorname{Var}_{\theta} \left( \sum_{i=1}^{n} \frac{\partial}{\partial \theta} \log L(X_{i}, \theta) \right) \\ &\overset{X_{i} \text{ i.i.d.}}{=} \sum_{i=1}^{n} \operatorname{Var}_{\theta} \left( \frac{\partial}{\partial \theta} \log L(X_{i}, \theta) \right) \\ &\overset{X_{i} \text{ i.i.d.}}{=} n \cdot \operatorname{Var}_{\theta} \left( \frac{\partial}{\partial \theta} \log L(X_{1}, \theta) \right) \\ &= n \cdot \operatorname{E}_{\theta} \left( \frac{\partial}{\partial \theta} \log L(X_{1}, \theta) \right)^{2} = n \cdot I(\theta). \end{aligned}$$

**Example 1.2.21.** Let  $X_i \sim N(\mu, \sigma^2)$ , i = 1, ..., n. For  $\theta = \mu$  the Fisher information is given by  $I(\mu) = \frac{1}{\sigma^2}$  assuming that  $\sigma^2$  is known. Indeed

$$L(X_1, \mu) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(X_1 - \mu)^2}{2\sigma^2}\right\},\$$
$$\log L(X_1, \mu) = -\log(\sqrt{2\pi}\sigma) - \frac{(X_1 - \mu)^2}{(2\sigma^2)},\$$
$$\frac{\partial \log L(X_1, \mu)}{\partial \mu} = -\frac{2(X_1 - \mu)}{2\sigma^2} \cdot (-1) = \frac{X_1 - \mu}{\sigma^2},\$$

Hence,

$$I(\mu) = \mathbb{E}_{\mu} \left( \frac{\partial \log L(X_1, \mu)}{\partial \mu} \right)^2 = \frac{1}{\sigma^4} \mathbb{E}_{\mu} (X_1 - \mu)^2 = \frac{1}{\sigma^4} \cdot \sigma^2 = \frac{1}{\sigma^2}.$$

By Remark 1.2.20 and [33, Theorem 7.3.2., 4)] with  $\hat{\mu} = \overline{X}_n$ , it holds that  $\operatorname{Var}_{\mu}\left(\frac{\partial}{\partial \mu} \log L(X_1, \dots, X_n, \mu)\right) = \frac{n}{\sigma^2} = \frac{1}{\operatorname{Var}_{\mu}(\hat{\mu})}$ . This means that little information about  $\mu$  (small values of  $I(\mu)$ ) leads to an increasing variance when estimating  $\mu$  and vice versa.

**Theorem 1.2.22.** Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables fulfilling conditions 1) to 6) and  $0 < I(\theta) < \infty$ ,  $\theta \in \Theta$ . Let  $\hat{\theta}(X_1, \ldots, X_n)$  be a weakly consistent ML estimator for  $\theta$ . Then, the ML estimator  $\hat{\theta}(X_1, \ldots, X_n)$  is also asymptotically normally distributed, in particular

 $\sqrt{n \cdot I(\theta)} \left( \hat{\theta}(X_1, \dots, X_n) - \theta \right) \xrightarrow[n \to \infty]{d} Y \sim N(0, 1).$ 

**Proof** Denote by  $l_n(\theta) = \log L(X_1, \dots, X_n, \theta)$  the log-likelihood function,  $\theta \in \Theta$ . Let  $l_n^{(k)}$  denote the k-th derivative of  $l_n$  with respect to  $\theta$ , i.e.

$$l_n^{(k)}(\theta) = \frac{d^k}{d\theta^k} l_n(\theta), \qquad k = 1, 2, 3.$$

Since  $\hat{\theta}$  is a ML estimator  $l_n^{(1)}(\hat{\theta}) = 0$  must hold. Considering the Taylor expansion of  $l_n^{(1)}(\hat{\theta})$  in a neighborhood of  $\theta$  yields

$$0 = l_n^{(1)}(\hat{\theta}) = l_n^{(1)}(\theta) + (\hat{\theta} - \theta) \cdot l_n^{(2)}(\theta) + (\hat{\theta} - \theta)^2 \cdot \frac{l_n^{(3)}(\theta^*)}{2},$$

where  $\theta^*$  is between  $\theta$  and  $\hat{\theta}$ . Note that

$$-(\hat{\theta} - \theta) \left( l_n^{(2)}(\theta) + (\hat{\theta} - \theta) \frac{l_n^{(3)}(\theta^*)}{2} \right) = l_n^{(1)}(\theta),$$

and consequently

$$\sqrt{n}(\hat{\theta} - \theta) = \frac{\frac{l_n^{(1)}(\theta)}{\sqrt{n}}}{-\frac{l_n^{(2)}(\theta)}{n} - (\hat{\theta} - \theta)\frac{l_n^{(3)}(\theta^*)}{2n}}.$$

By showing

1.  $\frac{l_n^{(1)}(\theta)}{\sqrt{n}} \xrightarrow[n \to \infty]{d} N(0, I(\theta)),$ 

2. 
$$-\frac{l_n^{(2)}(\theta)}{n} \xrightarrow[n \to \infty]{\text{a.s.}} I(\theta),$$

3.

$$(\hat{\theta} - \theta) \xrightarrow[n \to \infty]{P} 0$$
 and  $\frac{l_n^{(3)}(\theta^*)}{2n}$ 

is bounded, i.e.

there exists a 
$$c > 0$$
: 
$$\lim_{n \to \infty} P_{\theta} \left( \left| \frac{l_n^{(3)}(\theta^*)}{2n} \right| < c \right) = 1,$$

it can be followed that

$$(\hat{\theta} - \theta) \cdot \frac{l_n^{(3)}(\theta^*)}{2n} \xrightarrow[n \to \infty]{P} 0$$
, since  $\left| \frac{l_n^{(3)}(\theta^*)}{n} \right| \le c$  with high probability

and hence

$$\sqrt{n}(\hat{\theta} - \theta) = \frac{\frac{l_n^{(1)}(\theta)}{\sqrt{n}}}{-\frac{l_n^{(2)}(\theta)}{n} - (\hat{\theta} - \theta)\frac{l_n^{(3)}(\theta^*)}{2n}} \xrightarrow{n \to \infty} Z_1 \sim N\left(0, \frac{1}{I(\theta)}\right)$$

by Slutskys Theorem. Ultimately this yields  $\sqrt{n}\sqrt{I(\theta)}(\hat{\theta}-\theta) \xrightarrow[n\to\infty]{d} Y \sim N(0,1).$ 

1. The central limit Theorem implies

$$\frac{l_n^{(1)}(\theta)}{\sqrt{n}} = \frac{\sum_{i=1}^n \frac{\partial}{\partial \theta} \log L(X_i, \theta)}{\sqrt{n}} \xrightarrow[n \to \infty]{d} Y_1 \sim N\left(0, \underbrace{\operatorname{Var}_{\theta}\left(\frac{\partial}{\partial \theta} L(X_i, \theta)\right)}_{=I(\theta)}\right)$$

since  $\frac{\partial}{\partial \theta} \log L(X_i, \theta)$  are i.i.d. random variables with expectation 0 (cf. Remark 1.2.20).

2.

$$\begin{split} -\frac{1}{n}l_{n}^{(2)}(\theta) &= -\frac{1}{n}\sum_{i=1}^{n}\frac{\partial^{2}}{\partial\theta^{2}}\log L(X_{i},\theta) \\ &= \frac{1}{n}\sum_{i=1}^{n}\frac{\left(L^{(1)}(X_{i},\theta)\right)^{2} - L(X_{i},\theta) \cdot L^{(2)}(X_{i},\theta)}{\left(L(X_{i},\theta)\right)^{2}} \\ &= \frac{1}{n}\sum_{i=1}^{n}\left(\frac{L^{(1)}(X_{i},\theta)}{L(X_{i},\theta)}\right)^{2} - \frac{1}{n}\sum_{i=1}^{n}\frac{L^{(2)}(X_{i},\theta)}{L(X_{i},\theta)} \\ &\xrightarrow[n\to\infty]{\text{a.s.}} \operatorname{E}_{\theta}\left(\frac{L^{(1)}(X_{1},\theta)}{L(X_{1},\theta)}\right)^{2} - \operatorname{E}_{\theta}\left(\frac{L^{(2)}(X_{1},\theta)}{L(X_{1},\theta)}\right) = I(\theta) \end{split}$$

by the law of large numbers, where

$$L^{(k)}(X_i, \theta) = \frac{\partial^k}{\partial \theta^k} L(X_i, \theta)$$

is the k-th derivative of the likelihood function with respect to  $\theta$ , and

$$E_{\theta}\left(\frac{L^{(2)}(X_1,\theta)}{L(X_1,\theta)}\right) = \int_{B} \frac{\partial^2}{\partial \theta^2} L(x,\theta) dx \stackrel{5)}{=} \frac{d^2}{d\theta^2} \int_{B} L(x,\theta) dx = 0.$$

3. By the weak consistency of  $\hat{\theta}$  we have  $\hat{\theta} \xrightarrow[n \to \infty]{P} \theta$ . Following this it can be shown that

$$\frac{l_n^{(3)}(\theta^*)}{n}(\hat{\theta}-\theta) \xrightarrow[n \to \infty]{P} 0.$$

Note that  $\hat{\theta} \xrightarrow[n \to \infty]{P} \theta$  implies that for all  $\varepsilon > 0$ 

$$P\left(|\hat{\theta} - \theta| \le \varepsilon\right) \underset{n \to \infty}{\longrightarrow} 1$$
,

which means that  $|\hat{\theta} - \theta| \leq \sigma_{\theta}$ , with high probability  $\sigma_{\theta} > 0$ , as required in Condition 6. Thus, for all  $\theta$  with  $|\hat{\theta} - \theta| < \sigma_{\theta}$ 

$$\left| \frac{l_n^{(3)}(\theta)}{n} \right| \leq \frac{1}{n} \sum_{i=1}^n \underbrace{\left| \frac{\partial^3}{\partial \theta^3} \log L(X_i, \theta) \right|}_{\leq g_{\theta}(X_i)} \leq \frac{1}{n} \sum_{i=1}^n g_{\theta}(X_i) \xrightarrow[n \to \infty]{\text{a.s.}} E_{\theta} g_{\theta}(X_1) < \infty.$$

Consequently, there exists a constant c > 0 such that

$$P_{\theta}\left(\left|\frac{l_n^{(3)}(\theta^*)}{n}\right| < c\right) \underset{n \to \infty}{\longrightarrow} 1,$$

and hence

$$\frac{l_n^{(3)}(\theta^*)}{n}(\hat{\theta}-\theta) \underset{n\to\infty}{\overset{P}{\longrightarrow}} 0.$$

#### 1.2.4 Bayesian estimation

Let  $(X_1, \ldots, X_n)$  be a random sample, where  $X_i$  are i.i.d. random variables with distribution function  $F_{\theta}$ ,  $\theta \in \Theta$ . The distribution  $F_{\theta}$  can be either discrete or continuous. Additionally, let  $\theta$  be a realization of a random variable  $\tilde{\theta}$  with distribution  $Q(\cdot)$  on the measurable space  $(\Theta, \mathcal{B}_{\Theta})$ , which is either discrete with probability mass function  $q(\cdot)$  on absolutely continuous with probability density function  $q(\cdot)$ . As usual, both cases will be handled simultaneously with integration being replaced by summation in the discrete case.

**Definition 1.2.23.** The distribution  $Q(\cdot)$  is called *prior or apriori distribution* of the parameter  $\theta$  (of  $\tilde{\theta}$ ) (prior means "prior to the experiment  $(X_1, \ldots, X_n)$ ").

**Definition 1.2.24.** The posterior distribution of the parameter  $\theta$  (of  $\tilde{\theta}$ ) is given by the probability mass/density function  $q_{X_1,...,X_n}(\theta, X_1,...,X_n)$ , which is defined by

$$\begin{cases} P(\tilde{\theta} = \theta \mid X_1 = x_1, \dots, X_n = x_n), & \text{if } Q \text{ is discrete,} \\ f_{\tilde{\theta} \mid X_1, \dots, X_n}(\theta, x_1, \dots, x_n), & \text{if } Q \text{ is continuous.} \end{cases}$$

Here,

$$P(\tilde{\theta} = \theta \mid X_1 = x_1, \dots, X_n = x_n) = \frac{P(\tilde{\theta} = \theta, X_1 = x_1, \dots, X_n = x_n)}{P(X_1 = x_1, \dots, X_n = x_n)}$$

$$= \frac{P_{\theta}(X_i = x_i, i = 1, \dots, n) \cdot q(\theta)}{\sum_{\theta_1 \in \Theta} P_{\theta_1}(X_i = x_i, i = 1, \dots, n) \cdot q(\theta_1)}$$

by the *Bayes formula*, resp.

$$f_{\tilde{\theta}|X_1,\dots,X_n}(\theta,x_1,\dots,x_n) = \frac{f_{(\tilde{\theta},X_1,\dots,X_n)}(\theta,x_1,\dots,x_n)}{f_{X_1,\dots,X_n}(x_1,\dots,x_n)}$$
$$= \frac{L(x_1,\dots,x_n,\theta) \cdot q(\theta)}{\int_{\Theta} L(x_1,\dots,x_n,\theta_1) \cdot q(\theta_1) d\theta_1},$$

where  $L(x_1, \ldots, x_n, \theta)$  is the likelihood funtion defined in (1.7).

**Definition 1.2.25.** A loss function  $V: \Theta^2 \to \mathbb{R}_+$  is a  $\Theta^2$  measurable function.

Loss functions are used as follows: Denote by  $E_*V(\tilde{\theta}, a)$  the expected loss (mean risk), which occurs from estimating  $\theta$  with a, where  $E_*$  is the expectation with respect to the posterior distribution of  $\tilde{\theta}$ . Note that  $E_*V(\tilde{\theta}, a)$  is a function of a and  $x_1, \ldots, x_n$ , since the sample  $(x_1, \ldots, x_n)$  is an explicit part of the posteriori distribution. In particular, it holds that

$$E_*V(\tilde{\theta}, a) = \varphi(x_1, \dots, x_n, a).$$

**Definition 1.2.26.** An estimator  $\hat{\theta}$  is called a *Bayes estimator* of  $\theta$ , if

$$\hat{\theta}(x_1, \dots, x_n) = \underset{a}{\operatorname{argmin}} \, \mathcal{E}_* V(\tilde{\theta}, a) \tag{1.8}$$

exists and is unique.

#### Remark 1.2.27.

1. Sometimes  $\hat{\theta} \notin \Theta$ , which is attributable to  $\varphi(x_1, \dots, x_n, a)$  attaining its minimum outside of  $\Theta$ .

2. The name "Bayes approach" honors the English mathematician Thomas Bayes (1702–1761), who only introduced the idea behind the Bayes formula given by

$$P(B_i|A) = \frac{P(A|B_i) \cdot P(B_i)}{\sum_{i} P(A|B_j) \cdot P(B_j)}.$$
 (1.9)

The actual discovery of (1.9) was by Pierre-Simon Laplace (1749–1827) (end of the 18th century). The formula was explicitly used in the derivation of the *posterior distribution* of  $\tilde{\theta}$ .

3. The approach in Definition 1.2.26 is usually only realizable by numeric minimization. There are only very few cases where an analytic solution of the minimization problem stated in (1.8) can be computed.

**Example 1.2.28** (Quadratic loss function). If  $V(\theta_1, \theta_2) = (\theta_1 - \theta_2)^2$ , then

$$\underset{a}{\operatorname{argmin}} (\varphi(x_1, \dots, x_n, a)) = \underset{a}{\operatorname{argmin}} \left( \mathbb{E}_* (\tilde{\theta} - a)^2 \right)$$
$$= \underset{a}{\operatorname{argmin}} \left( \mathbb{E}_* \tilde{\theta}^2 - 2a \mathbb{E}_* \tilde{\theta} + a^2 \right)$$
$$= \mathbb{E}_* \tilde{\theta}.$$

The Bayes estimator of  $\hat{\theta}(x_1,\ldots,x_n)$  for  $\theta$  is thus given by  $\mathbf{E}_*\tilde{\theta}$ .

**Example 1.2.29** (Bernoulli distribution). Let  $(X_1, \ldots, X_n)$  be an i.i.d. random sample of random variables  $X_i \sim Bernoulli(p), p \in (0,1)$ . Furthermore let  $\tilde{p} \sim Beta(\alpha, \beta), \alpha, \beta > 0$  be the prior distribution, with probability mass function

$$q(p) = \frac{p^{\alpha - 1}(1 - p)^{\beta - 1}}{B(\alpha, \beta)} \cdot I_{[0,1]}(p).$$

The posterior distribution of  $\tilde{p}$  is then given by

$$q^*(p) = f_{\tilde{p}|X_1 = x_1, \dots, X_n = x_n}(p) = \frac{P_p(X_1 = x_1, \dots, X_n = x_n) \cdot q(p)}{\int_0^1 P_{p_1}(X_1 = x_1, \dots, X_n = x_n) \cdot q(p_1) \, dp_1}.$$

It is always possible to calculate the posterior distribution with respect to a function  $g(X_1, \ldots, X_n)$  instead of the vector  $(X_1, \ldots, X_n)$ .

Here,  $Y = g(X_1, ..., X_n) = \sum_{i=1}^n X_i$  denotes the number of successful trials within n experiments, where

$$X_i = \begin{cases} 1, & \text{with probability } p, \\ 0, & \text{else.} \end{cases}$$

Therefore,

$$\begin{split} q^*(p) &= f_{\tilde{p}|Y=k}(p) = \frac{P_p(Y=k) \cdot q(p)}{\int_0^1 P_{p_1}(Y=k)q(p_1) \, dp_1} \\ &\stackrel{Y \sim Bin(n,p),}{=} \frac{\binom{n}{k}p^k(1-p)^{n-k} \cdot (B(\alpha,\beta))^{-1} \cdot p^{\alpha-1}(1-p)^{\beta-1}}{\frac{\binom{n}{k}}{B(\alpha,\beta)} \cdot \int_0^1 p_1^{k+\alpha-1}(1-p_1)^{n-k+\beta-1} \, dp_1} \\ &= \frac{p^{k+\alpha-1}(1-p)^{n-k+\beta-1}}{B(k+\alpha,n-k+\beta)} \,, \qquad p \in [0,1] \,. \end{split}$$

holds for the posterior distribution with respect to Y. Hence, the posterior distribution of  $\tilde{p}$  under the condition Y = k is given by

$$Beta(k + \alpha, n - k + \beta).$$

For the Bayes estimator it holds that

$$\hat{p}(x_1, \dots, x_n) = \mathcal{E}_* \tilde{p} = \int_0^1 p \cdot q^*(p) \, dp = \frac{\int_0^1 p^{k+\alpha} (1-p)^{n-k+\beta-1} \, dp}{B(k+\alpha, n-k+\beta)}$$

$$= \frac{B(k+\alpha+1, n-k+\beta)}{B(k+\alpha, n-k+\beta)}$$

$$= \dots = \frac{k+\alpha}{\alpha+\beta+n}$$

$$= \frac{\sum_{i=1}^n x_i + \alpha}{\alpha+\beta+n}$$

$$= \frac{\alpha+n\bar{x}_n}{\alpha+\beta+n}.$$

Interpretation:

$$\hat{p}(X_1, \dots, X_n) = \underbrace{\frac{n}{\alpha + \beta + n}}_{=:c_1} \bar{X}_n + \underbrace{\frac{\alpha + \beta}{\alpha + \beta + n}}_{=:c_2} \cdot \frac{\alpha}{\alpha + \beta} = c_1 \cdot \bar{X}_n + c_2 \cdot \mathbf{E}_{apr} \tilde{\theta},$$

where  $c_1 + c_2 = 1$ . This means that the Bayes method is a middle ground between the estimator  $E_{apr}\tilde{\theta}$  (with no information about the sample  $(X_1, \ldots, X_n)$ ) and the estimator  $\bar{X}_n$  (with no information about the prior distribution of  $\tilde{p}$ ) for p.

# 1.2.5 Resampling methods for obtaining point estimators

Let  $(X_1, \ldots, X_n)$  be a random sample in a parametric model. The goal is to find an estimator  $\hat{\theta}$  for the parameter  $\theta$ . In order to construct this estimator, resampling methods will be applied, i.e. generating a new sample  $(X_1^*, \ldots, X_n^*)$  by randomly drawing from the old random sample  $(X_1, \ldots, X_n)$  independently with replacement. After resampling the sample mean, sample

variance and other estimators with respect to the new sample can be computed. In this case the dimension m of the parameter space  $\Theta$  is arbitrary. The following resampling methods are introduced:

- 1. Jackknife, which is supposed to imply its handiness in every situation
- 2. Bootstrap, which is supposed to imply its self-sufficiency
- 1. Jackknife methods for estimating the variance or the bias of estimators As an introductory example, consider  $\theta = EX = \mu$  or  $\theta = Var X = \sigma^2$  and the respective (unbiased) estimators  $\hat{\mu} = \bar{X}_n$  or  $\hat{\sigma}^2 = S_n^2$ .

It is already known that

$$\operatorname{Var} \hat{\mu} = \frac{\sigma^2}{n}, \qquad \operatorname{Var} \hat{\sigma}^2 = \frac{1}{n} \left( \mu_4' - \frac{n-3}{n-1} \sigma^4 \right).$$

Now an estimator of the variance of  $\hat{\mu}$  resp.  $\hat{\sigma}^2$  is desired. In order to do so, the plug-in methods are useful:

$$\widehat{\operatorname{Var}}\, \widehat{\mu} = \frac{S_n^2}{n} \,, \qquad \widehat{\operatorname{Var}}\, \widehat{\sigma}^2 = \frac{1}{n} \left( \widehat{\mu}_4' - \frac{n-3}{n-1} S_n^4 \right) \,,$$

where  $\hat{\mu}'_4$  is the fourth centered empirical moment.

In general there are no explicit formulas for  $\operatorname{Var} \hat{\theta}$  known. That is where the jackknife method comes into play.

• Let  $X_{[i]}$  be the random sample given by  $(X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_n)$ ,  $i = 1, \ldots, n$ . Let

$$\hat{\theta}(X_1,\ldots,X_n) = \varphi_n(X_1,\ldots,X_n),$$

and set

$$\hat{\theta}_{[i]} = \varphi_{n-1}(X_{[i]}), \quad \bar{\theta}_{[\cdot]} = \frac{1}{n} \sum_{i=1}^{n} \hat{\theta}_{[i]}$$

$$\widehat{\operatorname{Var}}_{jn}(\hat{\theta}) \stackrel{\text{def.}}{=} \frac{n-1}{n} \sum_{i=1}^{n} \left( \hat{\theta}_{[i]} - \bar{\theta}_{[\cdot]} \right)^{2}.$$

**Definition 1.2.30.** The estimator  $\bar{\theta}_{[\cdot]}$  resp.  $\widehat{\mathrm{Var}}_{jn}(\hat{\theta})$  is called a *jackknife estimator* for the expectation resp. variance of the estimator  $\hat{\theta}$  of  $\theta$ .

**Example 1.2.31.** Let  $\theta = \mu$ ,  $\hat{\theta} = \hat{\mu} = \bar{X}_n$ . Then

$$\varphi_n(x_1,\ldots,x_n) = \frac{1}{n} \sum_{i=1}^n x_i,$$

which implies that

$$\hat{\theta}_{[i]} = \frac{1}{n-1} \sum_{j \neq i} X_j = \frac{1}{n-1} \left( -X_i + \sum_{j=1}^n X_j \right)$$

$$= \frac{n}{n-1} \bar{X}_n - \frac{1}{n-1} X_i , \quad \forall i = 1, \dots, n ,$$

$$\bar{\theta}_{[\cdot]} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{[i]} = \frac{n}{n-1} \bar{X}_n - \frac{1}{n(n-1)} \sum_{i=1}^n X_i$$

$$= \frac{n \cdot \bar{X}_n}{n-1} - \frac{\bar{X}_n}{n-1} = \frac{n-1}{n-1} \bar{X}_n = \bar{X}_n .$$

Thus, a jackknife estimator for  $\mu$  is equal to  $\bar{X}_n$ . Construction of a jackknife estimator for the variance:

$$\widehat{\text{Var}}_{jn}(\widehat{\theta}) = \frac{n-1}{n} \sum_{i=1}^{n} \left( \frac{n}{n-1} \bar{X}_n - \frac{1}{n-1} X_i - \bar{X}_n \right)^2$$

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

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

$$= \frac{1}{n} S_n^2,$$

which is exactly the plug-in estimator for the variance of  $\hat{\mu}$ .

• jackknife for the bias of an estimator Let  $\hat{\theta}(X_1, ..., X_n)$  be an estimator for  $\theta$ . The bias of  $\hat{\theta}$  given by  $E_{\theta}\hat{\theta} - \theta = \text{Bias}(\hat{\theta})$ .

**Definition 1.2.32.** A jackknife estimator of the bias of  $\hat{\theta}$  is given by

$$\widehat{\mathrm{Bias}}_{jn}(\hat{\theta}) = (n-1)(\bar{\theta}_{[\cdot]} - \hat{\theta}).$$

The following examples show that the procedure above leads to a decreasing bias: The estimator

$$\tilde{\theta} = \hat{\theta} - \widehat{\text{Bias}}_{jn}(\hat{\theta}) = n\hat{\theta} - (n-1)\bar{\theta}_{[\cdot]}$$
 (1.10)

generally has a smaller bias than  $\hat{\theta}$ . Here

$$\hat{\theta}_{[i]} = \varphi_{n-1}(X_{[i]})$$
 and  $\bar{\theta}_{[\cdot]} = \frac{1}{n} \sum_{i=1}^{n} \hat{\theta}_{[i]}$ 

with

$$\hat{\theta}(X_1,\ldots,X_n)=\varphi_n(X_1,\ldots,X_n).$$

## Example 1.2.33.

- (a) Let  $\theta = \mathbf{E}X_i = \mu$ . Then  $\hat{\theta} = \bar{X}_n$  is an unbiased estimator for  $\mu$ . Now the question for the corrected bias estimator  $\tilde{\mu}$  arises? (It is not supposed to be any worse!) It holds that  $\bar{\theta}_{[\cdot]} = \bar{X}_n$  thus, the bias estimator of jackknife  $\widehat{\mathrm{Bias}}_{jn}(\hat{\theta}) = (n-1)(\bar{X}_n \bar{X}_n) = 0$ , and therefore  $\tilde{\theta} = \hat{\theta} 0 = \bar{X}_n$ . Hence, the jackknife method does not (at least in this example) add additional bias.
- (b)  $\theta = \sigma^2 = \text{Var}X_i$ ,  $\hat{\theta} = \hat{\sigma}^2 = \frac{1}{n}\sum_{i=1}^n (X_i \bar{X}_n)^2$  is a biased moment estimator of the variance. The question of how  $\tilde{\theta}$  looks like arises.

**Exercise 1.2.34.** Show that the bias corrected estimator  $\tilde{\theta}$  is an unbiased estimator of the variance:

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

It follows that the bias of  $\hat{\sigma}^2$  is completely removed by applying the jackknife method.

Idea of the proof: Show that

$$\widehat{\text{Bias}}_{jn}(\hat{\theta}) = -\frac{1}{n(n-1)} \sum_{i=1}^{n} (X_i - \bar{X}_n)^2.$$

Remark 1.2.35. The examples 1.2.33 a), b), which provided a jackknife estimator in analytic form are rather an exception. In most cases, the reduction of the bias is achieved by using Monte-Carlo methods on the basis of (1.10).

#### 2. Bootstrap estimator

The bootstrap method draws a new random sample  $(X_1^*, \ldots, X_n^*)$  from an approximate distribution  $\hat{F}$  of the random sample variables  $X_i$ ,  $i = 1, \ldots, n$ . Let  $E_*$  and  $Var_*$  be the expectation and variance with respect to the distribution  $P_*$  of  $(X_1^*, \ldots, X_n^*)$ . There are two possibilities for the construction of  $\hat{F}$ :

- i)  $\hat{F}(x) = \hat{F}_n(x)$ , which is the empirical distribution of  $X_i$ , if  $X_i$  are i.i.d.
- ii)  $\hat{F}$ , which is a parametric estimator of the parametric distribution F, of  $X_i$ . That means, if  $X_i \sim F_\theta$ ,  $i = 1, \ldots, n$  for a  $\theta \in \Theta$  and  $\hat{\theta} = \hat{\theta}(X_1, \ldots, X_n)$  an estimator for  $\theta$ , then  $\hat{F} = F_{\hat{\theta}}$  (plug-in method).

**Definition 1.2.36.** A bootstrap estimator for the expectation (resp. bias or variance) of the estimator  $\hat{\theta}(X_1, \dots, X_n)$  is given by

- (a)  $\hat{E}_{boot}(\hat{\theta}) = E_* \hat{\theta}(X_1^*, \dots, X_n^*).$
- (b)  $\widehat{\text{Bias}}_{boot}(\hat{\theta}) = \hat{\text{E}}_{boot}\hat{\theta} \hat{\theta}$ .
- (c)  $\widehat{\operatorname{Var}}_{boot}(\hat{\theta}) = \operatorname{Var}_*(\hat{\theta}(X_1^*, \dots, X_n^*)).$

**Example 1.2.37.** Let  $\theta = \mu = EX_i$  and  $\hat{F} = \hat{F}_n$  be the empirical distribution function. How is a random sample  $X_1^*, \dots, X_n^*$  with  $X_i^* \sim \hat{F}_n$  generated?

The empirical distribution function  $\hat{F}_n$  weighs every observation  $x_i$  of the original sample with a weight 1/n. As a consequence, it is sufficient to select one of the entries in  $(x_1, \ldots, x_n)$  (with probability 1/n, urn model "drawing with replacement"), in order to generate  $X_j^*$ ,  $j = 1, \ldots, n$ .

Bootstrap estimator for the expectation  $\hat{\mu} = \bar{X}_n$ :

$$\hat{E}_{boot}\hat{\mu} = E_* \left(\frac{1}{n} \sum_{i=1}^n X_i^*\right) \stackrel{X_i^* \text{ i.i.d.}}{=} \frac{1}{n} \cdot n E_*(X_1^*)$$
$$= \int x \, d\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n.$$

It follows that  $\widehat{\text{Bias}}_{boot}\hat{\mu} = 0$ . Moreover,

$$\widehat{\operatorname{Var}}_{boot}(\widehat{\mu}) = \operatorname{Var}_* \left( \frac{1}{n} \sum_{i=1}^n X_i^* \right)^{X_i^*} \stackrel{\text{u.i.v.}}{=} \frac{1}{n^2} \cdot n \cdot \operatorname{Var}_*(X_1^*)$$
$$= \frac{1}{n} \cdot \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{\widehat{\sigma}^2}{n} ,$$

is a plug-in estimator for  $Var \bar{X}_n = \sigma^2/n$ .

Monte-Carlo methods for constructing bootstrap estimators numerically:

What can be done, if there is no explicit expression of  $\widehat{\text{Var}}_{Boot}(\hat{\theta})$  (which is usually the case in statistics)?

Generate M independent random samples  $(X_{i1}^*, \ldots, X_{in}^*), i = 1, \ldots, M$  under i) or ii) by using Monte-Carlo simulation. Then,

$$\hat{\theta}_i = \hat{\theta}(X_{i1}^*, \dots, X_{in}^*), \quad i = 1, \dots, M \quad \text{and set} \quad \hat{\mathbf{E}}_{boot} \hat{\theta} \approx \frac{1}{M} \sum_{i=1}^M \hat{\theta}_i.$$

Similarily a bootstrap estimator for Bias  $\hat{\theta}$  and  $\operatorname{Var} \hat{\theta}$  is obtained:

$$\widehat{\text{Bias}}_{boot} \hat{\theta} \approx \hat{\text{E}}_{boot} \hat{\theta} - \hat{\theta}, \quad \widehat{\text{Var}}_{boot} \hat{\theta} \approx \frac{1}{M-1} \sum_{i=1}^{M} (\hat{\theta}_i - \hat{\text{E}}_{boot} \hat{\theta})^2.$$

Furthermore, the distribution function of  $X_{ij}^*$  can be determined by the empirical distribution function, i.e.

$$\hat{F}_{boot}(x) = \frac{1}{M} \sum_{i=1}^{M} \frac{1}{n} \sum_{j=1}^{n} I(X_{ij}^* \le x), \quad x \in \mathbb{R}.$$

Using the methods above the *Bootstrap confidence intervals* for  $\hat{\theta}$  can be constructed:

The quantiles  $\hat{F}_{\hat{\theta}}^{-1}(\alpha_1)$  and  $\hat{F}_{\hat{\theta}}^{-1}(1-\alpha_2)$  of the distribution of  $\hat{\theta}(X_1^*,\ldots,X_n^*)$  originating from the sample  $(\hat{\theta}_1,\ldots,\hat{\theta}_M)$  can be estimated empirically. Then

$$P\left(\hat{F}_{\hat{\theta}}^{-1}(\alpha_1) \le \hat{\theta}(X_1^*, \dots, X_n^*) \le \hat{F}_{\hat{\theta}}^{-1}(1 - \alpha_2)\right) \approx 1 - \alpha_1 - \alpha_2 = 1 - \alpha,$$

where  $\alpha = \alpha_1 + \alpha_2$  is sufficiently small. Note that it is desired that  $X_i^*$  are similarly distributed as the  $X_i$ , and hence

$$P\left(\hat{F}_{\hat{\theta}}^{-1}(\alpha_1) \le \hat{\theta}(X_1, \dots, X_n) \le \hat{F}_{\hat{\theta}}^{-1}(1 - \alpha_2)\right) \approx 1 - \alpha$$

holds.

## 1.3 Further quality properties of point estimators

## 1.3.1 Cramér-Rao inequality

Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables  $X_i$  with distribution function  $F_{\theta}$ ,  $\theta \in \Theta \subset \mathbb{R}$ , i.e., m = 1. Let  $\hat{\theta}(X_1, \ldots, X_n)$  be an estimator for  $\theta$ . If  $\hat{\theta}$  is unbiased, then the quality of another unbiased estimator  $\tilde{\theta}$  of  $\theta$  is determined by the its variance. That means, if  $\operatorname{Var}_{\theta} \tilde{\theta} < \operatorname{Var}_{\theta} \hat{\theta}$  then  $\tilde{\theta}$  is in a sense better. This section strives to answer the question, whether it is always possible to find a newer, better estimator  $\tilde{\theta}$  with decreasing variance. Under certain conditions this is not possible. The lower bound for  $\operatorname{Var}_{\theta} \hat{\theta}$  is given by the Cramér-Rao Theorem. Let  $L(x,\theta)$  be the likelihood function of  $X_i$ , i.e.

$$L(x,\theta) = \begin{cases} P_{\theta}(x), & \text{in the discrete case,} \\ f_{\theta}(x), & \text{in the absolutely continuous case,} \end{cases}$$

and  $L(x_1, \ldots, x_n, \theta) = \prod_{i=1}^n L(x_i, \theta)$  the likelihood function of the whole random sample  $(X_1, \ldots, X_n)$ . The conditions 1) to 5) for the asymptotically normal distribution on page 24 hold, where 5) holds for k = 1.

**Theorem 1.3.1** (Inequality of Cramér-Rao). Let  $\hat{\theta}(X_1, \ldots, X_n)$  be an estimator for  $\theta$  with the following properties:

- 1.  $E_{\theta}\hat{\theta}^2(X_1,\ldots,X_n) < \infty \quad \forall \theta \in \Theta.$
- 2. For all  $\theta \in \Theta$  exists  $\frac{d}{d\theta} \mathbb{E}_{\theta} \hat{\theta}(X_1, \dots, X_n)$ , given by

$$\begin{cases} \int_{\mathbb{R}^n} \hat{\theta}(x_1, \dots, x_n) \frac{\partial}{\partial \theta} L(x_1, \dots, x_n, \theta) dx_1 \dots dx_n, & \text{in the abs. cont. case }, \\ \sum_{x_1, \dots, x_n} \hat{\theta}(x_1, \dots, x_n) \frac{\partial}{\partial \theta} L(x_1, \dots, x_n, \theta), & \text{in the discrete case} \end{cases}$$

Then, a lower bound for the variance of  $\hat{\theta}$  is attained, i.e.

$$\operatorname{Var}_{\theta} \hat{\theta}(X_1, \dots, X_n) \ge \frac{\left(\frac{d}{d\theta} \operatorname{E}_{\theta} \hat{\theta}(X_1, \dots, X_n)\right)^2}{n \cdot I(\theta)}, \quad \theta \in \Theta,$$

where  $I(\theta)$  is the Fisher information defined in (1.5).

## **Proof** Let

$$\varphi_{\theta}(x_1, \dots, x_n) = \frac{\partial}{\partial \theta} \log L(x_1, \dots, x_n, \theta).$$

In Remark 1.2.20 it has been shown that

$$E_{\theta}\varphi_{\theta}(X_1,\ldots,X_n)=0$$
,  $Var_{\theta}\varphi_{\theta}(X_1,\ldots,X_n)=n\cdot I(\theta)$ .

Applying the Cauchy-Schwartz inequality to

Cov 
$$_{\theta}(\varphi_{\theta}(X_1,\ldots,X_n),\hat{\theta}(X_1,\ldots,X_n))$$

yields

Cov 
$$_{\theta}\left(\varphi_{\theta}(X_{1},\ldots,X_{n}),\hat{\theta}(X_{1},\ldots,X_{n})\right)$$
  

$$= \operatorname{E}_{\theta}\left(\varphi_{\theta}(X_{1},\ldots,X_{n})\cdot\hat{\theta}(X_{1},\ldots,X_{n})\right) - 0$$

$$\leq \sqrt{\operatorname{Var}_{\theta}\varphi_{\theta}(X_{1},\ldots,X_{n})}\sqrt{\operatorname{Var}_{\theta}\hat{\theta}(X_{1},\ldots,X_{n})}$$

Thus,

$$\operatorname{Var}_{\theta} \hat{\theta}(X_{1}, \dots, X_{n}) \geq \frac{\left( \underbrace{\operatorname{E}_{\theta} \left( \varphi_{\theta}(X_{1}, \dots, X_{n}) \cdot \hat{\theta}(X_{1}, \dots, X_{n}) \right)}^{=:A} \right)^{2}}{\operatorname{Var}_{\theta} \varphi_{\theta}(X_{1}, \dots, X_{n})} = \frac{A^{2}}{n \cdot I(\theta)}.$$

Now it suffices to show

$$A = \frac{d}{d\theta} \mathcal{E}_{\theta} \, \hat{\theta}(X_1, \dots, X_n) \,.$$

Only the absolutely continuous case will be shown (in the discrete case, replace the integrals with sums):

$$A = \int \frac{\partial}{\partial \theta} \log L(x_1, \dots, x_n, \theta) \cdot \hat{\theta}(x_1, \dots, x_n) \cdot L(x_1, \dots, x_n, \theta) dx_1 \dots dx_n$$

$$= \int \frac{\partial}{\partial \theta} L(x_1, \dots, x_n, \theta) \cdot \hat{\theta}(x_1, \dots, x_n) dx_1 \dots dx_n$$

$$\stackrel{\text{Cond. 2}}{=} \frac{d}{d\theta} \operatorname{E}_{\theta} \hat{\theta}(X_1, \dots, X_n).$$

Corollary 1.3.2. If  $\hat{\theta}$  is an unbiased estimator for  $\theta$  and the conditions of Theorem 1.3.1 are fulfilled, then

$$Var_{\theta} \hat{\theta}(X_1, \dots, X_n) \ge \frac{1}{n \cdot I(\theta)}$$
.

**Proof** Apply the Cramér-Rao inequality with

$$\frac{d}{d\theta} \left( \mathcal{E}_{\theta} \, \hat{\theta}(X_1, \dots, X_n) \right) = \frac{d}{d\theta} \theta = 1.$$

The following examples will show, that the estimator  $\bar{X}_n$  of the expectation  $\mu$  has the smallest variance within the class of all estimators  $\mu$  which fulfill the conditions of Theorem 1.3.1. Hence, the sample mean  $\bar{X}_n$  is the best unbiased estimator in this class for at least two families of distributions:

- Normal distribution and
- Poisson distribution.

## Example 1.3.3.

1. Let  $X_i \sim N(\mu, \sigma^2)$  and  $\hat{\mu} = \bar{X}_n$  be an estimator for  $\mu$ . Here,  $\hat{\mu}$  is unbiased with  $\text{Var}\hat{\mu} = \sigma^2/n$ . In the following it will be shown that the Cramér-Rao boundary for the variance of an unbiased estimator  $\hat{\theta}$  for  $\mu$  is also given by  $\sigma^2/n$ . In an initial step, the conditions of Theorem 1.3.1 will be validated. In order to show that

$$0 = \frac{d}{d\mu} \int_{\mathbb{R}} L(x,\mu) \, dx = \int_{\mathbb{R}} \frac{\partial}{\partial \mu} L(x,\mu) \, dx$$

with

$$L(x,\mu) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

consider

$$\begin{split} \frac{\partial}{\partial \mu} L(x,\mu) &= \frac{2(x-\mu)}{2\sigma^2} \cdot \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} = \frac{x-\mu}{\sigma^2} \cdot L(x,\mu)\,, \\ \int_{\mathbb{R}} \frac{\partial}{\partial \mu} L(x,\mu) \, dx &= \mathrm{E}\left(\frac{X-\mu}{\sigma^2}\right) = 0\,. \end{split}$$

For condition 2) in Theorem 1.3.1 it holds that

$$\frac{d}{d\mu} \mathbf{E} \bar{X}_n = \frac{d}{d\mu}(\mu) = 1$$

$$\stackrel{?}{=} \frac{1}{n} \int_{\mathbb{R}^n} (x_1 + \dots + x_n) \frac{\partial}{\partial \mu} \left( \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma}\right)^2} \right) dx_1 \dots dx_n.$$

Induction with repsect to n:

• Initial step n = 1:

$$\int_{\mathbb{R}} x \frac{\partial}{\partial \mu} L(x, \mu) \, dx = \int_{\mathbb{R}} \frac{x(x - \mu)}{\sigma^2} L(x, \mu) \, dx$$
$$= \frac{1}{\sigma^2} \left( \mathcal{E}_{\mu} X^2 - \mu^2 \right) = \frac{\operatorname{Var}_{\mu} X}{\sigma^2} = 1 \, .$$

• Induction hypothesis: For n it holds that

$$\int_{\mathbb{R}^n} (x_1 + \ldots + x_n) \cdot \frac{\partial}{\partial \mu} L(x_1, \ldots, x_n, \mu) \, dx_1 \ldots dx_n = n.$$

• Induction step  $n \to n+1$ :

$$A = \int_{\mathbb{R}^{n+1}} (x_1 + \dots + x_{n+1}) \frac{\partial}{\partial \mu} \underbrace{L(x_1, \dots, x_{n+1}, \mu)}_{=L(x_1, \dots, x_n, \mu) \cdot L(x_{n+1}, \mu)} dx_1 \dots dx_{n+1}$$

$$\stackrel{?}{=} n + 1.$$

For A it holds that

$$A = \int_{\mathbb{R}^{n+1}} (x_1 + \dots + x_n) \cdot \left( \frac{\partial}{\partial \mu} L(x_1, \dots, x_n, \mu) \cdot L(x_{n+1}, \mu) \right)$$

$$+ L(x_1, \dots, x_n, \mu) \cdot \frac{\partial}{\partial \mu} L(x_{n+1}, \mu) dx_1 \dots dx_n dx_{n+1}$$

$$+ \int_{\mathbb{R}^{n+1}} x_{n+1} \left( \frac{\partial}{\partial \mu} L(x_1, \dots, x_n, \mu) \cdot L(x_{n+1}, \mu) \right)$$

$$+ L(x_1, \dots, x_n, \mu) \cdot \frac{\partial}{\partial \mu} L(x_{n+1}, \mu) dx_1 \dots dx_n dx_{n+1}$$

$$= n \cdot \underbrace{\int_{\mathbb{R}} L(x_{n+1}, \mu) \, dx_{n+1}}_{=1}$$

$$+ \int_{\mathbb{R}^n} (x_1 + \dots + x_n) \cdot L(x_1, \dots, x_n, \mu) \, dx_1 \dots dx_n$$

$$\cdot \underbrace{\int \frac{\partial}{\partial \mu} L(x_{n+1}, \mu) \, dx_{n+1}}_{=0} + \int_{\mathbb{R}} x_{n+1} L(x_{n+1}, \mu) \, dx_{n+1}$$

$$= 0$$

$$\cdot \underbrace{\int_{\mathbb{R}^n} \frac{\partial}{\partial \mu} L(x_1, \dots, x_n, \mu) \, dx_1 \dots dx_n}_{=0}$$

$$+ \underbrace{\int_{\mathbb{R}} x_{n+1} \frac{\partial}{\partial \mu} L(x_{n+1}, \mu) \, dx_{n+1}}_{=\frac{d}{d\mu} \mathcal{E}_{\mu} X = \frac{d}{d\mu} \mu = 1}$$

$$\cdot \underbrace{\int_{\mathbb{R}^n} L(x_1, \dots, x_n, \mu) \, dx_1 \dots dx_n}_{=1} = n + 1.$$

Since all conditions are fulfilled, the bound can be computed by

$$\frac{1}{n \cdot I(\mu)}$$

with

$$I(\mu) = \mathcal{E}_{\mu} \left( \frac{\partial}{\partial \mu} \log L(X, \mu) \right)^{2}.$$

Example 1.2.21 implies that

$$I(\mu) = \frac{1}{\sigma^2} \implies n \cdot I(\mu) = \frac{n}{\sigma^2}$$

In summary:

$$\operatorname{Var}_{\mu} \hat{\theta} \ge \frac{1}{\frac{n}{\sigma^2}} = \frac{\sigma^2}{n} = \operatorname{Var}_{\mu} \bar{X}_n$$

holds for an arbitrary estimator  $\hat{\theta}$  for  $\mu$ , which fulfills the conditions of Theorem 1.3.1.

2. The second example will be an exercise.

**Exercise 1.3.4.** Let  $X_i \sim Poisson(\lambda)$ , i = 1, ..., n. Show that the Cramér-Rao bound given by

$$\frac{1}{n \cdot I(\lambda)} = \frac{\lambda}{n} = \operatorname{Var}_{\lambda} \bar{X}_n,$$

which means that  $\bar{X}_n$  is also the best unbiased estimator fulfilling the conditions of Theorem 1.3.1.

The following example will show that it is possible to construct an estimator with a smaller variance than the one provided by the Cramér-Rao bound, if the conditions of Theorem 1.3.1 are not fulfilled.

**Example 1.3.5.** Let  $X_i \sim U[0, \theta], \theta > 0$ . Then the condition "supp  $f_{\theta}(x) = [0, \theta]$  independent of  $\theta$ " is not met. Additionally

$$0 \neq \int_{\mathbb{R}} \frac{\partial}{\partial \theta} L(x, \theta) \, dx = \int_{0}^{\theta} \left(\frac{1}{\theta}\right)' \, dx = -\frac{1}{\theta^{2}} \cdot \theta = -\frac{1}{\theta}$$

holds. Let  $\hat{\theta}$  be an unbiased estimator of  $\theta$ , then Cramér-Rao would imply that  $\operatorname{Var}_{\theta} \hat{\theta} \geq (n \cdot I(\theta))^{-1}$ , where

$$\begin{split} I(\theta) &= \mathrm{E} \left( \frac{\partial}{\partial \theta} \log L(X, \theta) \right)^2 = \int_0^\theta \frac{1}{\theta} \left( \frac{\partial}{\partial \theta} \log \left( \frac{1}{\theta} \right) \right)^2 \, dx \\ &= \frac{1}{\theta} \int_0^\theta \, dx \cdot \left( -\frac{1}{\theta} \right)^2 = \frac{1}{\theta^2} \, . \end{split}$$

Thus

$$\operatorname{Var}_{\theta} \hat{\theta} \geq \frac{\theta^2}{n}$$

would hold. Consider

$$\hat{\theta}(X_1,\ldots,X_n) = \frac{n+1}{n} \max\{X_1,\ldots,X_n\} = \frac{n+1}{n} X_{(n)}.$$

In order to show that

$$E_{\theta} \hat{\theta}(X_1, \dots, X_n) = \theta$$
 and  $Var_{\theta} \hat{\theta}(X_1, \dots, X_n) < \frac{\theta^2}{n}$ ,

compute the k-th moments  $E_{\theta}X_{(n)}^{k}$ ,  $k \in \mathbb{N}$ . It holds that

$$\begin{split} F_{X_{(n)}}(x) &= F_{X_i}^n(x) = \begin{cases} \frac{x^n}{\theta^n} \,, & x \in [0,\theta] \,, \\ 1 \,, & x \ge \theta \,, \\ 0 \,, & x < 0 \,, \end{cases} \\ f_{X_{(n)}}(x) &= F_{X_{(n)}}'(x) = \frac{nx^{n-1}}{\theta^n} \cdot I(x \in [0,\theta]) \,, \\ \mathbf{E}_{\theta} X_{(n)}^k &= \int_0^\theta x^k \frac{nx^{n-1}}{\theta^n} \, dx = \frac{n}{\theta^n} \int_0^\theta x^{n+k-1} \, dx = \frac{n \cdot \theta^{n+k}}{\theta^n \cdot (n+k)} = \frac{n\theta^k}{n+k} \,. \end{split}$$

Thus,

$$E_{\theta} \hat{\theta} = \frac{n+1}{n} \cdot E_{\theta} X_{(n)} = \frac{n+1}{n} \cdot \frac{n\theta}{n+1} = \theta$$

which means that  $\hat{\theta}$  is unbiased. Furthermore,

$$\operatorname{Var}_{\theta} \hat{\theta} = \left(\frac{n+1}{n}\right)^{2} \cdot \operatorname{Var}_{\theta} X_{(n)} = \left(\frac{n+1}{n}\right)^{2} \cdot \left(\frac{n\theta^{2}}{n+2} - \frac{n^{2}\theta^{2}}{(n+1)^{2}}\right)$$
$$= \frac{(n+1)^{2}}{n^{2}} \cdot \frac{n(n+1)^{2} - n^{2}(n+2)}{(n+2)(n+1)^{2}} \cdot \theta^{2}$$
$$= \frac{\theta^{2}}{n(n+2)} (n^{2} + 2n + 1 - n^{2} - 2n) = \frac{\theta^{2}}{n(n+2)}.$$

Ultimately, it follows that

$$\operatorname{Var}_{\theta} \hat{\theta} = \frac{\theta^2}{n(n+2)} < \frac{\theta^2}{n}$$
.

## 1.3.2 Sufficiency

Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables  $X_i$  with distribution function  $F_{\theta}$ ,  $\theta \in \Theta \subseteq \mathbb{R}^m$ . If the whole information  $\{X_1 = x_1, \ldots, X_n = x_n\}$  passes to the estimator  $\hat{\theta}(X_1, \ldots, X_n)$  of  $\theta$ , then the function

$$\hat{\theta}: \mathbb{R}^n \to \mathbb{R}^m, \qquad m \ll n^4$$

causes a loss of information, since  $(X_1, \ldots, X_n)$  can (usually) not be reconstructed from  $\hat{\theta}(X_1, \ldots, X_n)$ . The class of so-called *sufficient* estimators minimize the loss of information in a stochastic sense:

#### Definition 1.3.6.

1. Let the random variables  $X_1, \ldots, X_n$  and  $\hat{\theta}(X_1, \ldots, X_n)$  be discrete. An estimator  $\hat{\theta}$  of the parameter  $\theta$  is called *sufficient*, if

$$P_{\theta}\left(X_{1}=x_{1},\ldots,X_{n}=x_{n}\,|\,\hat{\theta}(X_{1},\ldots,X_{n})=t\right)$$

does not depend on  $\theta$ , as long as  $x_1, \ldots, x_n$  and t are in the support of  $(X_1, \ldots, X_n)$  resp.  $\hat{\theta}(X_1, \ldots, X_n)$ .

2. Let  $X_1, \ldots, X_n$  and  $\hat{\theta}(X_1, \ldots, X_n)$  be absolutely continuous. Then, the estimator  $\hat{\theta}$  is called *sufficient* for  $\theta$ , if the probability

$$P\left((X_1,\ldots,X_n)\in B\,|\,\hat{\theta}(X_1,\ldots,X_n)=t\right)$$

does not depend on  $\theta \in \Theta$  for arbitrary  $B \in \mathcal{B}_{\mathbb{R}^n}$  and  $t \in \text{supp } f_{\hat{\theta}}$ , where  $f_{\hat{\theta}}$  is the probability density function of  $\hat{\theta}$ .

## Remark 1.3.7.

<sup>&</sup>lt;sup>4</sup>in classical statistics

1. In Definition 1.3.6, 2. it holds that

$$P\left(\hat{\theta}(X_1,\ldots,X_n)=t\right)=0, \quad \forall t,$$

because of the absolute continuity of  $\hat{\theta}$ . Therefore, the conditional probability (in contrary to Definition 1.3.6, 1.) is not understood in the classical sense, but as a conditional expectation. Conditional expectations were introduced in the lecture "Probability Theory and Stochastic Processes" (Section 1.1.4).

2. Consider the likelihood function

$$L_{\hat{\theta}}(x_1, \dots, x_n, \theta) = P_{\theta} \left( X_1 = x_1, \dots, X_n = x_n \, | \, \hat{\theta}(X_1, \dots, X_n) = t \right)$$

for discrete  $X_1, \ldots, X_n$ . Definition 1.3.6 implies that a new estimator for  $\theta$  cannot be obtained from this conditional likelihood function  $L_{\theta}(x_1, \ldots, x_n, \theta)$ , since it does not depend on  $\theta$ . That means, the estimator  $\hat{\theta}$  already provides all the information about  $\theta$  obtainable from  $(x_1, \ldots, x_n)$ .

3. Let  $g: \mathbb{R}^m \to \mathbb{R}^m$  be a bijective Borel measurable function and  $\hat{\theta}(X_1, \ldots, X_n)$  a sufficient estimator of  $\theta \in \Theta \subset \mathbb{R}^m$ . Then  $g(\hat{\theta}(X_1, \ldots, X_n))$  is also a sufficient estimator for  $\theta$ . This is due to the fact that

$$\left\{\omega \in \Omega : g\left(\hat{\theta}(X_1,\ldots,X_n)\right) = t\right\} = \left\{\omega \in \Omega : \hat{\theta}(X_1,\ldots,X_n) = g^{-1}(t)\right\},\,$$

for all  $t \in \mathbb{R}^m$ .

**Lemma 1.3.8.** Assume that the random variables  $X_1, \ldots, X_n$  and  $\hat{\theta}(X_1, \ldots, X_n)$  are either all discrete or absolutely continuous with likelihood functions

$$L(x_1,\ldots,x_n,\theta) = \begin{cases} P_{\theta}(X_1 = x_1,\ldots,X_n = x_n) \,, & \text{in the discrete case,} \\ f_{X_1,\ldots,X_n}(x_1,\ldots,x_n,\theta) \,, & \text{in the abs. cont. case,} \end{cases}$$
 
$$L_{\hat{\theta}}(t,\theta) = \begin{cases} P_{\theta}(\hat{\theta}(X_1,\ldots,X_n) = t) \,, & \text{in the discrete case,} \\ f_{\hat{\theta}}(t,\theta) \,, & \text{in the abs. cont, case.} \end{cases}$$

Denote the support of L by

supp 
$$L = \{(x_1, \dots, x_n) \in \mathbb{R}^n : L(x_1, \dots, x_n, \theta) > 0\}.$$

Then, the estimator  $\hat{\theta}$  is sufficient with respect to  $\theta$  if and only if the ratio

$$\frac{L(x_1, \dots, x_n, \theta)}{L_{\hat{\theta}}(\hat{\theta}(x_1, \dots, x_n), \theta)} \tag{1.11}$$

does not depend on  $\theta$  for all  $(x_1, \ldots, x_n) \in \text{supp } L$  such that  $\hat{\theta}(x_1, \ldots, x_n) \in \text{supp } L_{\hat{\theta}}$ .

**Proof** Only the discrete case will be shown in the following.

Let  $\hat{\theta}$  be sufficient for  $\theta$ . Then, it has to be verified that (1.11) does not depend on  $\theta$  for all  $(x_1, \ldots, x_n) \in \mathbb{R}$ ,  $t \in \mathbb{R}$  and  $\theta \in \Theta$ , such that  $(x_1, \ldots, x_n) \in \text{supp} L$ . It holds that:

$$P_{\theta}(X_{1} = x_{1}, \dots, X_{n} = x_{n} | \hat{\theta}(X_{1}, \dots, X_{n}) = t)$$

$$= \frac{P_{\theta}(X_{1} = x_{1}, \dots, X_{n} = x_{n}, \hat{\theta}(X_{1}, \dots, X_{n}) = t)}{P_{\theta}(\hat{\theta}(X_{1}, \dots, X_{n}) = t)}$$

$$= \begin{cases} 0, & \text{if } \hat{\theta}(x_{1}, \dots, x_{n}) \neq t \\ \frac{P_{\theta}(X_{1} = x_{1}, \dots, X_{n} = x_{n})}{P_{\theta}(\hat{\theta}(X_{1}, \dots, X_{n}) = \hat{\theta}(x_{1}, \dots, x_{n}))}, & \text{if } \hat{\theta}(x_{1}, \dots, x_{n}) = t. \end{cases}$$

Thus (1.11) does not depend on  $\theta$ .

Can be done in the same spirit as the previous argument (backwards).

## 

## Example 1.3.9.

1. Bernoulli distribution: Let  $X_i \sim Bernoulli(p)$ ,  $p \in [0, 1]$ , i = 1, ..., n,  $\hat{p} = \bar{X}_n$  be an unbiased estimator for p. In the following, it will be shown that  $\hat{p}$  is sufficient. It holds that

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

where  $Y \sim Bin(n, p)$ . By Remark 1.3.7 3. it is sufficient to show, that Y is a sufficient estimator for p. For  $x_i \in \{0, 1\}$  i = 1, ..., n compute

$$P(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}.$$

Define the likelihood function  $L_Y$  by

$$L_Y(y,p) = \binom{n}{y} p^y (1-p)^{n-y}, \qquad y = 0, \dots, n.$$

Replacing y with the sum  $\sum_{i=1}^{n} x_i$  yields

$$\frac{L(x_1,\ldots,x_n,p)}{L_Y(\sum_{i=1}^n x_i,p)} = \frac{p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}}{(\sum_{i=1}^n x_i) p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}} = \frac{1}{(\sum_{i=1}^n x_i)}.$$

The term above obviously does not depend on p, thus Lemma 1.3.8 implies, that Y and therefore also  $\hat{p}$  are sufficient.

2. Normal distribution with known variance: Let  $X_i \sim N(\mu, \sigma^2)$ , i = 1, ..., n, with known  $\sigma^2$ . Then,  $\hat{\mu} = \bar{X}_n$  is an unbiased estimator for  $\mu$ . In the following it will be shown, that  $\hat{\mu}$  is sufficient: Considering

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

and [33, Lemma 6.4.5] imply

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

Furthermore, note that  $\hat{\mu} \sim N(\mu, \sigma^2/n)$ , hence

$$L_{\hat{\mu}}(x,\mu) = \frac{\sqrt{n}}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{n}{2} \left(\frac{x-\mu}{\sigma}\right)^2\right),$$

$$\frac{L(x_1,\dots,x_n,\mu)}{L_{\hat{\mu}}(\bar{x}_n,\mu)} = \frac{\frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{\sum_{i=1}^n (x_i - \bar{x}_n)^2 + n(\bar{x}_n - \mu)^2}{2\sigma^2}\right)}{\frac{\sqrt{n}}{\sqrt{2\pi}\sigma} \cdot \exp\left(\frac{-n(\bar{x}_n - \mu)^2}{2\sigma^2}\right)}$$

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

which is independent of  $\mu$ . Lemma 1.3.8 implies that  $\hat{\mu} = \bar{X}_n$  is an sufficient estimator for  $\mu$ .

The Neyman-Fisher factorization theorem, which will be introduced below, implies that the estimator  $(\bar{X}_n, S_n^2)$  for  $(\mu, \sigma^2)$  with unknown variance is sufficient.

**Theorem 1.3.10** (Neyman-Fisher Factorization Theorem). Under the conditions of Lemma 1.3.8 it holds that  $\hat{\theta}(X_1, \ldots, X_n)$  is a sufficient estimator for  $\theta$  if and only if there exist two measurable functions  $g: \mathbb{R}^m \times \Theta \to \mathbb{R}$  and  $h: \mathbb{R}^n \to \mathbb{R}$ , such that the following factorization of the likelihood function  $L(x_1, \ldots, x_n, \theta)$  of the random sample  $(X_1, \ldots, X_n)$  holds:

$$L(x_1,\ldots,x_n,\theta)=g\left(\hat{\theta}(x_1,\ldots,x_n),\theta\right)\cdot h(x_1,\ldots,x_n)$$

for  $(x_1, \ldots, x_n) \in \text{supp } L, \theta \in \Theta$ .

**Proof** Only the discrete case will be shown.

1. If  $\hat{\theta}$  is sufficient, then Lemma 1.3.8 implies that

$$\underbrace{\frac{L(x_1,\ldots,x_n,\theta)}{L_{\hat{\theta}}(\hat{\theta}(x_1,\ldots,x_n),\theta)}}_{=g(\hat{\theta}(x_1,\ldots,x_n),\theta)} = h(x_1,\ldots,x_n)$$

does not depend on  $\theta$ . Thus, the factorization of Neyman-Fisher holds.

2. Let  $L(x_1, \ldots, x_n, \theta) = g(\hat{\theta}(x_1, \ldots, x_n), \theta) \cdot h(x_1, \ldots, x_n)$  for all  $(x_1, \ldots, x_n) \in \text{supp } L, \theta \in \Theta$ . Furthermore, define

$$C = \{(y_1, \dots, y_n) \in \mathbb{R}^n : \hat{\theta}(y_1, \dots, y_n) = \hat{\theta}(x_1, \dots, x_n)\}$$
$$= \hat{\theta}^{-1} \left(\hat{\theta}(x_1, \dots, x_n)\right),$$

then

$$\begin{split} \frac{P_{\theta}(X_{1} = x_{1}, \dots, X_{n} = x_{n})}{\underbrace{L_{\theta}(\hat{\theta}(x_{1}, \dots, x_{n}), \theta)}_{=P_{\theta}(\hat{\theta}(X_{1}, \dots, x_{n}), \theta)} &= \frac{g(\hat{\theta}(x_{1}, \dots, x_{n}), \theta) \cdot h(x_{1}, \dots, x_{n})}{\sum_{(y_{1}, \dots, y_{n}) \in C} P_{\theta}(X_{1} = y_{1}, \dots, X_{n} = y_{n})} \\ &= \frac{g(\hat{\theta}(x_{1}, \dots, x_{n}), \theta) \cdot h(x_{1}, \dots, x_{n})}{\sum_{(y_{1}, \dots, y_{n}) \in C} g(\underbrace{\hat{\theta}(y_{1}, \dots, y_{n})}_{=\hat{\theta}(x_{1}, \dots, x_{n})}, \theta) \cdot h(y_{1}, \dots, y_{n})}_{=\hat{\theta}(x_{1}, \dots, x_{n})} \\ &= \frac{h(x_{1}, \dots, x_{n})}{\sum_{(y_{1}, \dots, y_{n}) \in C} h(y_{1}, \dots, y_{n})}, \end{split}$$

does not depend on  $\theta$ . Thus,  $\hat{\theta}$  is sufficient by Lemma 1.3.8.

Example 1.3.11.

1. Poisson distribution: Let  $X_i \sim Poisson(\lambda)$ ,  $\lambda > 0$ ,  $\hat{\lambda} = \bar{X}_n$  be an unbiased estimator for  $\lambda$ . In the following it will be shown that  $\hat{\lambda}$  is sufficient. For  $x_i \in \{0, 1, 2, \ldots\}$ ,  $i = 1, \ldots, n$  it holds that

$$L(x_1, \dots, x_n, \lambda) = \prod_{i=1}^n e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} = \frac{e^{-\lambda n} \cdot \lambda^{\sum_{i=1}^n x_i}}{x_1! \cdot \dots \cdot x_n!} = \frac{e^{-n\lambda} \lambda^{n\bar{x}_n}}{x_1! \cdot \dots \cdot x_n!},$$
  
=  $g(\bar{x}_n, \lambda) \cdot h(x_1, \dots, x_n)$ ,

where  $g(\bar{x}_n, \lambda) = e^{-n\lambda} \cdot \lambda^{n\bar{x}_n}$ ,  $h(x_1, \dots, x_n) = \frac{1}{x_1! \dots \cdot x_n!}$ . Thus,  $\hat{\lambda} = \bar{X}_n$  is sufficient by Theorem 1.3.10. 2. Exponential distribution Let  $X_i \sim Exp(\lambda)$ ,  $\lambda > 0$ ,  $\hat{\lambda} = \bar{X}_n^{-1}$  be a moment estimator for  $\lambda$ , which is not unbiased but strongly consistent, since the strong law of large numbers implies that  $\bar{X}_n \xrightarrow[n \to \infty]{a.s.} EX_i = \frac{1}{\lambda}$ . In the following it will be shown, that  $\hat{\lambda}$  is sufficient. For  $x_1 \geq 0, \ldots, x_n \geq 0$  it holds that

$$L(x_1, \dots, x_n, \lambda) = \prod_{i=1}^n \lambda e^{-\lambda x_i} = \lambda^n e^{-\lambda \sum_{i=1}^n x_i} = \lambda^n e^{-\lambda n \bar{x}_n}$$
$$= \lambda^n e^{-\frac{\lambda n}{\bar{\lambda}}} = g\left(\hat{\lambda}, \lambda\right) \cdot \underbrace{h(x_1, \dots, x_n)}_{=1},$$

where  $g(\hat{\lambda}, \lambda) = \lambda^n e^{-\frac{\lambda n}{\hat{\lambda}}}$  and  $h(x_1, \dots, x_n) \equiv 1$ . Thus,  $\hat{\lambda}$  is sufficient by Theorem 1.3.10.

**Exercise 1.3.12.** Using Theorem 1.3.10 show that the estimator  $(\bar{X}_n, S_n^2)$  is sufficient for  $(\mu, \sigma^2)$  if the random sample  $(X_1, \ldots, X_n)$  is i.i.d. with distribution  $X_i \sim N(\mu, \sigma^2)$  for all i.

**Remark 1.3.13.** An advantage of the Neyman-Fisher Theorem is, that if one wants to determine whether an estimator  $\hat{\theta}$  is sufficient, the likelihood function of  $\hat{\theta}$  does not need to be known explicitly. This is particularly important if the estimator  $\hat{\theta}$  is rather complicated and the likelihood function cannot be computed.

## 1.3.3 Completeness

**Definition 1.3.14.** An estimator  $\hat{\theta}(X_1, ..., X_n)$  of the parameter  $\theta \in \Theta \subset \mathbb{R}^m$  is called *complete*, if for an arbitrary measurable function  $g: \mathbb{R}^m \to \mathbb{R}$  with  $E_{\theta}g(\hat{\theta}(X_1, ..., X_n)) = 0$ ,  $\theta \in \Theta$  it holds that

$$g\left(\hat{\theta}\left(X_{1},\ldots,X_{n}\right)\right)\equiv0$$
.  $P_{\theta}$  - a.s. for all  $\theta\in\Theta$ .

## Remark 1.3.15.

1. Let  $g_1, g_2 : \mathbb{R}^m \to \mathbb{R}$  be functions with

$$\mathrm{E}_{\theta}\left|g_{i}\left(\hat{\theta}\left(X_{1},\ldots,X_{n}\right)\right)\right|<\infty$$

 $\forall \theta \in \Theta \text{ and }$ 

$$\mathrm{E}_{\theta}g_{1}\left(\hat{\theta}\left(X_{1},\ldots,X_{n}\right)\right)=\mathrm{E}_{\theta}g_{2}\left(\hat{\theta}\left(X_{1},\ldots,X_{n}\right)\right),$$

where  $\hat{\theta}$  is complete. Definition 1.3.14 then implies

$$g_1\left(\hat{\theta}\left(X_1,\ldots,X_n\right)\right) = g_2\left(\hat{\theta}\left(X_1,\ldots,X_n\right)\right), \text{ a.s.}$$

(Take  $g = g_1 - g_2$ ).

Conclusion: The completeness as characteristic allows a comparison between the estimators  $g_1(\hat{\theta})$  and  $g_2(\hat{\theta})$  with respect to their almost surely equality.

2. If  $\hat{\theta}$  is a complete estimator for  $\theta$ , then  $g(\hat{\theta})$  is also a complete estimator for  $\theta$  for an arbitrary measurable function  $g: \mathbb{R}^m \to \mathbb{R}^m$ .

## Example 1.3.16.

1. Bernoulli distribution: Let  $X_i \sim Bernoulli(p)$ ,  $p \in [0,1]$ . In order to show that  $\hat{p} = \bar{X}_n$  is complete, let g be an arbitrary real valued function. It is sufficient to show that  $Y = \sum_{i=1}^{n} X_i$  is complete. It holds that  $Y \sim Bin(n,p)$ , which implies that

$$E_p g(Y) = \sum_{k=0}^n g(k) \binom{n}{k} p^k (1-p)^{n-k}.$$

Furthermore,  $E_p g(Y) = 0$  if and only if

$$\sum_{k=0}^{n} g(k) \binom{n}{k} \left( \underbrace{\frac{p}{1-p}}_{-t} \right)^{k} = p_n(t) = 0$$

for  $p \in (0,1)$ , so  $t \in (0,\infty)$ . The polynomial  $p_n(t)$  is of degree n, hence

$$g(k) \binom{n}{k} = 0$$
 for all  $k$   
 $\implies g(k) = 0, \quad k = 0, \dots, n$   
 $\implies g(Y) = 0 \quad P_p$ -a.s..

Therefore, Y is complete and  $\hat{p} = \bar{X}_n$  as well.

2. Uniform distribution: Let  $X_i \sim U[0,\theta]$ ,  $i=1,\ldots,n$ . It has already been shown that the estimator  $\hat{\theta}(X_1,\ldots,X_n) = \frac{n+1}{n}X_{(n)}$  is unbiased. In order to show its completeness, it is sufficient to show that  $X_{(n)} = \max_{i=1,\ldots,n} X_i$  is complete, i.e. all measurable functions  $g: \mathbb{R} \to \mathbb{R}$  with  $E_{\theta}g(X_{(n)}) = 0$  need to fulfill  $g(X_{(n)}) = 0$  almost surely.

The probability density function of  $X_{(n)}$  is given by  $f_{X_{(n)}}(x) = \frac{nx^{n-1}}{\theta^n} \cdot I_{[0,\theta]}(x)$  by Example 1.3.5. Hence, we can compute

$$0 = \frac{d}{d\theta} \mathcal{E}_{\theta} g(X_{(n)}) = \frac{d}{d\theta} \int_{0}^{\theta} g(x) f_{X_{(n)}}(x) dx = \frac{d}{d\theta} \frac{1}{\theta^{n}} \int_{0}^{\theta} n x^{n-1} g(x) dx$$
$$= -n \frac{1}{\theta^{n+1}} \int_{0}^{\theta} g(x) n x^{n-1} dx + \frac{1}{\theta^{n}} n \theta^{n-1} g(\theta) = -\frac{n}{\theta} \underbrace{\mathcal{E}_{\theta} g(X_{(n)})}_{=0} + \frac{n}{\theta} g(\theta)$$
$$= \frac{n}{\theta} g(\theta) = 0 \text{ for all } \theta > 0 \Longrightarrow g(x) = 0, \quad x > 0.$$

It follows that  $g(X_{(n)}) = 0$  holds almost surely.

## 1.3.4 Best unbiased estimator

Following [33, Definition 7.2.9.] note that for a random sample  $(X_1, \ldots, X_n)$  with i.i.d. random variables  $X_i \sim F_\theta$ ,  $\theta \in \Theta \subset \mathbb{R}$  (m = 1), the estimator  $\hat{\theta}(X_1, \ldots, X_n)$  is called *best unbiased estimator*, if

$$E_{\theta}\hat{\theta}^2(X_1,\ldots,X_n) < \infty$$
  $E_{\theta}\hat{\theta}(X_1,\ldots,X_n) = \theta, \quad \theta \in \Theta, \text{ and }$ 

the estimator  $\hat{\theta}$  has the smallest variance among all unbiased estimators.

**Lemma 1.3.17** (Uniqueness of the best unbiased estimator). If  $\hat{\theta}$  is a best unbiased estimator for  $\theta$ , then it is unique.

**Proof** Let  $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$  be a best unbiased estimator for  $\theta$  and  $\tilde{\theta}$  another best unbiased estimator for  $\theta$ . In the following it will be shown that both estimators coincide, i.e.  $\hat{\theta} = \tilde{\theta}$ .

Ex adverso: Assume, that  $\hat{\theta} \neq \tilde{\theta}$  and consider  $\theta^* = 1/2(\hat{\theta} + \tilde{\theta})$ . Obviously  $\theta^*$  is unbiased and its variance is given by

$$Var_{\theta}\theta^* = \frac{1}{4}Var_{\theta}(\hat{\theta} + \tilde{\theta}) = \frac{1}{4}Var_{\theta}\hat{\theta} + \frac{1}{4}Var_{\theta}\tilde{\theta} + \frac{1}{2}Cov_{\theta}(\hat{\theta}, \tilde{\theta}).$$

The Cauchy-Schwartz inequality implies  $|\operatorname{Cov}_{\theta}(\hat{\theta}, \tilde{\theta})| \leq \sqrt{\operatorname{Var}_{\theta}\hat{\theta} \cdot \operatorname{Var}_{\theta}\hat{\theta}} = \operatorname{Var}_{\theta}\hat{\theta}$  and therefore

$$\operatorname{Var}_{\theta}\theta^* \leq \frac{1}{2}\operatorname{Var}_{\theta}\hat{\theta} + \frac{1}{2}\operatorname{Var}_{\theta}\hat{\theta} = \operatorname{Var}_{\theta}\hat{\theta}.$$

Since  $\hat{\theta}$  is a best unbiased estimator, it follows  $\operatorname{Var}_{\theta}\theta^* = \operatorname{Var}_{\theta}\hat{\theta}$ , and consequently  $\varrho(\hat{\theta}, \tilde{\theta}) = 1$  implies that  $\hat{\theta}$  and  $\tilde{\theta}$  are linearly dependent, i.e. there exist some constants a and b, such that  $\hat{\theta} = a\tilde{\theta} + b$ . It holds that a = 1 since  $\operatorname{Var}_{\theta}\hat{\theta} = a^2\operatorname{Var}\tilde{\theta} = \operatorname{Var}_{\theta}\hat{\theta}$ . Moreover, b = 0, because  $\hat{\theta}$  and  $\tilde{\theta}$  are unbiased:  $\theta = \operatorname{E}_{\theta}\hat{\theta} = \operatorname{E}_{\theta}\tilde{\theta} + b = \theta + b$ . Ultimately,  $\hat{\theta} = \tilde{\theta}$ , which completes the proof.  $\square$ 

**Lemma 1.3.18.** A unbiased estimator  $\hat{\theta}$  with finite second moment is the best unbiased estimator for  $\theta$  if and only if  $\text{Cov }_{\theta}(\hat{\theta}, \varphi) = 0$ ,  $\theta \in \Theta$  for an arbitrary sample function  $\varphi : \mathbb{R}^n \to \mathbb{R}$  with  $\text{E}_{\theta}\varphi(X_1, \ldots, X_n) = 0$ ,  $\forall \theta \in \Theta$ .

## **Proof**

"\(\Righta\)" Let  $\hat{\theta}$  be the best unbiased estimator for  $\theta$  and  $\varphi(X_1, \ldots, X_n)$  a sample function with  $\mathcal{E}_{\theta}\varphi(X_1, \ldots, X_n) = 0, \forall \theta \in \Theta$ . It is sufficient to show  $\mathcal{C}$ ov  $\theta(\hat{\theta}, \varphi) = \mathcal{E}_{\theta}(\hat{\theta}\varphi) = 0, \theta \in \Theta$ .

Define  $\tilde{\theta} = \hat{\theta} + a\varphi$ ,  $a \in \mathbb{R}$ . In order to compute

$$\operatorname{Var}_{\theta}\tilde{\theta} = \operatorname{Var}_{\theta}\hat{\theta} + a^{2}\operatorname{Var}_{\theta}\varphi + 2a\operatorname{Cov}_{\theta}(\hat{\theta},\varphi)$$

for  $a \in \mathbb{R}$ , let  $g(a) = a^2 \operatorname{Var}_{\theta} \varphi + 2a \operatorname{Cov}_{\theta}(\varphi, \hat{\theta})$ . For  $\operatorname{Cov}_{\theta}(\varphi, \hat{\theta}) \neq 0$  there exists an  $a \in \mathbb{R}$  with g(a) < 0. Since  $\tilde{\theta}$  is an unbiased estimator for  $\theta$  ( $\operatorname{E}_{\theta}\tilde{\theta} = \operatorname{E}_{\theta}\hat{\theta} + a\operatorname{E}_{\theta}\varphi = \theta + 0 = \theta$ ) it holds that  $\operatorname{Var}_{\theta}\tilde{\theta} \geq \operatorname{Var}_{\theta}\hat{\theta}$  for all  $a \in \mathbb{R}$ . This is a contradiction to g(a) < 0 for an  $a \in \mathbb{R}$ . Thus,  $\operatorname{Cov}_{\theta}(\varphi, \hat{\theta}) = 0, \theta \in \Theta$ .

"\(\infty\)" Let  $\hat{\theta}$  be an unbiased estimator with  $E_{\theta}\hat{\theta}^2 < \infty, \theta \in \Theta$  and  $Cov_{\theta}(\varphi, \hat{\theta}) = 0, \theta \in \Theta$  if  $E_{\theta}\varphi = 0, \theta \in \Theta$ . Let  $\tilde{\theta}$  be another unbiased estimator for  $\theta$ . In order to show that  $Var_{\theta}\tilde{\theta} \geq Var_{\theta}\hat{\theta}$ , consider

$$\tilde{\theta} = \hat{\theta} + (\underbrace{\tilde{\theta} - \hat{\theta}}_{=:\varphi}), \quad E_{\theta}\varphi = E_{\theta}\tilde{\theta} - E_{\theta}\hat{\theta} = \theta - \theta = 0, \quad \forall \theta \in \Theta.$$

It follows that

$$Var_{\theta}\tilde{\theta} = Var_{\theta}\hat{\theta} + \underbrace{Var_{\theta}\varphi}_{>0} + 2\underbrace{Cov_{\theta}(\hat{\theta},\varphi)}_{=0} \ge Var_{\theta}\hat{\theta},$$

which implies, that  $\hat{\theta}$  is the best unbiased estimator for  $\theta$ .

**Theorem 1.3.19** (Lehmann-Scheffé). Let  $\hat{\theta}$  be an unbiased, complete and sufficient estimator for  $\theta$  with  $E_{\theta}\hat{\theta}^2 < \infty$  for all  $\theta \in \Theta$ . Then,  $\hat{\theta}$  is the best unbiased estimator for  $\theta$ .

**Proof** In order to make use of Lemma 1.3.18 it has to be shown that  $\operatorname{Cov}_{\theta}(\hat{\theta}, \varphi) = \operatorname{E}_{\theta}(\hat{\theta}\varphi) = 0, \theta \in \Theta$  for  $\operatorname{E}_{\theta}\varphi = 0, \theta \in \Theta$ . It holds that

$$E_{\theta}(\hat{\theta}\varphi) = E_{\theta}(E(\hat{\theta}\varphi|\hat{\theta})) \stackrel{\hat{\theta} \ \sigma(\hat{\theta})\text{-measurable}}{=} E_{\theta}(\hat{\theta} \cdot E_{\theta}(\varphi|\hat{\theta})) = E_{\theta}(\hat{\theta} \cdot g(\hat{\theta})) \stackrel{?}{=} 0,$$

for  $g(\hat{\theta}) = 0$  almost surely. Since  $\hat{\theta}$  is sufficient,  $g(t) = E_{\theta}(\varphi | \hat{\theta} = t)$  is independent of  $\theta$ .

Consider  $E_{\theta}g(\hat{\theta})$ . In order to show that  $g(\hat{\theta}) = 0$  for all  $\theta \in \Theta$ , it has to be shown that  $E_{\theta}g(\hat{\theta}) = 0$   $\theta \in \Theta$  since  $\hat{\theta}$  is already assumed to be complete.

$$E_{\theta}g(\hat{\theta}) = E_{\theta}(E_{\theta}(\varphi|\hat{\theta})) = E_{\theta}\varphi = 0$$

is assumed to hold, thus  $E_{\theta}(\varphi \hat{\theta}) = 0$  and  $\hat{\theta}$  is uncorrelated to  $\varphi : E_{\theta} \varphi = 0$ ,  $\theta \in \Theta$ , which implies that  $\hat{\theta}$  is the best unbiased estimator by Lemma 1.3.18.

**Theorem 1.3.20.** Let  $\hat{\theta}$  be an unbiased estimator for  $\theta$  and  $E_{\theta}\hat{\theta}^2 < \infty, \theta \in \Theta$ . Let  $\tilde{\theta}$  be a complete and sufficient estimator for  $\theta$ . Then, the estimator  $\theta^* = E(\hat{\theta} \mid \tilde{\theta})$  is the best unbiased estimator for  $\theta$ .

## Proof

1. It has to be shown that  $E_{\theta}\theta^{*2} < \infty \forall \theta \in \Theta$ . It holds that

$$E_{\theta}\left(\theta^{*2}\right) = E_{\theta}\left(E\left(\hat{\theta} \mid \tilde{\theta}\right)\right)^{2} \leq E_{\theta}\left(E\left(\hat{\theta}^{2} \mid \tilde{\theta}\right)\right) = E_{\theta}\hat{\theta}^{2} < \infty,$$

by Jensen's inequality for the conditional expectation, which states

$$f(\mathrm{E}(X \mid \mathcal{B})) \stackrel{\mathrm{f.s.}}{\leq} \mathrm{E}(f(X) \mid \mathcal{B})$$

for any random variable X,  $\sigma$ -algebra  $\mathcal{B}$  and convex function f.

- 2. It has to be shown that  $\theta^*$  is unbiased:  $E_{\theta}\theta^* = E_{\theta}(E(\hat{\theta} \mid \tilde{\theta})) = E_{\theta}\hat{\theta} = \theta, \ \theta \in \Theta$ , since  $\hat{\theta}$  is unbiased.
- 3. By Lemma 1.3.18, it is sufficient to show that  $E_{\theta}(\theta^*\varphi) = 0$  for  $\theta \in \Theta$ , if  $E_{\theta}\varphi = 0$ ,  $\theta \in \Theta$ .

$$\begin{split} E_{\theta}(\theta^*\varphi) &= E_{\theta}\big(\underbrace{E(\hat{\theta}\,|\,\tilde{\theta})}_{=g(\tilde{\theta}),\,\tilde{\theta}\,\,\mathrm{suf.}} \varphi\big) = E_{\theta}\big(g(\tilde{\theta})\varphi\big) = E_{\theta}\big(E\big(g(\tilde{\theta})\varphi\,|\,\tilde{\theta})\big) \\ &= g(\tilde{\theta}),\,\tilde{\theta}\,\,\mathrm{suf.} \\ &= E_{\theta}\big(g(\tilde{\theta})\cdot\underbrace{E(\varphi\,|\,\tilde{\theta})}_{=g_1(\tilde{\theta})}\big) = 0\,, \end{split}$$

if  $g_1(\tilde{\theta}) \stackrel{\text{a.s.}}{=} 0$ ,  $\theta \in \Theta$ . It needs to be shown that  $E_{\theta}g_1(\tilde{\theta}) = 0$ . Now,  $E_{\theta}g_1(\tilde{\theta}) = E_{\theta}(E(\varphi | \tilde{\theta})) = E_{\theta}\varphi = 0$  and the completeness of  $\tilde{\theta}$  imply (similarly to the proof of Theorem 1.3.19) that  $g_1(\tilde{\theta}) = 0$  almost surely.

**Lemma 1.3.21** (Blackwell-Rao inequality). Let  $\hat{\theta}$  be an unbiased estimator for  $\theta$  and  $E_{\theta}\hat{\theta}^2 < \infty$ ,  $\theta \in \Theta$ . Furthermore, let  $\tilde{\theta}$  be a sufficient estimator for  $\theta$ . Then, the unbiased estimator  $\theta^* := E_{\theta}(\hat{\theta} \mid \tilde{\theta})$  attains a variance which is smaller or equal to  $Var_{\theta}\hat{\theta}$ .

**Proof** See proof of Theorem 1.3.20. Here,  $\theta^*$  is unbiased, due to 2) in Theorem 1.3.20 and  $\operatorname{Var}_{\theta}\theta^* = \operatorname{E}_{\theta}\theta^{*2} - \theta^2 \leq \operatorname{E}_{\theta}\hat{\theta}^2 - \theta^2 = \operatorname{Var}_{\theta}\hat{\theta}$  due to 1) in Theorem 1.3.20.

**Remark 1.3.22.** The sufficiency of  $\tilde{\theta}$  is not mentioned explicitly in the proof of Lemma 1.3.21. It is still necessary in order to assure that  $\theta^* = E_{\theta}(\hat{\theta} \mid \tilde{\theta}) = g(\tilde{\theta})$  does not depend on  $\theta$ .

**Corollary 1.3.23.** If  $\hat{\theta}$  is a complete and sufficient estimator for  $\theta$  and there exists a function  $g: \mathbb{R} \to \mathbb{R}$  such that  $E_{\theta}g(\hat{\theta}) = \theta$ ,  $\forall \theta \in \Theta$ , then  $g(\hat{\theta})$  is the best unbiased estimator for  $\theta$ .

**Proof**  $g(\hat{\theta}) = E(g(\hat{\theta}) | \hat{\theta})$ , which is the best unbiased estimator by Theorem 1.3.20.

## 1.3.5 $\delta$ -Method

Let  $\hat{\theta}_n = \hat{\theta}_n(X_1, \dots, X_n)$  be an estimator of a parameter  $\theta \in \Theta \subseteq \mathbb{R}$  (m = 1), where  $(X_1, \dots, X_n)$  is a random sample of i.i.d. random variables  $X_j$  for  $j = 1, \dots, n$ ,  $X_j \sim F_{\theta}$ . Suppose that  $\hat{\theta}_n$  is asymptotically normal distributed, i.e. there exists a sequence of functions  $\{\sigma_n(\theta)\}_{n \in \mathbb{N}}$  with  $\sigma_n(\theta) > 0$  and  $\sigma_n(\theta) \xrightarrow[n \to \infty]{} 0$ ,  $\forall n \in \mathbb{N}$ ,  $\theta \in \Theta$  such that

$$\frac{\hat{\theta}_n - \theta}{\sigma_n(\theta)} \xrightarrow{d} Y \sim N(0, 1).$$

Let  $g: \Theta \to \mathbb{R}$  be a Borel measurable function. What can be said about the asymptotic normality of  $g(\hat{\theta}_n)$ ? In other words, this section aims to identify the sufficient conditions under which

$$\frac{g(\hat{\theta}_n) - g(\theta)}{\tilde{\sigma}_n(\theta)} \xrightarrow{d} Y \tag{1.12}$$

for another sequence  $\{\tilde{\sigma}_n(\theta)\}_{n\in\mathbb{N}}$  with  $\tilde{\sigma}_n(\theta) > 0$ ,  $n \in \mathbb{N}$ , and  $\tilde{\sigma}_n(\theta) \xrightarrow[n \to \infty]{} 0$ ,  $\theta \in \Theta$ . For linear  $g(\theta) = a \cdot \theta + b$ ,  $a, b \in \mathbb{R}$ , relation (1.12) obviously holds. When does (1.12) hold for more general functions g? There may be multiple reasons for the consideration of functions  $g(\hat{\theta}_n)$ . One of those lies in the *variance stabilization* which will be discussed at the end of this section. There a function g is considered, such that  $\tilde{\sigma}_n(\theta)$  does not depend on  $\theta$ . This makes the construction of asymptotic confidence regions for  $\theta$  much easier (cf. Section 2.2.3 for examples).

The following method of proving the asymptotic normality for  $g(\hat{\theta}_n)$  makes use of the Taylor series decomposition of a sufficiently smooth function g. It has been known since the early  $19^{\text{th}}$  century and first asymptotically described by J. Doob [11]. The name " $\delta$  method" alludes to the differential or increment  $dg(x) = g(x + \delta x) - g(x)$  which lies in the core of the method. Due to its very general nature, the results can be formulated for any asymptotically normal sequence of random variables  $\{Y_n\}_{n\in\mathbb{N}}$ , i.e. sequences with  $\frac{Y_n-\mu}{\sigma_n} \stackrel{d}{\longrightarrow} Y \sim N(0,1)$  for some  $\mu \in \mathbb{R}$  and a normalizing sequence  $\{\sigma_n\}_{n\in\mathbb{N}}$  with  $\sigma_n > 0$  for all  $n \in \mathbb{N}$  and  $\sigma_n \stackrel{\longrightarrow}{\longrightarrow} 0$ .

#### **Theorem 1.3.24.** Suppose that

$$\frac{Y_n - \mu}{\sigma_n} \xrightarrow{d} Y \sim N(0, 1) \tag{1.13}$$

for a sequence  $\{Y_n\}_{n\in\mathbb{N}}$ ,  $\mu$  and  $\{\sigma_n\}_{n\in\mathbb{N}}$  as above. Let  $g:\mathbb{R}\to\mathbb{R}$  be differentiable at  $x=\mu$  with  $g'(\mu)\neq 0$ . Then,

$$\frac{g(Y_n) - g(\mu)}{g'(\mu)\sigma_n} \xrightarrow[n \to \infty]{d} Y.$$

**Proof** First, show that (1.13) implies with  $\sigma_n \xrightarrow[n \to \infty]{} 0$  that

$$Y_n \xrightarrow[n \to \infty]{P} \mu. \tag{1.14}$$

Indeed, Slutsky's Theorem (cf. [32, Theorem 3.4.3.]) yields

$$\frac{Y_n - \mu}{\sigma_n} \xrightarrow[n \to \infty]{d} Y, \ \sigma_n \xrightarrow{a.s.} 0 \implies Y_n - \mu = \sigma_n \cdot \frac{Y - \mu}{\sigma_n} \xrightarrow[n \to \infty]{d} 0 \cdot Y = 0$$

$$\implies Y_n - \mu \xrightarrow[n \to \infty]{P} 0$$

by [32, Theorem 3.3.4.]. Introduce the function

$$h(x) = \begin{cases} \frac{g(x) - g(\mu)}{x - \mu} - g'(\mu), & x \neq \mu, \\ 0, & x = \mu. \end{cases}$$

Since g(x) is differentiable at  $x = \mu$ , h(x) is continuous at  $x = \mu$ . The Continuous Mapping Theorem (cf [32, Theorem 3.4.4.] implies

$$h(Y_n) \xrightarrow[n \to \infty]{P} h(\mu) = 0,$$

i.e.

$$\frac{g(Y_n) - g(\mu)}{Y_n - \mu} - g'(\mu) \underset{n \to \infty}{\overset{P}{\longrightarrow}} 0.$$

Multiplying both sides by  $\frac{Y_n - \mu}{\sigma_n}$  and using (1.13) in combination with Slutsky's Theorem implies that

$$\frac{h(Y_n)(Y_n - \mu)}{\sigma_n} = \frac{g(Y_n) - g(\mu)}{\sigma_n} - g'(\mu) \underbrace{\frac{Y_n - \mu}{\sigma_n}}_{n \to \infty} \xrightarrow[n \to \infty]{\overset{d}{\longrightarrow}} Y \sim N(0.1)$$

Hence  $\frac{g(Y_n)-g(\mu)}{\sigma_n} \xrightarrow[n\to\infty]{d} g'(\mu) \cdot Y$  as well and dividing by  $g'(\mu)$  yields the desired result.

**Remark 1.3.25.** If  $g \in C^1(B_\delta(\mu))$  for some  $\delta > 0$ , where

$$B_{\delta}(\mu) = \{x \in \mathbb{R} : |x - \mu| < \delta\},$$

the proof above can be simplified by using the Mean Value Theorem

$$g(Y_n) = g(\mu) + g'(\xi)(Y_n - \mu),$$

where  $\xi$  lies between  $\mu$  and  $Y_n$ . In addition, the Continuous Mapping Theorem together with (1.14) and the assumption  $g \in \mathcal{C}^1(B_\delta(\mu))$  yield  $g'(Y_n) \xrightarrow[n \to \infty]{P} g'(\mu)$ . By Slutsky's Theorem a modified version of (1.13) holds:

$$\frac{g(Y_n) - g(\mu)}{g'(Y_n)\sigma_n} \xrightarrow[n \to \infty]{d} Y.$$

**Example 1.3.26.** By [33, Theorem 7.4.4, 2)], the sample variance

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

is an asymptotically normally distributed estimator of  $\sigma^2 = \text{Var} X_j > 0$ :

$$\sqrt{n} \frac{S_n^2 - \sigma^2}{\sqrt{\mu_4' - \sigma_n^4}} \xrightarrow[n \to \infty]{d} Y \sim N(0, 1),$$

where  $\mu'_4 = \mathrm{E}(X_j - \mathrm{E}X_j)^4$ . One can show that the empirical standard deviation  $S_n$  is an asymptotically normal estimate of  $\sigma$ . Here,

$$g(x) = \sqrt{x},$$

$$g'(x) = \frac{1}{2\sqrt{x}},$$

$$\theta = \sigma^2 > 0,$$

$$g'(\sigma^2) = \frac{1}{2\sigma} > 0 \text{ and}$$

$$\sigma_n = \sqrt{\frac{\mu'_n - \sigma^4}{n}}$$

Following Theorem 1.3.24 it holds that

$$2\sigma\sqrt{n}\frac{S_n-\sigma}{\sqrt{\mu_4'-\sigma^4}}\xrightarrow[n\to\infty]{d}Y.$$

What happens if  $g'(\mu) = 0$  in Theorem 1.3.24? In this case, a higher order Taylor approximation should be used, as the following result shows.

**Theorem 1.3.27.** Assume that a sequence of random variables  $\{Y_n\}_{n \in \mathbb{N}}$  satisfies the conditions of Theorem 1.3.24. Let  $g : \mathbb{R} \to \mathbb{R}$  be  $m \geq 2$  times differentiable at  $\mu$  with  $g^{(j)}(\mu) = 0$ , j < m and  $g^{(m)}(\mu) \neq 0$ . Then,

$$m! \cdot \frac{g(Y_n) - g(\mu)}{g^{(m)}(\mu)\sigma_n^m} \xrightarrow[n \to \infty]{d} Y^m,$$

where  $Y \sim N(0, 1)$ .

**Proof** Use the function

$$h(x) = \begin{cases} m! \frac{g(x) - g(\mu)}{(x - \mu)^m} - g^{(m)}(\mu), & x \neq \mu, \\ 0, & x = \mu. \end{cases}$$

in the proof of Theorem 1.3.24.

**Example 1.3.28.** Suppose that  $\{Y_n\}_{n\in\mathbb{N}}$  is a sequence of random variables with

$$\frac{Y_n}{\sigma_n} \xrightarrow[n \to \infty]{d} Y \sim N(0,1)$$

for  $\sigma_n \xrightarrow[n \to \infty]{} 0$  with  $\sigma_n > 0$  for all  $n \in \mathbb{N}$ . Apply Theorem 1.3.27 to  $g(x) = \log^2(1+x)$  and  $m=2, \mu=0$ :

$$g'(x) = \frac{2\log(1+x)}{1+x}, \ g'(0) = 0,$$

$$g''(x) = \frac{\frac{2}{1+x}(1+x) - 2\log(1+x)}{(1+x)^2} = 2\frac{1 - \log(1+x)}{(1+x)^2} = 2\frac{\log\left(\frac{e}{1+x}\right)}{(1+x)^2},$$

$$g''(0) = 2 > 0.$$

Then,

$$2 \cdot \frac{\log^2(1+Y_n)}{2\sigma_n^2} = \frac{1}{\sigma_n^2} \log^2(1+Y_n) \xrightarrow{d} Y^2 \sim \chi_1^2.$$

As already mentioned above, it might be advantageous for some applications in the asymptotic theory of confidence intervals and statistical tests to eliminate the dependence of the asymptotic variance  $\sigma_n(\theta)$  from the parameter  $\theta$ . In other words, find a transformation g of the estimate  $\hat{\theta}$  such that  $\tilde{\sigma}_n(\theta)$  in (1.12) does not depend on  $\theta$  anymore. This device is known as variance stabilization. By Theorem 1.3.24 a function  $g: \mathbb{R} \to \mathbb{R}$  with  $g'(\theta) \neq 0$  such that  $g'(\theta) \cdot \sigma_n(\theta)$  depends only on  $n \in \mathbb{N}$  has to be found. Let  $\sigma_n(\theta) = \sigma(\theta) \cdot v_n$ , with  $v_n \to 0$ . Then it suffices to solve the ordinary differential equation

$$g'(\theta) = \frac{c}{\sigma(\theta)}, c \text{ constant.}$$
 (1.15)

If  $g'(\theta) = 0$ , Theorem 1.3.27 can be applied here accordingly.

### Example 1.3.29.

1. Consider a random sample  $(X_1, \ldots, X_n)$  of centered i.i.d. random variables with  $\mu_4 = \mathrm{E} X_j^4 < \infty$  and  $\sigma^2 = \mathrm{Var} X_j > 0$ . Since  $\mu = \mathrm{E} X_j = 0$ , consider the estimate  $\tilde{S}_n^2 = \frac{1}{n} \sum_{j=1}^n X_j^2$  of  $\sigma^2$ . Assume that  $\mu_4$  is known. By [33, Theorem 7.4.4, 2)], it holds that

$$\sqrt{n} \frac{\tilde{S}_n^2 - \sigma^2}{\sqrt{\mu_4 - \sigma^4}} \xrightarrow[n \to \infty]{d} Y \sim N(0, 1).$$

By stabilizing the asymptotic variance in this case,

$$g'(\sigma^2) = \frac{1}{\sqrt{\mu_4 - (\sigma^2)^2}}$$

has to be solved. The solution is given by  $g(x) = \arcsin\left(\frac{x}{\sqrt{\mu_4}}\right)$  and thus

$$\sqrt{n}\left(\arcsin\left(\frac{\tilde{S}_n^2}{\sqrt{\mu_4}}\right) - \arcsin\left(\frac{\sigma^2}{\sqrt{\mu_4}}\right)\right) \xrightarrow[n \to \infty]{d} Y \sim N(0, 1).$$

2. Let  $(X_1, \ldots, X_n)$  be a sample of i.i.d. random variables with  $X_j \sim \text{Bernoulli}(p)$ , for  $p \in (0, 1)$ . By [33, Theorem 7.3.2, a)], it holds that

$$\frac{\sqrt{n}(\hat{p}_n - p)}{\sqrt{p(1-p)} \xrightarrow[n \to \infty]{d}} Y \sim N(0,1),$$

where  $\hat{p}_n = \bar{X}_n$ . Similarly to 1), the variance stabilising transform g is given by  $g(p) = 2\arcsin(\sqrt{p})$ , since  $g'(p) = \frac{1}{\sqrt{p(1-p)}}$ . Applying Theorem 1.3.24 yields

$$2\sqrt{n}(\arcsin(\sqrt{\hat{p}_n}) - \arcsin(\sqrt{p})) \xrightarrow[n \to \infty]{d} Y \sim N(0, 1)$$
 (1.16)

3. Let  $(X_1, \ldots, X_n)$  be a sample of i.i.d. Poisson $(\lambda)$  distributed random variables with  $\lambda > 0$ . For  $\theta = \lambda$ ,  $\hat{\lambda}_n = \bar{X}_n$  it holds that

$$\sqrt{n} \frac{\hat{\lambda}_n - \lambda}{\sqrt{\lambda}} \xrightarrow[n \to \infty]{d} Y \sim N(0, 1).$$

The variance stabilizing transform g is then given by  $g(x) = 2\sqrt{x}$  because of  $g'(\lambda) = \frac{1}{\sqrt{\lambda}}$ . In summary, we get

$$2\sqrt{n}\left(\sqrt{\hat{\lambda}_n} - \sqrt{\lambda}\right) \xrightarrow[n \to \infty]{d} Y \sim N(0, 1). \tag{1.17}$$

**Remark 1.3.30.** The  $\delta$ -method can be extended to the asymptotic normality of (functions g of) d-dimensional random vectors  $\{Y_n\}_{n\in\mathbb{N}}$ , for  $d\geq 2$ . See [30, Section 3.3] for more details. It can be used to prove the asymptotic normality of the empirical Bravais-Pearson correlation coefficient

$$\rho_{XZ} = \frac{\sum_{j=1}^{n} X_{j} Z_{j} - n \bar{X}_{n} \bar{Z}_{n}}{S_{n,X} S_{n,Z}}$$

of i.i.d. random samples  $(X_1, \ldots, X_n)$  and  $(Z_1, \ldots, Z_n)$ , where  $S_{n,X}^2$  and  $S_{n,Z}^2$  are their sample variances. Similarly, the *empirical coefficient of variation*  $\frac{S_n}{X_n}$  of one i.i.d. sample  $(X_1, \ldots, X_n)$  can be shown to be asymptotically normal with

$$\sqrt{n} \left( \frac{S_n}{\bar{X}_n} - \frac{\sigma}{\mu} \right) \xrightarrow[n \to \infty]{d} Y \sim N \left( 0, \frac{\sigma_*^2 \mu^2}{4\sigma^2} \right),$$

where  $\bar{X}_n$  is the sample mean and  $S_n^2$  the sample variance, cf. [1]. Here

$$\sigma_*^2 = \frac{\mu_4}{\mu^4} - \left(\frac{\mu_2}{\mu^2}\right)^2 + 4\left(\frac{\mu_2}{\mu^2}\right)^3 - \frac{4\mu_2\mu_3}{\mu^5}$$

is a function of the first four moments of  $X_j$  which are assumed to be finite.

# Chapter 2

# Confidence Intervals

## 2.1 Introduction

This chapter will focus on the formal definition of confidence intervals. We will gain a deeper understanding of how they work and what they are used for. In particular, this chapter will cover *one-sample problems* and *two-sample problems*.

Recall the assumptions of parametric models: Let  $(X_1, ..., X_n)$  be a random sample with  $X_i \sim F_{\theta}$ , i = 1, ..., n, and  $F_{\theta} \in \{F_{\theta} : \theta \in \Theta\}$ , where  $\{F_{\theta} : \theta \in \Theta\}$  is some parametric family with  $\Theta \subset \mathbb{R}$ .

Each point estimator of  $\theta$  provides a value for the parameter vector. It would also be beneficial to have information about the accuracy of the estimator, i.e., a neighborhood which contains  $\theta$  with a certain probability  $1-\alpha$ . Here  $\alpha$  denotes a significance level, which indicates the probability of  $\theta$  being outside the predetermined neighborhood. Typical values are  $\alpha=0.01;0.05;0.1$ . For m=1 the neighborhood is an interval called confidence interval and the probability  $1-\alpha$  is called coverage probability or confidence level. It is always desired to achieve a high confidence level, e.g.,  $1-\alpha=0,99;0,95;0,9$  are typical values.

**Definition 2.1.1.** Let  $1 - \alpha$  be a confidence probability and  $\underline{\theta}: \mathbb{R}^n \to \overline{\mathbb{R}} = \mathbb{R} \cup \{\pm \infty\}, \ \overline{\theta}: \mathbb{R}^n \to \overline{\mathbb{R}}$  be two measurable sample functions with the property

$$\underline{\theta}(x_1, \dots x_n) \leq \overline{\theta}(x_1, \dots, x_n) \quad \forall (x_1, \dots x_n) \in \mathbb{R}^n.$$

If

1. 
$$P_{\theta}\left(\theta \in \left[\underline{\theta}(X_1, \dots, X_n), \overline{\theta}(X_1, \dots X_n)\right]\right) \geq 1 - \alpha, \quad \theta \in \Theta,$$

2. 
$$\inf_{\theta \in \Theta} P_{\theta} \left( \theta \in \left[ \underline{\theta}(X_1, \dots, X_n), \overline{\theta}(X_1, \dots, X_n) \right] \right) = 1 - \alpha \text{ and }$$

3. 
$$\lim_{n \to \infty} P_{\theta} \left( \theta \in \left[ \underline{\theta}(X_1, \dots, X_n), \overline{\theta}(X_1, \dots, X_n) \right] \right) = 1 - \alpha, \quad \theta \in \Theta,$$

then 
$$I = \left[\underline{\theta}(X_1, \dots, X_n), \overline{\theta}(X_1, \dots X_n)\right]$$
 is called

- 1. confidence interval,
- 2. minimal confidence interval,
- 3. asymptotic confidence interval,

with confidence level  $1-\alpha$ . Here,  $l_{\theta}(X_1, \ldots X_n) = \overline{\theta}(X_1, \ldots X_n) - \underline{\theta}(X_1, \ldots X_n)$  denotes the *length* of the confidence interval. It is desired to construct an interval, which has a relatively short length but a high confidence level, i.e.,  $1-\alpha = .99$ .

In Example 1.2.14, the construction of a confidence interval was introduced. This methodology can be generalized as follows.

- 1. Find a statistic  $T(X_1, \ldots, X_n, \theta)$  which
  - depends on  $\theta$  and
  - underlies a known (test) distribution F (possibly asymptotic as  $n \to \infty$ ).
- 2. Determine the quantiles  $F^{-1}(\alpha_1)$  and  $F^{-1}(1-\alpha_2)$  of the distribution F for the niveaus  $\alpha_1$  and  $1-\alpha_2$ , such that  $\alpha_1+\alpha_2=\alpha$ .
- 3. Solve (if possible) the inequality  $F^{-1}(\alpha_1) \leq T(X_1, \ldots, X_n, \theta) \leq F^{-1}(1 \alpha_2)$  w.r.t.  $\theta$ . The respective solution (if the statistic T in  $\theta$  is monotonically increasing)  $I = [T^{-1}(F^{-1}(\alpha_1)), T^{-1}(F^{-1}(1 \alpha_2))]$  is a confidence interval for  $\theta$  with confidence level  $1 \alpha$ , because

$$P_{\theta}(\theta \in I) = P_{\theta} \left( T_{\theta}^{-1}(F^{-1}(\alpha_{1})) \leq \theta \leq T^{-1}(F^{-1}(1 - \alpha_{2})) \right)$$

$$= P_{\theta} \left( F^{-1}(\alpha_{1}) \leq T_{\theta}(X_{1}, \dots, X_{n}, \theta) \leq F^{-1}(1 - \alpha_{2}) \right)$$

$$= F(F^{-1}(1 - \alpha_{2})) - F(F^{-1}(\alpha_{1}))$$

$$= 1 - \alpha_{2} - \alpha_{1}$$

$$= 1 - \alpha \text{ for all } \theta \in \Theta.$$

For asymptotic confidence intervals the notation  $\lim_{n\to\infty}$  is introduced:  $\lim_{n\to\infty} P_{\theta}(\theta \in I) = \ldots = 1-\alpha$ . Here  $T_{\theta}^{-1}$  denotes the inverse of  $T(X_1, \ldots, X_n, \theta)$  w.r.t.  $\theta$ . A corresponding picture can be found in Figure 2.1.

## Definition 2.1.2.

1. If  $\alpha_1 = \alpha_2 = \alpha/2$ , then the confidence interval given by

$$I = \left[T^{-1}\left(F^{-1}\left(\frac{\alpha}{2}\right)\right), T^{-1}\left(F^{-1}\left(1-\frac{\alpha}{2}\right)\right)\right]$$

is called *symmetric*.

2. If  $\alpha_1 = 0$ , i.e.,  $\underline{\theta}(X_1, \dots, X_n) = -\infty$ , then the confidence interval  $\left(-\infty, \overline{\theta}(X_1, \dots, X_n)\right]$  is called *one sided*.

Analogously, if  $\alpha_2 = 0$  i.e.,  $\overline{\theta}(X_1, \dots, X_n) = +\infty$ , then the confidence interval is given by  $[\underline{\theta}(X_1, \dots, X_n), +\infty)$ .

From now on, mostly symmetric confidence intervals will be constructed. More general, non symmetric confidence intervals can easily be constructed similarly.

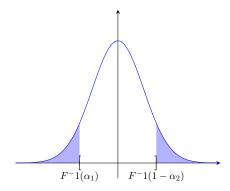


Figure 2.1: asymptotic confidence interval

**Remark 2.1.3.** One can observe, that the process of constructing a confidence interval is similar to constructing a test. In Definition 2.1.2,  $T(X_1, \ldots, X_n)$  is called *test statistic*. Generally, a statistical test for every confidence interval can be constructed, but not the other way around.

# 2.2 One-sample problems

This section provides examples of confidence intervals for parameters of known distributions using the algorithm above.

#### 2.2.1 Normal distribution

Let  $X_1, \ldots, X_n$  be i.i.d. random sample with  $X_i \sim N(\mu, \sigma^2)$ ,  $i = 1, \ldots, n$ .

## Confidence interval for the expectation $\mu$

• Known variance  $\sigma^2$ : Under the assumption that  $\sigma^2$  is known, [33, Theorem 7.3.2] implies that an exact confidence interval for  $\mu$  with confidence level  $1 - \alpha$  can be constructed. Since  $\overline{X}_n \sim N(\mu, \sigma^2/n)$ ,

$$T(X_1, \dots, X_n, \mu) = \sqrt{n} \frac{\overline{X}_n - \mu}{\sigma} \sim \mathcal{N}(0, 1)$$

Let  $z_{\alpha_1}$  and  $z_{1-\alpha_2}$  be quantiles of the  $\mathcal{N}(0,1)$  distribution, such that  $\alpha_1 + \alpha_2 = \alpha$ . Then,  $1 - \alpha$  is the given confidence level and

$$1 - \alpha = P\left(z_{\alpha_1} \le T(X_1, \dots, X_n, \mu) \le z_{1-\alpha_2}\right)$$

$$= P\left(z_{\alpha_1} \le \sqrt{n} \frac{\overline{X}_n - \mu}{\sigma} \le z_{1-\alpha_2}\right)$$

$$\stackrel{(-z_{\alpha_1} = z_{1-\alpha_1})}{=} P\left(\overline{X}_n - \frac{z_{1-\alpha_2}\sigma}{\sqrt{n}} \le \mu \le \overline{X}_n + \frac{z_{1-\alpha_1}\sigma}{\sqrt{n}}\right).$$

Hence,  $\left[\underline{\theta}(X_1,\ldots,X_n),\overline{\theta}(X_1,\ldots,X_n)\right]$  with

$$\underline{\theta}(X_1,\ldots,X_n) = \overline{X}_n - z_{1-\alpha_2} \frac{\sigma}{\sqrt{n}}$$

and

$$\overline{\theta}(X_1,\ldots,X_n) = \overline{X}_n + z_{1-\alpha_1} \frac{\sigma}{\sqrt{n}},$$

is a confidence interval for  $\mu$  with confidence level  $1 - \alpha$ . Its length is  $l_{\mu}(X_1, \ldots, X_n) = \frac{\sigma}{\sqrt{n}} (z_{1-\alpha_2} + z_{1-\alpha_1}).$ 

If  $n \to \infty$ , then  $l_{\mu}(X_1, \ldots, X_n) \to 0$  which means, that if the amount of available information increases, i.e.,  $n \to \infty$ , the precision of the estimation also increases.

If the underlying distribution is symmetric i.e.,  $\alpha_1 = \alpha_2 = \alpha/2$ , then

$$\underline{\theta}(X_1, \dots, X_n) = \overline{X}_n - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}},$$

$$\overline{\theta}(X_1, \dots, X_n) = \overline{X}_n + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}},$$

and

$$l_{\mu}(X_1,\ldots,X_n) = \frac{2\sigma}{\sqrt{n}} z_{1-\alpha/2}.$$

If the length  $\varepsilon>0$  is predetermined, the number of necessary observations n for achieving the desired precision can be calculated by solving

$$\frac{2\sigma}{\sqrt{n}}z_{1-\alpha/2} \le \varepsilon \tag{2.1}$$

for n, which yields

$$n \ge \left(\frac{2\sigma z_{1-\alpha/2}}{\varepsilon}\right)^2.$$

For  $\alpha_1 = 0$  or  $\alpha_2 = 0$  one sided intervals like  $\left(-\infty, \overline{X}_n + z_{1-\alpha} \frac{\sigma}{\sqrt{n}}\right]$ , resp.  $\left[\overline{X}_n - z_{1-\alpha} \frac{\sigma}{\sqrt{n}}, +\infty\right)$  can be constructed.

• Unknown variance  $\sigma^2$ : Using [33, Theorem 7.4.10.], the following confidence interval with confidence level  $1 - \alpha \in (0, 1)$  for the expectation  $\mu$  of a normally distributed random sample  $(X_1, \ldots, X_n)$  with unknown variance  $\sigma^2$  can be constructed.

$$P\left(\mu \in \left[\bar{X}_n - \frac{t_{n-1,1-\alpha/2}}{\sqrt{n}}S_n, \bar{X}_n + \frac{t_{n-1,1-\alpha/2}}{\sqrt{n}}S_n\right]\right) = 1 - \alpha,$$

since

$$P\left(\sqrt{n}\frac{\bar{X}_{n} - \mu}{S_{n}} \in \left[\underbrace{t_{n-1,\alpha/2}}_{\text{bc. of the sym. of } t \text{ dist.}}, t_{n-1,1-\alpha/2}\right]\right) = (2.2)$$

$$= -t_{n-1,1-\alpha/2} \text{ bc. of the sym. of } t \text{ dist.}$$

$$= F_{t_{n-1}}(t_{n-1,1-\alpha/2}) - F_{t_{n-1}}(t_{n-1,\alpha/2})$$

$$= 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha,$$

where  $t_{n-1,\alpha}$  is the  $\alpha$  quantile of the  $t_{n-1}$  distribution. By solving for  $\mu$  in (2.2) the remaining part can be shown.

Note that the length  $l_{\mu}(X_1, \dots X_n) = \frac{2S_n}{\sqrt{n}} t_{n-1,1-\alpha/2}$  of the confidence interval is a random variable. Thus, the expected length

$$\mathbb{E} l_{\mu}(X_1, \dots X_n) = \frac{2}{\sqrt{n}} \mathbb{E} S_n t_{n-1, 1-\alpha/2}$$

yields an answer to the question about the required number of observations n for a predetermined precision  $\varepsilon > 0$  (cf. Equation (2.1)).

## Confidence interval for the variance $\sigma^2$

• Known expectation  $\mu$ : Consider the estimator  $\tilde{S}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$  for  $\sigma^2$ . [33, Theorem 7.4.8 1.] implies  $\frac{n\tilde{S}_n^2}{\sigma^2} \sim \chi_n^2$ . Define  $T(X_1, \dots, X_n, \sigma^2) := \frac{n\tilde{S}_n^2}{\sigma^2}$ , then

$$P\left(\chi_{n,\alpha_2}^2 \le \frac{n\tilde{S_n}^2}{\sigma^2} \le \chi_{n,1-\alpha_1}^2\right) = P\left(\frac{n\tilde{S}_n^2}{\chi_{n,1-\alpha_1}^2} \le \sigma^2 \le \frac{n\tilde{S}_n^2}{\chi_{n,\alpha_2}^2}\right) = 1 - \alpha.$$

Thus,  $\left[\frac{n\tilde{S}_n^2}{\chi_{n,1-\alpha_1}^2}, \frac{n\tilde{S}_n^2}{\chi_{n,\alpha_2}^2}\right]$  is a confidence interval for  $\sigma^2$  with level  $1-\alpha$ , where  $\alpha = \alpha_1 + \alpha_2$ . The expected length is given by

$$\mathbb{E} l_{\sigma^2} = n\sigma^2 \left( \frac{1}{\chi_{n,\alpha_2}^2} - \frac{1}{\chi_{n,1-\alpha_1}^2} \right).$$

• Unknown expectation  $\mu$ : Similarly to the construction above, [33, Theorem 7.4.8. 1.] implies that  $\left[\frac{(n-1)S_n^2}{\chi_{n-1,1-\alpha_1}^2}, \frac{(n-1)S_n^2}{\chi_{n-1,\alpha_2}^2}\right]$  is a confidence interval for  $\sigma^2$  with confidence level  $1-\alpha$ , where  $\alpha=\alpha_1+\alpha_2$ . Note that  $\frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2$  for the sample variance  $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n \left(X_i - \overline{X}_n\right)^2$ . The expected length is

$$\mathbb{E} l_{\sigma^2} = (n-1)\sigma^2 \left( \frac{1}{\chi_{n-1,\alpha_2}^2} - \frac{1}{\chi_{n-1,1-\alpha_1}^2} \right).$$

## 2.2.2 Confidence intervals and stochastic inequalities

An alternative approach for obtaining confidence intervals is applying stochastic inequalities. Let, for example,  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables with  $\mathbb{E} X_i = \mu$ ,  $\operatorname{Var} X_i = \sigma^2 \in (0, \infty)$ , then the Tschebyschew inequality can be used to construct a simple confidence interval for  $\mu$ :

$$P(|\overline{X}_n - \mu| > \varepsilon) \le \frac{\operatorname{Var} \overline{X}_n}{\varepsilon^2} = \frac{\sigma^2}{n\varepsilon^2} = \alpha.$$

Then, for  $\varepsilon = \frac{\sigma}{\sqrt{n\alpha}}$ 

$$1 - \alpha \le P\left(|\overline{X}_n - \mu| \le \varepsilon\right)$$

$$= P\left(-\frac{\sigma}{\sqrt{n\alpha}} \le -\overline{X}_n + \mu \le \frac{\sigma}{\sqrt{n\alpha}}\right)$$

$$= P\left(\overline{X}_n - \frac{\sigma}{\sqrt{n\alpha}} \le \mu \le \overline{X}_n + \frac{\sigma}{\sqrt{n\alpha}}\right)$$

holds. The confidence interval  $\left[\overline{X}_n - \frac{\sigma}{\sqrt{n\alpha}}, \overline{X}_n + \frac{\sigma}{\sqrt{n\alpha}}\right]$  for  $\mu$  with known variance  $\sigma^2$  is independent of the underlying distribution of  $X_i$  since no assumptions have been made.

More precise confidence intervals can be constructed by using the *Hoeffding inequality*:

**Theorem 2.2.1** (Hoeffding inequality). Let  $Y_1, \ldots, Y_n$  be independent random variables with  $\mathbb{E} Y_i = 0, a_i \leq Y_i \leq b_i$  a.s.,  $i = 1, \ldots, n$ . For all  $\varepsilon > 0$ ,

$$P\left(\sum_{i=1}^{n} Y_i \ge \varepsilon\right) \le \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^{n} (b_i - a_i)^2}\right)$$

holds.
(without proof)

Assume that  $X_1, \ldots, X_n$  are i.i.d. with  $X_i \sim \text{Bernoulli}(p), p \in (0,1)$ . Next, we will show how to construct a confidence interval for p.

Corollary 2.2.2. Let  $X_1, \ldots, X_n$  be i.i.d. Bernoulli (p) random variables. Then

$$P(|\overline{X}_n - p| > \varepsilon) \le 2e^{-2n\varepsilon^2}, \quad \varepsilon > 0.$$

**Proof** 

$$\overline{X}_n - p = \frac{1}{n} \sum_{i=1}^n \underbrace{(X_i - p)}_{Y_i}, Y_i \in [-p, 1 - p],$$

holds, which means that  $a_i = -p$ ,  $b_i = 1 - p$ ,  $b_i - a_i = 1$ ,  $i = 1, \ldots, n$ ,  $\mathbb{E} Y_i = p - p = 0$ . Then,

$$P_{p}\left(\left|\overline{X}_{n}-p\right|>\varepsilon\right) = P_{p}\left(\left|\sum_{i=1}^{n}Y_{i}\right| \geq \varepsilon n\right)$$

$$= P_{p}\left(\sum_{i=1}^{n}Y_{i} \geq \varepsilon n\right) + P_{p}\left(\sum_{i=1}^{n}(-Y_{i}) \geq \varepsilon n\right)$$
(Theorem 2.2.1)
$$< 2e^{-\frac{2\varepsilon^{2}n^{2}}{n}} = 2e^{-2\varepsilon^{2}n}.$$

where Theorem 2.2.1 is applied to  $\{Y_i\}$  as well as  $\{-Y_i\}$ .

**Remark 2.2.3.** Let  $\alpha > 0$  and  $\varepsilon_n = \sqrt{\frac{1}{2n} \log \frac{2}{\alpha}}$ . Applying Corollary 2.2.2 with  $\varepsilon_n$  yields  $P_p\left(|\overline{X}_n - p| > \varepsilon_n\right) \le \alpha$ , and thus  $P_p\left(|\overline{X}_n - p| \le \varepsilon_n\right) \ge 1 - \alpha$ . Hence,

$$\left[\overline{X}_n - \sqrt{\frac{1}{2n}\log\frac{2}{\alpha}}, \, \overline{X}_n + \sqrt{\frac{1}{2n}\log\frac{2}{\alpha}}\right]$$

is a confidence interval for p with level  $1 - \alpha$ .

## 2.2.3 Asymptotic confidence intervals

The idea behind asymptotic confidence intervals is relatively simple, as it can be explained by using the example of an asymptotically normal distributed estimator  $\hat{\theta}$  for a parameter  $\theta$ . Let  $(X_1, \ldots, X_n)$  be an i.i.d. random sample with  $X_i \sim F_{\theta}$ ,  $\theta \in \Theta \subseteq \mathbb{R}$ . Let  $\hat{\theta}_n = \hat{\theta}(X_1, \ldots, X_n)$  be an estimator for  $\theta$ , that is asymptotically normal distributed. If  $\hat{\theta}_n$  is unbiased for every  $n \in \mathbb{N}$ , then

$$\frac{\hat{\theta}_n - \theta}{\hat{\sigma}_n} \stackrel{d}{\longrightarrow} Y \sim \mathcal{N}(0, 1),$$

where  $\hat{\sigma}_n$  is a consistent estimator for the asymptotic variance of  $\hat{\theta}_n$ . Furthermore,

$$\lim_{n \to \infty} P_{\theta} \left( z_{\alpha/2} \le \frac{\hat{\theta}_n - \theta}{\hat{\sigma}_n} \le z_{1-\alpha/2} \right)$$

$$= \lim_{n \to \infty} P_{\theta} \left( \theta \in \left[ \hat{\theta}_n - z_{1-\alpha/2} \hat{\sigma}_n, \, \hat{\theta}_n + z_{1-\alpha/2} \hat{\sigma}_n \right] \right) = 1 - \alpha.$$

Thus,

$$\left[\hat{\theta}_n - z_{1-\alpha/2}\hat{\sigma}_n, \, \hat{\theta}_n + z_{1-\alpha/2}\hat{\sigma}_n\right]$$

is an asymptotic confidence interval for  $\theta$  with level  $1-\alpha$ . This approach can be applied to the following two examples:

## • Bernoulli distribution

Let  $X_i \sim \text{Bernoulli}(p)$ , i = 1, ..., n. Then  $\theta = p$  and  $\hat{\theta}_n = \hat{p}_n = \overline{X}_n$ . Moreover,  $\mathbb{E}_p \hat{p}_n = p$ ,  $\text{Var}_p \hat{p}_n = \frac{p(1-p)}{n}$ . Let  $\hat{\sigma}^2 = \frac{1}{n} \hat{p}_n (1 - \hat{p}_n) = \frac{\overline{X}_n}{n} (1 - \overline{X}_n)$  be the Plug-In estimator for  $\sigma^2$ . Then the central limit theorem [33, Theorem 5.2.2.] and Slutsky's theorem [32, Theorem 3.4.1] imply

$$\sqrt{n} \frac{\overline{X}_n - p}{\sqrt{\overline{X}_n (1 - \overline{X}_n)}} \xrightarrow{n \to \infty} Y \sim \mathcal{N}(0, 1).$$

Thus,

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

is an asymptotic confidence interval for p with confidence level  $1-\alpha$ . Since  $p \in [0, 1]$  is supposed to hold, one considers

$$\underline{p}(X_1, \dots, X_n) = \max \left\{ 0, \, \overline{X}_n - z_{1-\alpha/2} \sqrt{\frac{\overline{X}_n(1 - \overline{X}_n)}{n}} \right\}$$

and

$$\overline{p}(X_1,\ldots,X_n) = \min \left\{ 1, \ \overline{X}_n + z_{1-\alpha/2} \sqrt{\frac{\overline{X}_n(1-\overline{X}_n)}{n}} \right\}.$$

## Remark 2.2.4.

1. Another confidence interval for the parameter p of the Bernoulli distribution can be obtained by considering an application of the central limit theorem

$$\lim_{n \to \infty} P_p \left( -z_{1-\alpha/2} \le \sqrt{n} \frac{\overline{X}_n - p}{\sqrt{p(1-p)}} \le z_{1-\alpha/2} \right) = 1 - \alpha$$

and solving the quadratic inequality for p.

Exercise 2.2.5 Solve the inequality!

2. Using the variance stabilization from Example 1.3.29, 2., the relation (1.16) can be used to construct a confidence interval for p with sufficiently large n.

$$P\left(-z_{1-\frac{\alpha}{2}} \le 2\sqrt{n}(\arcsin\sqrt{\bar{X}_n} - \arcsin\sqrt{p}) \le z_{1-\frac{\alpha}{2}}\right) \approx 1 - \alpha$$

holds, hence

$$\arcsin\sqrt{\bar{X}_n} - \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}} \le \arcsin\sqrt{p} \le \arcsin\sqrt{\bar{X}_n} + \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}}.$$

With probability  $1 - \alpha$ 

$$\begin{split} \sqrt{p} &\in \left[ \sin \left( \arcsin \sqrt{\bar{X}_n} - \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}} \right), \sin \left( \arcsin \sqrt{\bar{X}_n} + \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}} \right) \right] \Rightarrow \\ p &\in \left[ \sin^2 \left( \arcsin \sqrt{\bar{X}_n} - \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}} \right), \sin^2 \left( \arcsin \sqrt{\bar{X}_n} + \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}} \right) \right] \end{split}$$

holds. As 
$$\sqrt{\bar{X}_n} \xrightarrow[n \to \infty]{a.s.} \sqrt{p} \in (0,1), \xrightarrow{z_{1-\frac{\alpha}{2}}} \xrightarrow[n \to \infty]{} 0$$
 and since

$$\sin\left(\arcsin\sqrt{\bar{X}_n} \pm \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}}\right) > 0$$

for sufficiently large n, the terms

$$\max\left\{0,\arcsin\sqrt{\bar{X}_n} - \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}}\right\}$$

and

$$\min\left\{\frac{\pi}{2}, \arcsin\sqrt{\bar{X}_n} + \frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}}\right\}$$

do not need to be considered here.

## • Poisson distribution:

Let  $X_i \sim \text{Poisson}(\lambda)$ , i = 1, ..., n, then  $\theta = \lambda$ ,  $\hat{\theta}_n = \hat{\lambda} = \overline{X}_n$ . Since  $\mathbb{E}_{\lambda} X_i = \text{Var}_{\lambda} X_i = \lambda$ , the central limit theorem [33, Theorem 5.2.2.] can be applied:

$$\sqrt{n} \frac{\overline{X}_n - \lambda}{\sqrt{\lambda}} \xrightarrow[n \to \infty]{d} Y \sim \mathcal{N}(0, 1)$$

Since  $\overline{X}_n$  is strongly consistent for  $\lambda$ , Slutsky's theorem [32, Theorem 3.4.1] implies

$$\sqrt{n} \frac{\overline{X}_n - \lambda}{\sqrt{\overline{X}_n}} \xrightarrow[n \to \infty]{d} Y \sim \mathcal{N}(0, 1).$$

Thus, a asymptotic confidence interval

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

for the parameter  $\lambda$  with level  $1 - \alpha$  can be obtained.

## Remark 2.2.6.

1. Similarly to Remark 2.2.4, solving the quadratic inequality

$$\lim_{n \to \infty} P_{\lambda} \left( \sqrt{n} \frac{\overline{X}_n - \lambda}{\sqrt{\lambda}} \in [-z_{1-\alpha/2}, z_{1-\alpha/2}] \right) = 1 - \alpha$$

for  $\lambda$  leads to an alternative asymptotic confidence interval for  $\lambda$ . **Exercise 2.2.7.** Solve this quadratic inequality.

2. Since  $\lambda > 0$ , the lower bound can be adjusted to

$$\underline{\lambda}(X_1, \dots, X_n) = \max \left\{ 0, \, \overline{X}_n - z_{1-\alpha/2} \sqrt{\frac{\overline{X}_n}{n}} \right\}$$

3. Using the variance stabilization transformation from Example 1.3.29, 3.

$$P\left(-z_{1-\frac{\alpha}{2}} \le \sqrt{n}(\sqrt{\bar{x}_n} - \sqrt{\lambda}) \le z_{1-\frac{\alpha}{2}}\right) \underset{n \to \infty}{\longrightarrow} 1 - \alpha$$

holds. For n sufficiently large, the asymptotic confidence interval for  $\lambda$  with confidence level  $1 - \alpha$  is given by

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

Since  $\sqrt{\bar{X}_n} \xrightarrow[n \to \infty]{a.s.} \sqrt{\lambda}$  by the strong law of large numbers and

$$\frac{z_{1-\frac{\alpha}{2}}}{\sqrt{n}} \underset{n \to \infty}{\longrightarrow} 0,$$

for sufficiently large n

$$\sqrt{\bar{X}_n} - \frac{z_{1-\alpha}2}{2\sqrt{n}} > 0$$

holds. Hence  $\max\left\{0,\sqrt{\bar{X}_n}-\frac{z_{1-\frac{\alpha}{2}}}{2\sqrt{n}}\right\}$  does not need to be taken care of.

# 2.3 Two-sample problems

In this section, some characteristics or parameters of two different samples will be compared by constructing confidence intervals for simple functions of those parameters.

Consider two random samples  $Y_1 = (X_{11}, \ldots, X_{1n_1})$  and  $Y_2 = (X_{21}, \ldots, X_{2n_2})$  of random variables  $X_{i1}, \ldots, X_{in_i}$ , i = 1, 2, which are, within the sample  $Y_i$  i.i.d. with  $X_{ij} \stackrel{d}{=} X_i$ ,  $j = 1, \ldots n_i$ , i = 1, 2. Assume for the prototype random variable  $X_i \sim F_{\theta_i}$ ,  $\theta_i \in \Theta \subset \mathbb{R}^m$ . In general it will not be assumed that  $Y_1$  and  $Y_2$  are independent. If they are dependent, the random samples  $Y_1$  and  $Y_2$  are called related samples. Consider a function  $g : \mathbb{R}^{2m} \to \mathbb{R}$  of the parameter vectors  $\theta_1$  and  $\theta_2$ . In this lecture the cases m = 1, 2,  $g(\theta_1, \theta_2) = \theta_{1j} - \theta_{2j}$  and  $g(\theta_1, \theta_2) = \frac{\theta_{1j}}{\theta_{2j}}$  will mostly be covered, where  $\theta_i = (\theta_{i1}, \ldots, \theta_{im})$ , i = 1, 2. The goal is to construct a (possibly asymptotic) confidence interval for  $g(\theta_1, \theta_2)$  by using  $(Y_1, Y_2)$ .

As it turns out, the approach will be similar to Section 2.2. A statistic  $T(Y_1, Y_2, g(\theta_1, \theta_2))$  is desired, that has a (possibly asymptotic) test distribution F and explicitly depends on  $g(\theta_1, \theta_2)$ .

By solving the inequality  $F_{\alpha_1}^{-1} \leq T(Y_1, Y_2, g(\theta_1, \theta_2)) \leq F_{1-\alpha_2}^{-1}$  for  $g(\theta_1, \theta_2)$  a (possibly asymptotic) confidence interval with level  $1 - \alpha$ ,  $\alpha = \alpha_1 + \alpha_2$  can be obtained.

#### 2.3.1 Normally distributed samples

Assume, that  $X_i \sim N(\mu_i, \sigma_i^2)$ , i = 1, 2.

# Confidence interval for the difference $\mu_1 - \mu_2$ with known variance $\sigma_1^2$ and $\sigma_2^2$ and independent random samples

Let  $Y_1$  and  $Y_2$  be independent and  $\sigma_1^2, \sigma_2^2$  known. Consider the parameter function  $g(\mu_1, \mu_2) = \mu_1 - \mu_2$  and

$$\overline{X}_{in_i} = \frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij}, i = 1, 2$$

the sample mean of  $Y_1$  and  $Y_2$ . Then,  $\overline{X}_{in_i} \sim N(\mu_i, \frac{\sigma_i^2}{n_i})$ , i = 1, 2. [33, Theorem 7.3.2, 4] implies that  $\overline{X}_{1n_1}$  and  $\overline{X}_{2n_2}$  are independent. The stability of the normal distribution implies

$$\overline{X}_{1n_1} - \overline{X}_{2n_2} \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$$

and normalizing yields

$$T(Y_1, Y_2, \mu_1 - \mu_2) = \frac{\overline{X}_{1n_1} - \overline{X}_{2n_2} - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim \mathcal{N}(0, 1).$$

The confidence interval

$$\left[\overline{X}_{1n_1} - \overline{X}_{2n_2} - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \, \overline{X}_{1n_1} - \overline{X}_{2n_2} + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right]$$

for  $\mu_1 - \mu_2$  with level  $1 - \alpha$  then results.

# Confidence interval for the quotient $\sigma_1^2/\sigma_2^2$ with unknown expected values $\mu_1$ and $\mu_2$ and independent random samples

Let  $Y_1$  and  $Y_2$  be independent and  $g(\sigma_1, \sigma_2) = \frac{\sigma_1^2}{\sigma_2^2}$ . Construct a statistic  $T(Y_1, Y_2, \frac{\sigma_1^2}{\sigma_2^2})$ . Let

$$S_{in_i}^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} \left( X_{ij} - \overline{X}_{in_i} \right)^2, i = 1, 2$$

be the sample variances of  $Y_1$  and  $Y_2$ . Then, applying [33, Theorem 7.4.8.] yields  $\frac{(n_i-1)S_{in_i}^2}{\sigma_i^2} \sim \chi_{n_i-1}^2$ , i=1,2. Since  $S_{in_i}^2$ , i=1,2 are independent, the definition of the F distribution implies

$$T\left(Y_1, Y_2, \frac{\sigma_1^2}{\sigma_2^2}\right) = \frac{\frac{(n_2 - 1)S_{2n_2}^2}{(n_2 - 1)\sigma_2^2}}{\frac{(n_1 - 1)S_{1n_1}^2}{(n_1 - 1)\sigma_2^2}} = \frac{S_{2n_2}^2}{S_{1n_1}^2} \cdot \frac{\sigma_1^2}{\sigma_2^2} \sim F_{n_2 - 1, n_1 - 1}.$$

Thus, the confidence interval

$$\left[\frac{S_{1n_1}^2}{S_{2n_2}^2}F_{n_2-1,\,n_1-1,\,\alpha_1},\,\frac{S_{1n_1}^2}{S_{2n_2}^2}F_{n_2-1,\,n_1-1,\,1-\alpha_2}\right]$$

for  $\frac{\sigma_1^2}{\sigma_2^2}$  with level  $1 - \alpha$  is obtained.

# Confidence interval for the difference $\mu_1 - \mu_2$ of expected values with dependent samples

Let  $Y_1$  and  $Y_2$  be linked, i.e.,  $X_1 - X_2 \sim N(\mu_1 - \mu_2, \sigma^2)$  for an unknown  $\sigma^2 > 0$ ,  $n_1 = n_2 = n$ . Since  $X_{ij}, j = 1, ..., n$  are i.i.d.,  $Z_j = X_{1j} - X_{2j} \sim N(\mu_1 - \mu_2, \sigma^2), j = 1, ..., n$  holds.

The goal is to construct a confidence interval for  $\mu_1 - \mu_2$ . Consider the random samples  $(Z_1, \ldots, Z_n)$  and the results of Section 2.2, then the confidence interval

$$\left[\overline{Z}_n - t_{n-1,1-\frac{\alpha}{2}} \frac{S_n}{\sqrt{n}}, \, \overline{Z}_n + t_{n-1,1-\frac{\alpha}{2}} \frac{S_n}{\sqrt{n}}\right]$$

for  $\mu_1 - \mu_2$  with level  $1 - \frac{\alpha}{2}$  is obtained. Here,

$$\overline{Z}_n = \frac{1}{n} \sum_{j=1}^n Z_j = \frac{1}{n} \sum_{j=1}^n (X_{1j} - X_{2j}) = \overline{X}_{1n} - \overline{X}_{2n}$$

and

$$S_n^2 = \frac{1}{n-1} \sum_{j=1}^n \left( Z_j - \overline{Z}_n \right)^2 = \frac{1}{n-1} \sum_{j=1}^n \left( X_{1j} - X_{2j} - \overline{X}_{1n} + \overline{X}_{2n} \right)^2.$$

### 2.3.2 Poisson distributed random samples

Assume that the random samples  $Y_1$  and  $Y_2$  are independent and  $X_i \sim \text{Poisson}(\lambda_i)$ , i = 1, 2. The goals is to construct confidence intervals for

$$g(\lambda_1, \lambda_2) = \lambda_1 - \lambda_2,$$
  

$$g(\lambda_1, \lambda_2) = \frac{n_2 \lambda_2}{n_1 \lambda_1 + n_2 \lambda_2} = \frac{\lambda_2}{\rho \lambda_1 + \lambda_2},$$

where  $\rho = \frac{n_1}{n_2} = \text{const for } n_1, n_2 \to \infty$ .

# Asymptotic confidence interval for $\lambda_1 - \lambda_2$

In order to obtain an asymptotically  $\mathcal{N}(0,1)$  distributed statistic  $T(Y_1, Y_2, \lambda_1 - \lambda_2)$ , the central limit theorem of Ljapunow (cf. [33, Theorem 4.2.13]) will be used.

**Lemma 2.3.1.** For  $n_1, n_2 \to \infty$  with  $0 < c_1 \le n_1/n_2 \le c_2 < \infty$ ,

$$\frac{\overline{X}_{1n_1} - \overline{X}_{2n_2} - \lambda_1 + \lambda_2}{\sqrt{\frac{\lambda_1}{n_1} + \frac{\lambda_2}{n_2}}} \xrightarrow[n_1, n_2 \to \infty]{d} Y \sim \mathcal{N}(0, 1)$$

holds.

**Proof** Define the random variable

$$Z_{nk} = \begin{cases} \frac{X_{1k} - \lambda_1}{n_1 \sqrt{\frac{\lambda_1}{n_1} + \frac{\lambda_2}{n_2}}}, & k = 1, \dots, n_1\\ -\frac{X_{2k - n_1} - \lambda_2}{n_2 \sqrt{\frac{\lambda_1}{n_1} + \frac{\lambda_2}{n_2}}}, & k = n_1 + 1, \dots, n_1 + n_2 \end{cases}$$

where  $n = n_1 + n_2$ . Then,  $\mathbb{E} Z_{nk} = 0$  for all k = 1, ..., n, and

$$0 < \sigma_{nk}^2 = \operatorname{Var} Z_{nk} = \begin{cases} \frac{\operatorname{Var} X_{1k}}{n_1^2 \left(\frac{\lambda_1}{n_1} + \frac{\lambda_2}{n_2}\right)} = \frac{\lambda_1}{n_1^2 \left(\frac{\lambda_1}{n_1} + \frac{\lambda_2}{n_2}\right)}, & k = 1, \dots, n_1, \\ \frac{\lambda_2}{n_2^2 \left(\frac{\lambda_1}{n_1} + \frac{\lambda_2}{n_2}\right)}, & k = n_1 + 1, \dots, n, \end{cases}$$

Thus,

$$\sum_{k=1}^{n} \sigma_{nk}^{2} = \left(\frac{\lambda_{1}}{n_{1}^{2}} n_{1} + \frac{\lambda_{2}}{n_{2}^{2}} n_{2}\right) \frac{1}{\frac{\lambda_{1}}{n_{1}} + \frac{\lambda_{2}}{n_{2}}} = 1.$$

Furthermore, for  $\delta > 0$  and  $n_1, n_2 \to \infty$ 

$$\lim_{n \to \infty} \sum_{k=1}^{n} \mathbb{E}(|Z_{nk}|)^{2+\delta} = \lim_{n_1, n_2 \to \infty} \left( \frac{\mathbb{E}(|X_{11} - \lambda_1|^{2+\delta})}{n_1^{1+\delta} \left(\frac{\lambda_1}{n_1} + \frac{\lambda_2}{n_2}\right)^{(2+\delta)/2}} + \frac{\mathbb{E}(|X_{21} - \lambda_2|)^{2+\delta}}{n_2^{1+\delta} \left(\frac{\lambda_1}{n_1} + \frac{\lambda_2}{n_2}\right)^{(2+\delta)/2}} \right)$$

$$= 0$$

holds. The Ljapunow condition is therefore met and [33, Theorem 4.2.13] implies

$$\sum_{k=1}^{n} Z_{nk} \xrightarrow[n_1, n_2 \to \infty]{d} Y \sim \mathcal{N}(0, 1).$$

Finally, 
$$\sum_{k=1}^{n} Z_{nk} = \frac{\overline{X}_{1n_1} - \overline{X}_{2n_2} - \lambda_1 + \lambda_2}{\sqrt{\frac{\lambda_1}{n_1} + \frac{\lambda_2}{n_2}}}$$
, which completes the proof.

The strong law of large numbers implies  $\overline{X}_{in_i} \xrightarrow{f.s.} \lambda_i$ , i=1,2 and using Slutskys theorem then yields

$$T(Y_1, Y_2, \lambda_1 - \lambda_2) = \frac{\overline{X}_{1n_1} - \overline{X}_{2n_2} - \lambda_1 + \lambda_2}{\sqrt{\overline{X}_{1n_1}/n_1 + \overline{X}_{n_2}/n_2}} \xrightarrow[n_1, n_2 \to \infty]{d} Y \sim \mathcal{N}(0, 1).$$

The asymptotic confidence interval for  $\lambda_1 - \lambda_2$  with level  $1 - \alpha$  is thus given

$$\left[\overline{X}_{1n_{1}} - \overline{X}_{2n_{2}} - z_{1-\alpha/2} \sqrt{\frac{\overline{X}_{1n_{1}}}{n_{1}} + \frac{\overline{X}_{2n_{2}}}{n_{2}}}, \ \overline{X}_{1n_{1}} - \overline{X}_{2n_{2}} + z_{1-\alpha/2} \sqrt{\frac{\overline{X}_{1n_{1}}}{n_{1}} + \frac{\overline{X}_{2n_{2}}}{n_{2}}}\right]$$

# Asymptotic confidence interval for $\frac{n_2\lambda_2}{n_1\lambda_1+n_2\lambda_2}$

Let  $\rho$  be some constant,  $n_1/n_2=\rho$  and  $g(\lambda_1,\lambda_2)=\frac{n_2\lambda_2}{n_1\lambda_1+n_2\lambda_2}=\frac{\lambda_2}{\rho\lambda_1+\lambda_2}\stackrel{\mathrm{Def.}}{=} p$ . The goal is to construct an asymptotic confidence interval for p. Consider the statistic

$$T(Y_1, Y_2, p) = \frac{S_{2n_2} - p(S_{1n_1} + S_{2n_2})}{\sqrt{\hat{p}(1 - \hat{p})(S_{1n_1} + S_{2n_2})}},$$

where  $S_{in_i} = \sum_{i=1}^{n_i} X_{ij}, i = 1, 2$  and

$$\hat{p} = \frac{S_{2n_2}}{S_{1n_1} + S_{2n_2}} = \frac{n_2 \overline{X}_{2n_2}}{n_1 \overline{X}_{1n_1} + n_2 \overline{X}_{2n_2}} \xrightarrow[n_1, n_2 \to \infty]{f.s.} p$$

is a consistent estimator for p (by the strong law of large numbers). If it can be shown that  $T(Y_1, Y_2, p) \xrightarrow[n_1, n_2 \to \infty]{d} Y \sim \mathcal{N}(0, 1)$ , then

$$\lim_{n_1, n_2 \to \infty} P\left(-z_{1-\alpha/2} \le \frac{\frac{S_{2n_2}}{S_{1n_1} + S_{2n_2}} - p}{\sqrt{S_{1n_1} \cdot S_{2n_2}}} \cdot (S_{1n_1} + S_{2n_2})^{3/2} \le z_{1-\alpha/2}\right) = 1 - \alpha,$$

which yields the asymptotic confidence interval

$$\left[\underline{\theta}(Y_1,Y_2), \, \overline{\theta}(Y_1,Y_2), \, \right]$$

for p with level  $1-\alpha$ , where

$$\underline{\theta}(\lambda_1, \lambda_2) = \frac{S_{2n_2}}{S_{1n_1} + S_{2n_2}} - z_{1-\alpha/2} \cdot \sqrt{\frac{S_{1n_1} \cdot S_{2n_2}}{\left(S_{1n_1} + S_{2n_2}\right)^3}}$$

and

$$\overline{\theta}(\lambda_1, \lambda_2) = \frac{S_{2n_2}}{S_{1n_1} + S_{2n_2}} + z_{1-\alpha/2} \cdot \sqrt{\frac{S_{1n_1} \cdot S_{2n_2}}{\left(S_{1n_1} + S_{2n_2}\right)^3}}.$$

Since 0 , the boundaries can be adjusted as follows:

$$\underline{\theta}^*(Y_1, Y_2) = \max\{0, \underline{\theta}(Y_1, Y_2)\},$$
  
$$\overline{\theta}^*(Y_1, Y_2) = \min\{1, \overline{\theta}(Y_1, Y_2)\}.$$

Now the asymptotically normal distribution of  $T(Y_1, Y_2, p)$  will be shown. It results from Slutskys theorem and the following Lemma:

Lemma 2.3.2. One has the following property

$$\frac{S_{2n_2} - p(S_{1n_1} + S_{2n_2})}{\sqrt{p(1-p)(S_{1n_1} + S_{2n_2})}} \xrightarrow[n_1 \to \infty]{d} Y \sim \mathcal{N}(0,1)$$

**Proof** In order to show the assertion, a version of the central limit theorem for sums of random variables with random numbers as sum limits (cf. [33, Theorem 4.2.2]) is used. Let  $N_n = S_{1n_1} + S_{2n_2}$  be a sequence of nonnegative random variables, then the sum is monotonically increasing. Let then  $a_{n_2} = n_1\lambda_1 + n_2\lambda_2$ . Obviously

$$\begin{split} \frac{N_n}{a_{n_2}} &= \frac{S_{1n_1}}{n_1\lambda_1 + n_2\lambda_2} + \frac{S_{2n_2}}{n_1\lambda_1 + n_2\lambda_2} \\ &= \frac{\overline{X}_{1n_1}}{\lambda_1 + \rho^{-1}\lambda_2} + \frac{\overline{X}_{2n_2}}{\rho\lambda_1 + \lambda_2} \\ \underset{n_1, \overrightarrow{n_2} \to \infty}{\underbrace{\frac{f.s.}{\lambda_1 + \rho^{-1}\lambda_2} + \frac{\lambda_2}{\rho\lambda_1 + \lambda_2}}} \\ &= \frac{\rho\lambda_1}{\rho\lambda_1 + \lambda_2} + \frac{\lambda_2}{\rho\lambda_1 + \lambda_2} = 1 \end{split}$$

holds. Furthermore:

$$P(S_{2n_2} = k \mid N_n = m) = \frac{P(S_{2n_2} = k, S_{1n_1} + S_{2n_2} = m)}{P(S_{1n_1} + S_{2n_2} = m)}$$

$$= \frac{P(S_{2n_2} = k, S_{1n_1} = m - k)}{P(S_{1n_1} + S_{2n_2} = m)}$$

$$= \frac{e^{-n_2\lambda_2} \frac{(\lambda_2 n_2)^k}{k!} \cdot e^{-n_1\lambda_1} \frac{(n_1\lambda_1)^{m-k}}{(m-k)!}}{e^{-n_1\lambda_1 - n_2\lambda_2} \frac{(n_1\lambda_1 + n_2\lambda_2)^m}{m!}}$$

$$= \frac{m!}{(m-k)!k!} \left(\frac{n_2\lambda_2}{n_1\lambda_1 + n_2\lambda_2}\right)^m \left(\frac{n_1\lambda_1}{n_1\lambda_1 + n_2\lambda_2}\right)^{m-k}$$

$$= \binom{m}{k} p^k (1-p)^{m-k}$$

which means that  $S_{2n_2} \mid \{N_n = m\} \sim \text{Bin}(m,p)$ . Then,  $\frac{S_{2n_2} - mp}{\sqrt{mp(1-p)}} \mid \{N_n = m\} \stackrel{d}{=} \frac{S_m - mp}{\sqrt{mp(1-p)}}$ , where  $S_m = \sum_{i=1}^m Z_i$  is a sum of identically distributed  $Z_i \sim \text{Bernoulli}(p)$ . [33, Theorem 4.2.2] implies

$$\frac{S_{N_n} - N_n p}{\sqrt{N_n p(1-p)}} \xrightarrow{d} Y \sim \mathcal{N}(0,1) \iff \frac{S_{2n_2} - N_n p}{\sqrt{N_n p(1-p)}} \xrightarrow{d} Y \sim \mathcal{N}(0,1).$$

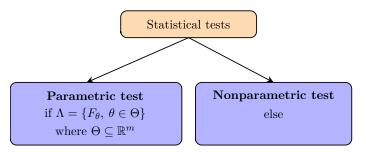
# Chapter 3

# Testing Statistical Hypotheses

In [33, Chapter 7], some tests like the Kolmogorow-Smironow test were introduced. This chapter, however, focuses on introducing tests for statistical significance formally.

# 3.1 General philosophy of testing

Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables  $X_i$  with distribution function  $F \in \Lambda$ , where  $\Lambda$  is some class of distributions. Let  $(x_1, \ldots, x_n)$  be a realization of the random sample  $(X_1, \ldots, X_n)$ . In statistical testing, hypotheses with respect to the nature of a (unknown) distribution F are posed and tested. Generally, two concepts are distinguished:



Parametric tests check, whether a parameter  $\theta$  attains certain values (e.g.  $\theta = 0$ ). Popular nonparametric tests are the so-called "goodness-of-fit tests", which check whether the distribution F is equal to a predetermined distribution  $F_0$ .

In an initial step, the term Hypotheses needs to be formalized. The set  $\Lambda$  of admissible distributions F is divided into two disjoint sets  $\Lambda_0$  and  $\Lambda_1$ ,  $\Lambda_0 \cup \Lambda_1 = \Lambda$ . The assertion

"The null hypothesis  $H_0: F \in \Lambda_0$  is tested against the alternative  $H_1: F \in \Lambda_1$ "

means, that we aim to assign the distribution function of the random variable  $X_i$  to  $\Lambda_0$  or  $\Lambda_1$ , based on the explicit realization  $(x_1, \ldots, x_n)$ . The process of assigning the distribution of  $X_i$  involves a decision rule

$$\varphi: \mathbb{R}^n \to [0,1],$$

which is a statistic with the following interpretation:

The sample space  $\mathbb{R}^n$  is divided into 3 disjoint sets  $K_0, K_{01}$   $K_1$ , such that  $\mathbb{R}^n = K_0 \cup K_{01} \cup K_1$ , where

$$K_0 = \varphi^{-1}(\{0\}) = \{x \in \mathbb{R}^n : \varphi(x) = 0\},$$

$$K_1 = \varphi^{-1}(\{1\}) = \{x \in \mathbb{R}^n : \varphi(x) = 1\},$$

$$K_{01} = \varphi^{-1}((0,1)) = \{x \in \mathbb{R}^n : 0 < \varphi(x) < 1\}.$$

Thus  $H_0: F \in \Lambda_0$  is

- rejected, if  $\varphi(x) = 1$ , i.e.,  $x \in K_1$ ,
- not rejected, if  $\varphi(x) = 0$ , i.e.,  $x \in K_0$ .

If  $\varphi(x) \in (0,1)$ , i.e.,  $x \in K_{01}$ , then  $\varphi(x)$  is interpreted as a Bernoulli probability and a random variable  $Y \sim \text{Bernoulli}(\varphi(x))$  is generated with

$$Y = \begin{cases} 1 & \Longrightarrow H_0 \text{ is rejected} \\ 0 & \Longrightarrow H_0 \text{ is not rejected} \end{cases}$$

If  $K_{01} \neq \emptyset$ , the decision rule is called randomized. If  $K_{01} = \emptyset$ , i.e.,  $\mathbb{R}^n = K_0 \cup K_1$  the tests are called non-randomized.  $K_0$  and  $K_1$  are called acceptance region and rejection region (critical region) of  $H_0$  respectively.  $K_{01}$  is called randomization region.

#### Remark 3.1.1.

1. One deliberately says " $H_0$  is not rejected", instead of " $H_0$  is accepted", since statistical inference can generally not make positive decisions rather than negative decisions. The issue above is a general philosophical problem with respect to the falsifiability of hypotheses or scientific theories, which can generally not be at odds with the truth. (cf. wissenschaftliche Erkenntnistheorie by Karl Popper (1902-1994)).

2. Randomized tests are generally of a more theoretic nature (cf. Section 3.3). In practice, most non-randomized rules are used, which leads to a decision with respect to  $H_0$  based on the explicit sample  $(x_1, \ldots, x_n)$  alone. Here,  $\varphi(x) = I_{K_1}, x = (x_1, \ldots, x_n) \in \mathbb{R}^n$  holds.

In the following paragraph, non-randomized tests are considered in order to return to the more general approach in Section 3.3.

**Definition 3.1.2.** The non-randomized test rule  $\varphi : \mathbb{R}^n \to \{0,1\}$  provides a *(non-randomized) statistical test with significance level*  $\alpha$ , if for  $F \in \Lambda_0$ 

$$P_F(\varphi(X_1,\ldots,X_n)=1)=P(H_0 \text{ reject } \mid H_0 \text{ true }) \leq \alpha.$$

### Definition 3.1.3.

1. If  $H_0$  is rejected, even though  $H_0$  is correct, then a type I error has occurred. The probability

$$\alpha_n(F) = P_F(\varphi(x_1, \dots, x_n) = 1), \quad F \in \Lambda_0,$$

is called *Probability of a type I error*. This probability is supposed to be lower than the significance level  $\alpha$ .

2. A type II error occurs, if a wrong hypothesis  $H_0$  is not rejected. Here

$$\beta_n(F) = P_F(\varphi(x_1, \dots, x_n) = 0), \quad F \in \Lambda_1,$$

is called Probability of a type II error.

A summary of all possible errors can be found in the following matrix, which is called *confusion matrix*:

	$H_0$ true	$H_0$ false
reject $H_0$	Error of type I with probability $\alpha_n(F) \leq \alpha$	right decision
not rejecting $H_0$	right decision	Error of type II with probability $\beta_n(F)$

Here  $\alpha_n$  and  $\beta_n$  are aimed to be small, which is contrary to the fact that a decreasing value of  $\alpha$  increases the probability of mistakes of type II.

#### Definition 3.1.4.

1. The function

$$G_n(F) = P_F(\varphi(X_1, \dots, X_n) = 1), \quad F \in \Lambda$$

is called performance function (or power function) of a test  $\varphi$ .

2. The constraint  $G_n$  on  $\Lambda_1$  is called *power* of the test  $\varphi$ . With respect to the constraint it holds

$$\begin{cases} G_n(F) = \alpha_n(F) \le \alpha, F \in \Lambda_0, \\ G_n(F) = 1 - \beta_n(F), F \in \Lambda_1. \end{cases}$$

# Example 3.1.5. Parametric tests.

What does a parametric test look like? The parameter space  $\Theta$  is given by  $\Theta_0 \cup \Theta_1$ , where  $\Theta_0 \cap \Theta_1 = \emptyset$ . Then

$$\Lambda_0 = \{ F_\theta : \theta \in \Theta_0 \},$$
  
$$\Lambda_1 = \{ F_\theta : \theta \in \Theta_1 \}.$$

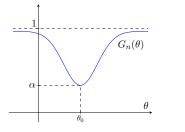
 $P_F$  is replaced by  $P_{\theta}$ . Furthermore  $\alpha_n, G_n$  and  $\beta_n$  are defined on  $\Theta$  instead of  $\Lambda$ .

Which hypotheses  $H_0$  and  $H_1$  are popular among parametric tests? The case  $\Theta = \mathbb{R}$  is discussed below, but it should be noted that a more general choice of  $\Theta$  is also possible.

- 1.  $H_0: \theta = \theta_0 \text{ vs. } H_1: \theta \neq \theta_0$
- 2.  $H_0: \theta \ge \theta_0 \text{ vs. } H_1: \theta < \theta_0$
- 3.  $H_0: \theta \le \theta_0 \text{ vs. } H_1: \theta > \theta_0$
- 4.  $H_0: \theta \in [a, b] \text{ vs. } H_1: \theta \notin [a, b]$

In the first case the, parametric test is called *two-sided* and in the second and third case *one-sided* (*right-* resp. *left-sided*). The fourth case is called *interval hypothesis*  $H_0$ .

Considering a one-sided or two-sided test, the power function may look like the one displayed in Figure 3.1 (a) or 3.1 (b), resp.



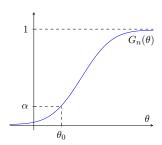


Figure 3.1: Performance function

In general models (not necessarily parametric), the ideal power function can be illustrated schematically, as in Figure 3.2.

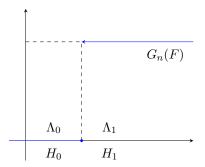


Figure 3.2: Schematic illustration of an ideal power function

- Definition 3.1.3, errors of type I and II and the rejection rule imply that the hypotheses  $H_0$  and  $H_1$  can not be treated symmetrically, since only the probability of errors of type I is controlled. That is the reason why statisticians mostly formulate the hypothesis of interest as  $H_1$  instead of  $H_0$ , because if one decides that  $H_0$  can be rejected, it can be assured that the probability of a false decision is below the significance level  $\alpha$ .
- How is a statistical, non-randomized test constructed in practice? The construction of the rejection rule  $\varphi$  is very similar to constructing confidence intervals:
  - 1. Find a test statistic  $T: \mathbb{R}^n \to \mathbb{R}$ , which has a certain test distribution (perhaps asymptotically for  $n \to \infty$ ) under  $H_0$ .
  - 2. Define  $B_0 = [t_{\alpha_1}, t_{1-\alpha_2}]$ , where  $t_{\alpha_1}$  and  $t_{1-\alpha_2}$  are quantiles of the test distribution of T with  $\alpha_1 + \alpha_2 = \alpha \in [0, 1]$ .
  - 3. If  $T(X_1, \ldots, X_n) \in \mathbb{R} \setminus B_0 = B_1$ , then set  $\varphi(X_1, \ldots, X_n) = 1$  and reject  $H_0$ . Else, set  $\varphi(X_1, \ldots, X_n) = 0$ .
- If the distribution of T can only be determined asymptotically, then  $\varphi$  is called *asymptotic test*.
- Most of the times, even the asymptotic distribution of T is unknown. In this case, the so called *Monte-Carlo tests* come into play. In those tests the quantiles  $t_{\alpha}$  are determined approximatively by conducting a large number of Monte-Carlo simulations of T (under  $H_0$ ):

If  $t^i$ ,  $i=1,\ldots,m$  takes the values of T in m independent simulations, i.e.  $t^i=T(x_1^i,\ldots,x_n^i)$ , where  $x_j^i$  are independent realizations of  $X_j\sim F\in \Lambda_0$  for  $j=1,\ldots,n,\ i=1,\ldots,m$ , then  $t_\alpha\approx t^{(\lfloor\alpha\cdot m\rfloor),1}$  with  $t^{(1)},\ldots,t^{(m)}$  the order statistics and  $\alpha\in[0,1]$ .

 $<sup>^{1}\</sup>mathrm{set}\ t^{(0)} = -\infty$ 

Remark 3.1.6. It is easy to see that by using an arbitrary confidence interval

$$I_{\theta} = \left[ I_1^{\theta}(X_1, \dots, X_n), I_2^{\theta}(X_1, \dots, X_n) \right]$$

with confidence level  $1 - \alpha$  for a parameter  $\theta \in \mathbb{R}$ , a test for  $\theta$  can be constructed. The hypotheses  $H_0: \theta = \theta_0$  vs.  $H_1: \theta \neq \theta_0$  are tested under the following decision rule:

$$\varphi(X_1, \dots, X_n) = 1$$
, if  $\theta_0 \notin \left[ I_1^{\theta_0}(X_1, \dots, X_n), I_2^{\theta_0}(X_1, \dots, X_n) \right]$ .

The significance level of the test is  $\alpha$ .

**Example 3.1.7.** Normal distribution, testing the expected value with known variance. Let

$$X_1, \ldots, X_n \sim \mathcal{N}(\mu, \sigma^2)$$

with known variance  $\sigma^2$ . A confidence interval for  $\mu$  is given by

$$I^{\mu} = [I_1^{\mu}(X_1, \dots, X_n), I_2^{\mu}(X_1, \dots, X_n)] = \left[ \overline{X}_n - \frac{z_{1-\alpha/2} \cdot \sigma}{\sqrt{n}}, \overline{X}_n + \frac{z_{1-\alpha/2} \cdot \sigma}{\sqrt{n}} \right]$$

(cf. Section 2.2.1). Hence,  $H_0: \mu = \mu_0$  (versus the alternative  $H_1: \mu \neq \mu_0$ ), is rejected, if

$$|\mu_0 - \overline{X}_n| > \frac{z_{1-\alpha/2} \cdot \sigma}{\sqrt{n}}.$$

In the language of testing, the above can be rewritten as

$$\varphi(x_1,\ldots,x_n)=I\left((x_1,\ldots x_n)\in K_1\right),\,$$

where

$$K_1 = \left\{ (x_1, \dots, x_n) \in \mathbb{R}^n : |\mu_0 - \overline{x}_n| > \frac{z_{1-\alpha/2} \cdot \sigma}{\sqrt{n}} \right\}$$

is the rejection region. For the test statistic  $T(X_1, \ldots, X_n)$ 

$$T(X_1, \dots, X_n) = \frac{\overline{X}_n - \mu_0}{\sigma} \sqrt{n} \sim \mathcal{N}(0, 1)$$

under  $H_0$  it holds  $\alpha_n(\mu) = \alpha$ .

The power function (cf. Figure 3.3) can be calculated as follows

$$\begin{split} G_n(\mu) &= P_\mu \left( |\mu_0 - \overline{X}_n| > \frac{z_{1-\alpha/2}}{\sqrt{n}} \right) = 1 - P_\mu \left( \left| \overline{X}_n - \mu_0 \right| \leq \frac{\sigma z_{1-\alpha/2}}{\sqrt{n}} \right) \\ &= 1 - P_\mu \left( \left| \sqrt{n} \frac{\overline{X}_n - \mu}{\sigma} + \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right| \leq z_{1-\alpha/2} \right) \\ &= 1 - P_\mu \left( -z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \leq \sqrt{n} \frac{\overline{X}_n - \mu}{\sigma} \leq z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right) \\ &= 1 - \Phi \left( z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right) + \Phi \left( -z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right) \\ &= \Phi \left( -z_{1-\alpha/2} + \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right) + \Phi \left( -z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right). \end{split}$$

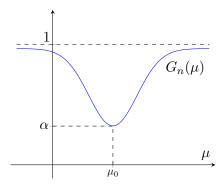


Figure 3.3: Performance function of a two-sided test of the expected value of a normal distribution with known variance.

The "yes-no" decision of tests is mainly viewed as too rough. That is why it is desirable to obtain a finer measure for the data with respect to the hypotheses  $H_0$  and  $H_1$ . The so-called *p-value* solves the problem above and it is luckily included in most statistic software.

**Definition 3.1.8.** Let  $(x_1, \ldots, x_n)$  be an explicit sample, i.e., a realization of  $(X_1, \ldots, X_n)$  and  $T(X_1, \ldots, X_n)$  the test statistic which was used to construct the decision rule  $\varphi$ . The p-value of the test  $\varphi$  is the smallest significance level to the value  $t = T(x_1, \ldots, x_n)$  which leads to a rejection of  $H_0$ .

In the example of a one-sided test with rejection region  $B_1 = (t, \infty)$ , the rule of thumb for p is given by

$$p = "P(T(X_1, ..., X_n) \ge t \mid H_0)",$$

where the quotation marks imply that the term is not a probability in the classical sense, rather than a conditional probability, which will be defined more precisely later.

Using the p-value, the rejection rule changes: The hypothesis  $H_0$  is rejected with a significance level  $\alpha$ , if  $\alpha \geq p$ . Previously, the significance of a test (rejection of  $H_0$ ) was determined with respect to the following table:

p-value	interpretation	
$p \le 0,001$	very strongly significant	
$0,001$	strongly significant	
$0,01$	weakly significant	
p > 0,05	not significant	

Since the p-value can be calculated with ease nowadays, one can use the p-value directly to decide which significance level is sufficient for the underlying test.

#### Remark 3.1.9.

- 1. The significance level must not depend on p. By doing so, the general philosophy of testing is jeopardized!
- 2. The p-value is not a probability rather than a random variable since it depends on  $(X_1, \ldots, X_n)$ . The expression

$$p = P\left(T(X_1, \dots, X_n) \ge t \mid H_0\right),\,$$

in Definition 3.1.8 for the *p*-value of an one-sided test with test statistic T can be interpreted as an exceedance probability. The exceedance probability is defined with respect to  $t = T(x_1, \ldots, x_n)$  or more extreme values in order to be close to the hypothesis  $H_1$  while repeating the random experiment under  $H_0$ :

$$p = P(T(X'_1, ..., X'_n) \ge T(x_1, ..., x_n) \mid H_0),$$

where  $(X'_1, \ldots, X'_n) \stackrel{d}{=} (X_1, \ldots, X_n)$ . If instead of the explicit sample  $(x_1, \ldots, x_n)$  the random sample  $(X_1, \ldots, X_n)$  is used, then

$$p = p(X_1, \dots, X_n) = P(T(X_1', \dots, X_n') \ge T(X_1, \dots, X_n) \mid H_0, X_1, \dots, X_n).$$

- 3. For other hypotheses  $H_0$ , the p-values may look different. For example:
  - (a) A symmetric two-sided test has an acceptance region

$$B_0 = \left[ -t_{1-\alpha/2}, \, t_{1-\alpha/2} \right]$$

for  $H_0$ . Therefore

$$p = P(|T(X'_1, \dots, X'_n)| \ge T(X_1, \dots, X_n) | H_0, X_1, \dots, X_n).$$

(b) A left-sided test with  $B_0 = [t_\alpha, \infty]$  results in

$$p = P(T(X'_1, \dots, X'_n) \le T(X_1, \dots, X_n) | H_0, X_1, \dots, X_n).$$

4. The behavior of the p-value can be evaluated using the following lemma

**Lemma 3.1.10.** If the distribution function F of T is continuous and monotonically increasing (e.g., the distribution T is absolutely continuous with continuous probability density function for example), then  $p \sim U[0, 1]$ .

**Proof** The result will be shown for right-sided tests only.

$$P(p \le \alpha \mid H_0) = P(\overline{F}_T(T(X_1, \dots, X_n)) \le \alpha \mid H_0)$$

$$= P(F_T(T(X_1, \dots, X_n)) \ge 1 - \alpha \mid H_0)$$

$$= P(U \ge 1 - \alpha) = 1 - (1 - \alpha) = \alpha, \quad \alpha \in [0, 1],$$

since  $F_T(T(X_1, ..., X_n)) \stackrel{d}{=} U \sim U[0, 1]$ , and  $F_T$  is absolutely continuous.

**Exercise 3.1.11.** Show that for an arbitrary random variable X with continuous and monotonically increasing distribution function  $F_X$ 

$$F_X(X) \sim U[0,1]$$

holds.

If the distribution of T with domain  $\{t_1, \ldots, t_n\}$ ,  $t_i < t_j$ , is discrete for i < j, then the distribution of p is also discrete. In particular, it does not hold that  $p \sim U[0,1]$ . In this case  $F_p(x)$  is a step function which touches the line y = u at the points  $u_k = \sum_{i=1}^k P(T(X_1, \ldots, X_n) = t_i)$ ,  $k = 1, \ldots, n$  (cf. Figure 4).

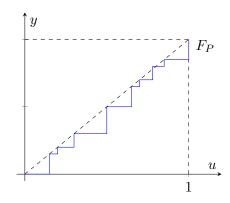


Figure 3.4: Distribution of p for discrete T

## Definition 3.1.12.

1. If the power  $G_n(\cdot)$  of a test  $\varphi$  with significance level  $\alpha$  satisfies the inequality

$$G_n(F) \ge \alpha, \quad F \in \Lambda_1,$$

then the test is called *unbiased*.

2. Let  $\varphi$  and  $\varphi^*$  be two tests with significance level  $\alpha$  and power functions  $G_n(\cdot)$  and  $G_n^*(\cdot)$ . The test  $\varphi$  is said to be *more powerful* than  $\varphi^*$  if its power is larger:

$$G_n(F) \ge G_n^*(F) \quad \forall F \in \Lambda_1.$$

3. The test  $\varphi$  is called consistent, if  $G_n(F) \xrightarrow[n \to \infty]{} 1$  for all  $F \in \Lambda_1$ .

#### Remark 3.1.13.

1. The power of a one-sided test is mostly larger than the two-sided version:

**Example 3.1.14.** Consider the Gauss-test for the expected value of the normal distribution if the variance is known. The two-sided test

$$H_0: \mu = \mu_0 \text{ vs. } H_1: \mu \neq \mu_0.$$

implies that the power function is given by

$$G_n(\mu) = \Phi\left(-z_{1-\alpha/2} + \sqrt{n}\frac{\mu - \mu_0}{\sigma}\right) + \Phi\left(-z_{1-\alpha/2} - \sqrt{n}\frac{\mu - \mu_0}{\sigma}\right).$$

The one-sided test  $\varphi^*$  of the hypotheses

$$H_0^*: \mu \leq \mu_0 \text{ vs. } H_1^*: \mu > \mu_0$$

attains a power function given by

$$G_n^*(\mu) = \Phi\left(-z_{1-\alpha} + \sqrt{n}\frac{\mu - \mu_0}{\sigma}\right).$$

Since  $G_n(\mu) \xrightarrow[n\to\infty]{} 1$ ,  $G_n^*(\mu) \xrightarrow[n\to\infty]{} 1$  both tests are consistent. In the case above,  $\varphi^*$  is more powerful than  $\varphi$ . Moreover, both tests are unbiased (cf. Figure 3.1.14).

- 2. For testing interval hypotheses  $H_0: \theta \in [a, b]$  vs.  $H_1: \theta \notin [a, b]$  with confidence level  $\alpha$  the following methodology can be used: Test
  - (a)  $H_0^a: \theta \geq a$  vs.  $H_1^a: \theta < a$  with significance level  $\alpha/2$ ,
  - (b)  $H_0^b: \theta \leq b$  vs.  $H_1^b: \theta > b$  with significance level  $\alpha/2$ .

 $H_0$  is not rejected if  $H_0^a$  and  $H_0^b$  are not rejected. The probability for a type I error is  $\alpha$ . The power of those tests is generally low.

3. As a rule of thumb, it holds that an increase in parameters that need to be estimated with respect to the test statistic leads to a decrease in power.

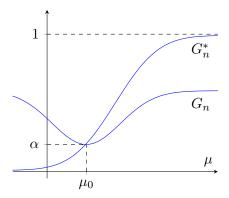


Figure 3.5: Power function of an one-sided and two-sided test for the expected value of a normal distribution.

# 3.2 Non-randomized tests

# 3.2.1 Parametric significance tests

This section provides examples of tests that can mostly be obtained from their corresponding confidence intervals for the parameters of distributions.

- 1. Tests for the parameters of a normal distribution  $\mathcal{N}(\mu, \sigma^2)$ 
  - (a) Test of  $\mu$  with unknown variance
    - Hypotheses:  $H_0: \mu = \mu_0$  vs.  $H_1: \mu \neq \mu_0$ .
    - Test statistic:

$$T(X_1, \dots, X_n) = \frac{\overline{X}_n - \mu_0}{S_n} \sim t_{n-1} \quad | H_0$$

• Decision rule:

$$\varphi(X_1, \dots, X_n) = 1$$
, if  $|T(X_1, \dots, X_n)| > t_{n-1, 1-\alpha/2}$ .

- (b) Test of  $\sigma^2$  with known  $\mu$ 
  - Hypotheses:  $H_0: \sigma^2 = \sigma_0^2$  vs.  $H_1: \sigma^2 \neq \sigma_0^2$ .
  - Test statistic:

$$T(X_1, \dots, X_n) = \frac{n\tilde{S}_n^2}{\sigma_0^2} \sim \chi_n^2 \quad | H_0$$

with 
$$\tilde{S}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$$
.

• Decision rule:

$$\varphi(X_1, \dots, X_n) = 1$$
, if  $T(X_1, \dots, X_n) \notin \left[\chi_{n,\alpha/2}^2, \chi_{n,1-\alpha/2}^2\right]$ .

• Performance function:

$$G_n(\sigma^2) = 1 - P_{\sigma^2} \left( \chi_{n,\alpha/2}^2 \le \frac{n\tilde{S}_n^2}{\sigma_0^2} \le \chi_{n,1-\alpha/2}^2 \right)$$

$$= 1 - P_{\sigma^2} \left( \frac{\chi_{n,\alpha/2}^2 \sigma_0^2}{\sigma^2} \le \frac{n\tilde{S}_n^2}{\sigma^2} \le \frac{\chi_{n,1-\alpha/2}^2 \sigma_0^2}{\sigma^2} \right)$$

$$= 1 - F_{\chi_n^2} \left( \chi_{n,1-\alpha/2}^2 \frac{\sigma_0^2}{\sigma^2} \right) + F_{\chi_n^2} \left( \chi_{n,\alpha/2}^2 \frac{\sigma_0^2}{\sigma_2} \right)$$

- (c) Test for  $\sigma^2$  with unknown  $\mu$ 
  - Hypotheses:  $H_0: \sigma^2 = \sigma_0^2$  vs.  $H_1: \sigma^2 \neq \sigma_0^2$ .
  - Test statistic:

$$T(X_1, \dots, X_n) = \frac{(n-1)S_n^2}{\sigma_0^2} \sim \chi_{n-1}^2 \mid H_0,$$

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

• Decision rule:

$$\varphi(X_1,\ldots,X_n)=1$$
, if  $T(X_1,\ldots,X_n)\notin \left[\chi_{n-1,\alpha/2}^2,\chi_{n-1,1-\alpha/2}^2\right]$ .

#### Exercise 3.2.1.

- i. Find  $G_n(\cdot)$  for the one-sided versions of the tests above.
- ii. Show that the one-sided tests are unbiased, contrary to the two-sided tests being biased.

#### 2. Asymptotic tests

Considering asymptotic tests, the test statistic distribution can only be estimated (for large n). In the same spirit, the confidence level  $\alpha$  is obtained. Its construction is mostly based on limit theorems.

The general methodology is introduced via the Wald test (named after the statistician Abraham Wald (1902-1980)):

- Let  $(X_1, ..., X_n)$  be a random sample and  $X_i$  be i.i.d. for i = 1, ..., n, with  $X_i \sim F_\theta$ ,  $\theta \in \Theta \subseteq \mathbb{R}$ .
- $H_0: \theta = \theta_0$  vs.  $H_1: \theta \neq \theta_0$  is tested. Let  $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n)$  be an asymptotically normal distributed estimator for  $\theta$ . Let

$$\frac{\hat{\theta}_n - \theta_0}{\hat{\sigma}_n} \xrightarrow[n \to \infty]{d} Y \sim \mathcal{N}(0, 1) \quad | H_0,$$

where  $\hat{\sigma}_n^2$  is a consistent estimator for the variance of  $\hat{\theta}_n$ . The test statistic is given by

$$T(X_1,\ldots,X_n) = \frac{\hat{\theta}_n(X_1,\ldots,X_n) - \theta_0}{\hat{\sigma}_n}$$

• The decision rule is:  $H_0$  is rejected, if

$$|T(X_1, ..., X_n)| > z_{1-\alpha/2}$$
, where  $z_{1-\alpha/2} = \Phi^{-1}(1 - \alpha/2)$ .

This decision rule should only be applied for large n. The probability of a type I error is equal to  $\alpha$ , since  $P(|T(X_1,\ldots,X_n)|>z_{1-\alpha/2}\mid H_0)\underset{n\to\infty}{\longrightarrow}\alpha$  because of the asymptotically normal distribution of T.

The power function of the test is asymptotically given by

$$\lim_{n \to \infty} G_n(\theta) = 1 - \Phi\left(z_{1-\alpha/2} + \frac{\theta_0 - \theta}{\sigma}\right) + \Phi\left(-z_{1-\alpha/2} + \frac{\theta_0 - \theta}{\sigma}\right),\,$$

where  $\hat{\sigma}_n^2 \xrightarrow[n \to \infty]{P} \sigma^2$ .

Special cases of the Wald test are asymptotic tests for the expected value of Poisson or Bernoulli distributed random samples.

#### **Example 3.2.2**

#### (a) Bernoulli distribution

Let  $X_i \sim \text{Bernoulli}(p), p \in (0,1)$  be i.i.d. random variables.

- Hypotheses:  $H_0: p = p_0 \text{ vs. } H_1: p \neq p_0.$
- Test statistic:

$$T(X_1, \dots, X_n) = \begin{cases} \sqrt{n} \frac{\overline{X}_n - p_0}{\sqrt{\overline{X}_n} (1 - \overline{X}_n)}, & \text{if } \overline{X}_n \neq 0, 1, \\ 0, & \text{otherwise.} \end{cases}$$

Under  $H_0$ ,  $T(X_1, \ldots, X_n) \xrightarrow[n \to \infty]{d} Y \sim \mathcal{N}(0, 1)$  holds.

#### (b) Poisson distribution

Let  $X_i \sim \text{Poisson}(\lambda)$ ,  $\lambda > 0$  be i.i.d. random variables.

- Hypotheses:  $H_0: \lambda = \lambda_0$  vs.  $H_1: \lambda \neq \lambda_0$ .
- Test statistic:

$$T(X_1, \dots, X_n) = \begin{cases} \sqrt{n} \frac{\overline{X}_n - \lambda_0}{\sqrt{\overline{X}_n}}, & \text{if } \overline{X}_n > 0, \\ 0, & \text{otherwise.} \end{cases}$$

Under  $H_0, T(X_1, \ldots, X_n) \xrightarrow[n \to \infty]{d} Y \sim N(0, 1)$  holds.

#### 3. Two sample problems

Let

$$Y_1 = (X_{11}, \dots, X_{1n_1}), Y_2 = (X_{21}, \dots, X_{2n_2}), n = \max\{n_1, n_2\}$$

be two random samples. Assume that  $X_{ij}$  are independent for  $j = 1, \ldots, n_i, X_{ij} \sim F_{\theta_i}, i = 1, 2.$ 

- (a) Test for the equality of two expected values for normal distributed random samples
  - with known variance Let  $X_{ij} \sim \mathcal{N}(\mu_i, \sigma_i^2)$ ,  $i = 1, 2, j = 1, ..., n_i$ . Here,  $\sigma_1^2, \sigma_2^2$  are known and  $X_{ij}$  are independent for all i, j.

The hypotheses are given by  $H_0: \mu_1 = \mu_2$  vs.  $H_1: \mu_1 \neq \mu_2$ . Consider the test statistic

$$T(Y_1, Y_2) = \frac{\overline{X}_{1n_1} - \overline{X}_{2n_2}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}.$$

Under  $H_0$ ,  $T(Y_1, Y_2) \sim \mathcal{N}(0, 1)$  holds. The decision rule is given by:  $H_0$  is rejected if  $|T(Y_1, Y_2)| > z_{1-\alpha/2}$ .

• with unknown but equal variances Let  $X_{ij} \sim \mathcal{N}(\mu_i, \sigma_i^2)$ ,  $i = 1, 2, j = 1, \dots, n_i$ . Here,  $\sigma_1^2, \sigma_2^2$  are unknown,  $\sigma_1^2 = \sigma_2^2$  and  $X_{ij}$  are independent for all i, j. The hypotheses are given by  $H_0: \mu_1 = \mu_2$  vs.  $H_1: \mu_1 \neq \mu_2$ . Consider the test statistic

$$T(Y_1, Y_2) = \frac{\overline{X}_{1n_1} - \overline{X}_{2n_2}}{S_{n_1 n_2}} \sqrt{\frac{n_1 n_2}{n_1 + n_2}},$$

where  $S_{n_1n_2}^2$  is given by

$$\frac{1}{n_1 + n_2 - 2} \cdot \left( \sum_{j=1}^{n_1} \left( X_{1j} - \overline{X}_{1n_1} \right)^2 + \sum_{j=1}^{n_2} \left( X_{2j} - \overline{X}_{2n_2} \right)^2 \right).$$

It can be shown that under  $H_0$   $T(Y_1, Y_2) \sim t_{n_1+n_2-2}$  holds. The decision rule is then given by:  $H_0$  is rejected if  $|T(Y_1, Y_2)| > t_{n_1+n_2-2,1-\alpha/2}$ .

(b) Test for the equality of the expected value for linked random samples

Let 
$$Y_1 = (X_{11}, \dots, X_{1n})$$
 and  $Y_2 = (X_{21}, \dots, X_{2n}), n_1 = n_2 = n$ ,

$$Z_j = X_{1j} - X_{2j} \sim \mathcal{N}(\mu_1 - \mu_2, \sigma^2), j = 1, \dots, n,$$

be i.i.d. with  $\mu_i = \mathbb{E} X_{ij}$ , i = 1, 2. The hypotheses are given by:  $H_0: \mu_1 = \mu_2$  vs.  $H_1: \mu_1 \neq \mu_2$  with unknown variance  $\sigma^2$ . The test statistic is given by

$$T(Z_1, \dots, Z_n) = \sqrt{n} \frac{\overline{Z}_n}{S_n},$$

where

$$S_n^2 = \frac{1}{n-1} \sum_{j=1}^n \left( Z_j - \overline{Z}_n \right)^2.$$

Under  $H_0, T(Z_1, \ldots, Z_n) \sim t_{n-1}$  holds. The decision rule is given by:  $H_0$  is rejected, if  $|T(Z_1, \ldots, Z_n)| > t_{n-1,1-\alpha/2}$ .

# (c) Test for the equality of variances for independent Gaussian random samples

Let  $Y_1 = (X_{11}, \ldots, X_{1n_1})$  and  $Y_2 = (X_{21}, \ldots, X_{2n_2})$  be i.i.d. with  $X_{ij} \sim \mathcal{N}(\mu_i, \sigma_i^2)$ , where  $\mu_i$  and  $\sigma_i^2$  are both unknown. The hypotheses are  $H_0: \sigma_1^2 = \sigma_2^2$  vs.  $H_1: \sigma_1^2 \neq \sigma_2^2$ . The test statistic is given by

$$T(Y_1, Y_2) = \frac{S_{2n_2}^2}{S_{1n_1}^2},$$

where

$$S_{in_i}^2 = \frac{1}{n_i - 1} \sum_{j=1}^n \left( X_{ij} - \overline{X}_{in_i} \right)^2, i = 1, 2.$$

Under  $H_0$ ,  $T(Y_1, Y_2) \sim F_{n_2-1,n_1-1}$  holds. The decision rule is then given by:  $H_0$  is rejected, if

$$T(Y_1, Y_2) \notin \left[ F_{n_2-1, n_1-1, \alpha/2}, F_{n_2-1, n_1-1, 1-\alpha/2} \right].$$

## (d) Asymptotic two sample tests

• for Bernoulli distributed random samples Let  $X_{ij} \sim \text{Bernoulli}(p_i), j = 1, ..., n_i, p_i \in (0, 1), i = 1, 2.$ The hypotheses are given by  $H_0: p_1 = p_2$  vs.  $H_1: p_1 \neq p_2$ . The test statistic is then given by

$$T(Y_1, Y_2) = \begin{cases} \frac{\overline{X}_{1n_1} - \overline{X}_{2n_2}}{\sqrt{\frac{\overline{X}_{1n_1}(1 - \overline{X}_{1n_1})}{n_1} + \frac{\overline{X}_{2n_2}(1 - \overline{X}_{2n_2})}{n_2}}} \\ 0, & \overline{X}_{1n_1} = \overline{X}_{2n_2} \in \{0, 1\} \end{cases}$$

Under  $H_0$ ,  $T(Y_1, Y_2) \xrightarrow[n_1, n_2 \to \infty]{d} Y \sim \mathcal{N}(0, 1)$  holds. The decision rule is then given by:  $H_0$  is rejected if  $|T(Y_1, Y_2)| > z_{1-\alpha/2}$ . This is a test with asymptotic confidence level  $\alpha$ .

#### for Poisson distributed random samples

Let  $X_{ij}$  be independent,  $X_{ij} \sim \text{Poisson}(\lambda_i)$ ,  $\lambda_i > 0$ , i = 1, 2. The hypotheses are:  $H_0: \lambda_1 = \lambda_2$  vs.  $H_1: \lambda_1 \neq \lambda_2$  and the test statistic is given by:

$$T(Y_1, Y_2) = \begin{cases} \frac{\overline{X}_{1n_1} - \overline{X}_{2n_2}}{\sqrt{\frac{\overline{X}_{1n_1}}{n_1} + \frac{\overline{X}_{2n_2}}{n_2}}} \\ 0, \quad \overline{X}_{1n_1} = \overline{X}_{2n_2} = 0 \end{cases}$$

The decision rule is then given by:  $H_0$  is rejected, if  $|T(Y_1, Y_2)| > z_{1-\alpha/2}$ . This is a test with asymptotic confidence level  $\alpha$ .

Remark 3.2.3. Asymptotic tests must only be used for large samples, since for small samples the, asymptotic significance level can not be assured.

## 3.3 Randomized test

In this section, classical results of Neyman-Pearson with respect to the terminology of *most powerful tests* are presented. Here, randomized tests play a considerably important role.

#### 3.3.1 Fundamentals

Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables  $X_i$  and  $(x_1, \ldots, x_n)$  a realization of  $(X_1, \ldots, X_n)$ . Assume that that the sample space  $(B, \mathcal{B})$  is either given by  $(\mathbb{R}^n, \mathcal{B}_{\mathbb{R}^n})$  or  $(\mathbb{N}_0^n, \mathcal{B}_{\mathbb{N}_0^n})$  depending on whether the distribution of  $X_i$ ,  $i = 1, \ldots, n$  is either absolutely continuous or discrete. If the random variables  $X_i$  are discrete, the domain is assumed to be  $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$ . The domain is equipped with a measure  $\mu$ , where

$$\mu = \begin{cases} \text{ Lebesgue measure on } \mathbb{R}, & \text{if } B = \mathbb{R}^n, \\ \text{ Counting measure on } \mathbb{N}_0, & \text{if } B = \mathbb{N}_0^n. \end{cases}$$

Thus

$$\int g(x)\mu(dx) = \begin{cases} \int_{\mathbb{R}} g(x)dx, & \text{in the absolutely continuous case,} \\ \sum_{x \in \mathbb{N}_0} g(x), & \text{in the discrete case,} \end{cases}$$

holds. Moreover, assume that  $X_i \sim F_\theta$ ,  $\theta \in \Theta \subseteq \mathbb{R}^m$ , i = 1, ..., n (parametric model). For  $\Theta = \Theta_0 \cup \Theta_1$ ,  $\Theta_0 \cap \Theta_1 = \emptyset$  the hypotheses are  $H_0 : \theta \in \Theta_0$  vs.  $H_1 : \theta \in \Theta_1$ , which are tested via the randomized test

$$\varphi(x) = \begin{cases} 1, & x \in K_1, \\ \gamma \in (0,1), & x \in K_{01}, \\ 0, & x \in K_0. \end{cases} \quad x = (x_1, \dots, x_n),$$

If  $x \in K_{01}$ , then a random variable  $Y \sim \text{Bernoulli}(\gamma)$  can be used in order to decide whether  $H_0$  is rejected (Y = 1) or not (Y = 0).

#### Definition 3.3.1.

1. The power (or performance) function of a randomized test  $\varphi$  is given by

$$G_n(\theta) = G_n(\varphi, \theta) = \mathbb{E}_{\theta} \varphi(X_1, \dots, X_n), \ \theta \in \Theta.$$

2. The test  $\varphi$  has the significance level  $\alpha \in [0,1]$  if  $G_n(\varphi,\theta) \leq \alpha$ , for all  $\theta \in \Theta_0$ . The number

$$\sup_{\theta \in \Theta_0} G_n(\varphi, \theta)$$

is called *scope* of the test  $\varphi$ . It obviously holds that the scope of an  $\alpha$  confidence level test is smaller than or equal to  $\alpha$ .

- 3. Let  $\Psi(\alpha)$  be the set of all test with confidence level  $\alpha$ . The test  $\varphi_1 \in \Psi(\alpha)$  is called *(uniformly) more powerful* than the test  $\varphi_2 \in \Psi(\alpha)$  if  $G_n(\varphi_1, \theta) \geq G_n(\varphi_2, \theta)$ ,  $\theta \in \Theta_1$ , i.e., if  $\varphi_1$  has a larger power.
- 4. A test  $\varphi^* \in \Psi(\alpha)$  is called *(uniform) most powerful test* in  $\Psi(\alpha)$  if

$$G_n(\varphi^*, \theta) \geq G_n(\varphi, \theta)$$
, for all tests  $\varphi \in \Psi(\alpha)$ ,  $\theta \in \Theta_1$ .

#### Remark 3.3.2.

1. Definition 3.3.1 1. is a generalization of Definition 3.1.4, since for  $\varphi(x) = I(x \in K_1)$ ,

$$G_n(\varphi, \theta) = \mathbb{E}_{\theta} \varphi(X_1, \dots, X_n)$$
$$= P_{\theta} ((X_1, \dots, X_n) \in K_1)$$
$$= P_{\theta} (\text{reject } H_0), \theta \in \Theta$$

holds.

2. A most powerful test  $\varphi^*$  in  $\Psi(\alpha)$  does not always exist. It only exists under certain conditions on  $P_{\theta}$ ,  $\Theta_0$ ,  $\Theta_1$  and  $\Psi(\alpha)$ .

#### 3.3.2 Neyman-Pearson test for simple hypotheses

In this section, simple hypotheses of the form

$$H_0: \theta = \theta_0 \quad \text{vs.} \quad H_1: \theta = \theta_1,$$
 (3.1)

where  $\theta_0, \theta_1 \in \Theta$ ,  $\theta_1 \neq \theta_0$  are considered.

Therefore,  $\Theta_0 = \{\theta_0\}$ ,  $\Theta_1 = \{\theta_1\}$ . Assume that  $F_{\theta_i}$  has a probability density function  $g_i(x)$  with respect to  $\mu$  for i = 0, 1. From now on the notation  $P_0 = P_{\theta_0}$ ,  $P_1 = P_{\theta_1}$ ,  $\mathbb{E}_0 = \mathbb{E}_{\theta_0}$ ,  $\mathbb{E}_1 = \mathbb{E}_{\theta_1}$  will be used. Let  $f_i(x) = \prod_{j=1}^n g_i(x_j)$ ,  $x = (x_1, \dots, x_n)$ , i = 0, 1 be the density of the random sample under  $H_0$  resp.  $H_1$ .

**Definition 3.3.3.** A Neyman-Pearson test (NP test) of simple hypotheses as in (3.1) is given by the rule

$$\varphi(x) = \varphi_K(x) = \begin{cases} 1, & \text{if } f_1(x) > K f_0(x), \\ \gamma, & \text{if } f_1(x) = K f_0(x), \\ 0, & \text{if } f_1(x) < K f_0(x), \end{cases}$$
(3.2)

for constants K > 0 and  $\gamma \in [0, 1]$ .

#### Remark 3.3.4.

- 1. Sometimes K = K(x) and  $\gamma = \gamma(x)$  are seen as functions of x and not as constants.
- 2. The rejection region of the Neyman-Pearson tests  $\varphi_K$  is

$$K_1 = \{x \in B : f_1(x) > K f_0(x)\}.$$

3. The scope of the Neyman-Pearson tests  $\varphi_K$  is given by

$$\mathbb{E}_{0} \varphi_{K}(X_{1}, \dots, X_{n}) = P_{0}(f_{1}(X_{1}, \dots, X_{n}))$$

$$> Kf_{0}(X_{1}, \dots, X_{n}))$$

$$+ \gamma P_{0}(f_{1}(X_{1}, \dots, X_{n})) = Kf_{0}(X_{1}, \dots, X_{n})).$$

4. Definition 3.3.3 can be given equivalently by defining the test statistic

$$T(x) = \begin{cases} \frac{f_1(x)}{f_0(x)}, & x \in B : f_0(x) > 0, \\ \infty, & x \in B : f_0(x) = 0. \end{cases}$$

Then the new test given by

$$\tilde{\varphi}_K(x) = \begin{cases} 1, & \text{if } T(x) > K, \\ \gamma, & \text{if } T(x) = K, \\ 0, & \text{if } T(x) < K, \end{cases}$$

can be introduced, which is for  $P_0$ - and  $P_1$ -almost all  $x \in B$  equivalent to  $\varphi_k$ .  $\varphi_K(x) = \tilde{\varphi}_K(x) \, \forall x \in B \setminus C$  holds, where  $C = \{x \in B : f_0(x) = f_1(x) = 0\}$  has  $P_0$ - resp.  $P_1$ - measure zero.

Using this new formulation, the scope of  $\varphi$  resp.  $\tilde{\varphi}_K$  is given by

$$\mathbb{E}_0 \, \tilde{\varphi}_K = P_0(T(X_1, \dots, X_n) > K) + \gamma \cdot P_0 \left( T(X_1, \dots, X_n) = K \right).$$

**Theorem 3.3.5.** (Optimality theorem) Let  $\varphi_K$  be a Neyman-Pearson test for a K > 0 and  $\gamma \in [0, 1]$ , then  $\varphi_K$  is the most powerful test with confidence level  $\alpha = \mathbb{E}_0 \varphi_K$  of its scope.

**Proof** Let  $\varphi \in \Psi(\alpha)$ , i.e.,  $\mathbb{E}_0(\varphi(X_1,\ldots,X_n)) \leq \alpha$ . In order to show that  $\varphi_K$  is more powerful than  $\varphi$ , it is sufficient to show for simple hypotheses  $H_0$  and  $H_1$ , that  $\mathbb{E}_1 \varphi_K(X_1,\ldots,X_n) \geq \mathbb{E}_1 \varphi(X_1,\ldots,X_n)$ . Define the sets:

$$M^{+} := \{x \in B : \varphi_{K}(x) > \varphi(x)\}$$
  

$$M^{-} := \{x \in B : \varphi_{K}(x) < \varphi(x)\}$$
  

$$M^{=} := \{x \in B : \varphi_{K}(x) = \varphi(x)\}$$

It obviously holds that

$$x \in M^+ \Rightarrow \varphi_K(x) > 0 \Rightarrow f_1(x) \ge K f_0(x),$$
  
 $x \in M^- \Rightarrow \varphi_K(x) < 1 \Rightarrow f_1(x) \le K f_0(x)$   
 $B = M^+ \cup M^- \cup M^=.$ 

Hence

$$\mathbb{E}_{1}(\varphi_{K}(X_{1},\ldots,X_{n})-\varphi(X_{1},\ldots,X_{n})) = \int_{B}(\varphi_{K}(x)-\varphi(x))f_{1}(x)\mu(dx)$$

$$= \left(\int_{M^{+}} + \int_{M^{-}} + \int_{M^{-}}\right)(\varphi_{K}(x)-\varphi(x))f_{1}(x)\mu(dx)$$

$$\geq \int_{M^{+}}(\varphi_{K}(x)-\varphi(x))Kf_{0}(x)\mu(dx)$$

$$+ \int_{M^{-}}(\varphi_{K}(x)-\varphi(x))Kf_{0}(x)\mu(dx)$$

$$= \int_{B}(\varphi_{K}(x)-\varphi(x))Kf_{0}(x)\mu(dx)$$

$$= K\left(\mathbb{E}_{0}\varphi_{K}(X_{1},\ldots,X_{n})-\mathbb{E}_{0}\varphi(X_{1},\ldots,X_{n})\right)$$

$$\geq K(\alpha-\alpha) = 0,$$

since both tests obtain confidence level  $\alpha$ .

### Remark 3.3.6.

- 1. As  $\gamma$  does not appear in the proof, the same result holds for  $\gamma(x) \neq$  const.
- 2. The proof implies the inequality given by

$$\int_{B} \left( \varphi_K(x) - \varphi(x) \right) \left( f_1(x) - K f_0(x) \right) \mu(dx) \ge 0$$

if K is constant, resp.

$$\mathbb{E}_1 \left( \varphi_K(X_1, \dots, X_n) - \varphi(X_1, \dots, X_n) \right) \ge \int_B \left( \varphi_K(x) - \varphi(x) \right) K(x) f_0(x) \mu(dx)$$

in the general case.

# Theorem 3.3.7. (Fundamental lemma of Neyman-Pearson)

- 1. For an arbitrary  $\alpha \in (0,1)$ , there exists a Neyman-Pearson test  $\varphi_K$  with scope  $\alpha$ , which is by Theorem 3.3.5 the most powerful  $\alpha$  confidence level test.
- 2. If  $\varphi$  is also a most powerful test with confidence level  $\alpha$ , then  $\varphi(x) = \varphi_K(x)$  for  $\mu$ -almost all  $x \in K_0 \cup K_1 = \{x \in B : f_1(x) \neq K f_0(x)\}$  and  $\varphi_K$  of part 1.

**Proof** 1. For  $\varphi_K(x)$  it holds that

$$\varphi_K(x) = \begin{cases} 1, & \text{if } x \in K_1 = \{x : f_1(x) > K \cdot f_0(x)\}, \\ \gamma, & \text{if } x \in K_{01} = \{x : f_1(x) = K \cdot f_0(x)\}, \\ 0, & \text{if } x \in K_0 = \{x : f_1(x) < K \cdot f_0(x)\}. \end{cases}$$

The scope of  $\varphi_K$  is given by

$$P_0(T(X_1,...,X_n) > K) + \gamma P_0(T(X_1,...,X_n) = K) = \alpha,$$
 (3.3)

where

$$T(x_1, ..., x_n) = \begin{cases} \frac{f_1(x_1, ..., x_n)}{f_0(x_1, ..., x_n)}, & \text{if } f_0(x_1, ..., x_n) > 0, \\ \infty, & \text{otherwise.} \end{cases}$$

The goal is to find a K > 0 and a  $\gamma \in [0,1]$ , such that equation (3.3) holds. Let  $\tilde{F}_0(x) = P_0(T(X_1, \dots, X_n) \leq x)$ ,  $x \in \mathbb{R}$  be the distribution function of T. Since  $T \geq 0$ , it holds that  $\tilde{F}_0(x) = 0$ , if x < 0. Furthermore,  $P_0(T(X_1, \dots, X_n) < \infty) = 1$ , which means that  $\tilde{F}^{-1}(\alpha) \in [0, \infty)$ ,  $\alpha \in (0, 1)$ . Equation (3.3) can then be rewritten as

$$1 - \tilde{F}_0(K) + \gamma \left( \tilde{F}_0(K) - \tilde{F}_0(K) - \right) = \alpha, \tag{3.4}$$

where  $\tilde{F}_0(K-) = \lim_{x \to K} \tilde{F}_0(x)$ .

Let 
$$K = \tilde{F}_0^{-1}(1 - \alpha)$$
, then:

- (a) If K is a point of continuity of  $\tilde{F}_0$ , then Equation (3.4) is satisfied for all  $\gamma \in [0, 1]$ , for example for  $\gamma = 0$ .
- (b) If K is not a point of continuity of  $\tilde{F}_0$ , then  $\tilde{F}_0(K) \tilde{F}_0(K-) > 0$ , which implies that

$$\gamma = \frac{\alpha - 1 + \tilde{F}_0(K)}{\tilde{F}_0(K) - \tilde{F}_0(K-)}.$$

Therefore, a Neyman-Pearson test with confidence level  $\alpha$  exists.

2. Define  $M^{\neq} := \{x \in B : \varphi(x) \neq \varphi_K(x)\}$ . It has to be shown that

$$\mu\left((K_0 \cup K_1) \cap M^{\neq}\right) = 0.$$

Consider

$$\mathbb{E}_1 \varphi(X_1, \dots, X_n) - \mathbb{E}_1 \varphi_K(X_1, \dots, X_n) = 0 \quad (\varphi \text{ and } \varphi_K \text{ are most powerful tests})$$

$$\mathbb{E}_0 \varphi(X_1, \dots, X_n) - \mathbb{E}_0 \varphi_K(X_1, \dots, X_n) \leq 0 \quad (\varphi \text{ and } \varphi_K \text{ are } \alpha\text{-tests}$$
with scope  $\varphi_K = \alpha$ )

$$\Rightarrow \int_{B} (\varphi - \varphi_K) \cdot (f_1 - K \cdot f_0) \, d\mu \ge 0.$$

In Remark 3.3.6 it has been shown, that

$$\int_{B} (\varphi - \varphi_{K})(f_{1} - K \cdot f_{0})d\mu \leq 0$$

$$\Rightarrow \int_{B} (\varphi - \varphi_{K})(f_{1} - K \cdot f_{0})d\mu = 0 = \int_{M} (\varphi - \varphi_{K})(f_{1} - K \cdot f_{0})d\mu.$$

$$M^{\neq} \cap (K_{0} \cup K_{1})$$

 $\mu(M^{\neq} \cap (K_0 \cup K_1)) = 0$  holds if the integrand  $(\varphi - \varphi_K)(f_1 - K \cdot f_0) > 0$  on  $(K_0 \cup K_1) \cap M^{\neq}$ . We need to show that

$$(\varphi_K - \varphi)(f_1 - Kf_0) > 0 \text{ für } x \in (K_0 \cup K_1) \cap M^{\neq}. \tag{3.5}$$

Now,

$$f_1 - K f_0 > 0 \Rightarrow \varphi_K - \varphi > 0,$$
  
 $f_1 - K f_0 < 0 \Rightarrow \varphi_K - \varphi < 0,$ 

holds, since

$$f_1(x) > K f_0(x) \Rightarrow \varphi_K(x) = 1$$
  
and with  $\varphi(x) < 1$  we get  $\varphi_K(x) - \varphi(x) > 0$  on  $M^{\neq}$ .  
 $f_1(x) < K f_0(x) \Rightarrow \varphi_K(x) = 0$   
and with  $\varphi(x) > 0$  we get  $\varphi_K(x) - \varphi(x) < 0$  on  $M^{\neq}$ .

Thus inequality (3.5) holds and finally

$$\mu\left((K_0 \cup K_1) \cap M^{\neq}\right) = 0.$$

**Remark 3.3.8.** If  $\varphi$  and  $\varphi_K$  are most powerful  $\alpha$ -tests, then they are  $P_0$ -resp.  $P_1$ - almost surely equal.

# Example 3.3.9. (Neyman-Pearson test for the parameter of the Poisson distribution)

Let  $(X_1, \ldots, X_n)$  be a random sample with  $X_i \sim \text{Poisson}(\lambda)$ ,  $\lambda > 0$ , where  $X_i$  are i.i.d. for  $i = 1, \ldots, n$ . The hypotheses  $H_0: \lambda = \lambda_0$  vs.  $H_1: \lambda = \lambda_1$  need to be tested. Here

$$g_i(x) = e^{-\lambda_i} \frac{\lambda_i^x}{x!}, x \in \mathbb{N}_0, i = 0, 1,$$

$$f_i(x) = f_i(x_1, \dots, x_n) = \prod_{j=1}^n g_i(x_j) = \prod_{j=1}^n e^{-\lambda_i} \frac{\lambda_i^{x_j}}{x_j!} = e^{-n\lambda_i} \cdot \frac{\lambda_i^{\sum_{j=1}^n x_j}}{(x_1! \cdot \dots \cdot x_n!)}, i = 0, 1.$$

The Neyman-Pearson test statistic is given by

$$T(x_1,\ldots,x_n) = \begin{cases} \frac{f_1(x)}{f_0(x)} = e^{-n(\lambda_1 - \lambda_0)} \cdot (\lambda_1/\lambda_0)^{\sum_{j=1}^n x_j}, & \text{if } x_1,\ldots,x_n \in \mathbb{N}_0, \\ \infty, & \text{otherwise.} \end{cases}$$

The Neyman-Pearson decision rule is given by

$$\varphi_K(x_1, \dots, x_n) = \begin{cases} 1, & \text{if } T(x_1, \dots, x_n) > K, \\ \gamma, & \text{if } T(x_1, \dots, x_n) = K, \\ 0, & \text{if } T(x_1, \dots, x_n) < K. \end{cases}$$

Choose  $K > 0, \gamma \in [0, 1]$ , such that  $\varphi_K$  has scope  $\alpha$ . In order to do so, solve

$$\alpha = P_0(T(X_1, \dots, X_n) > K) + \gamma P_0(T(X_1, \dots, X_n) = K)$$

for  $\gamma$  and K.

$$P_0(T(X_1, \dots, X_n) > K) = P_0(\log T(X_1, \dots, X_n) > \log K)$$

$$= P_0\left(-n(\lambda_1 - \lambda_0) + \log\left(\frac{\lambda_1}{\lambda_0}\right) \sum_{j=1}^n X_j > \log K\right)$$

$$= P_0\left(\sum_{j=1}^n X_j > A_K\right),$$

where

$$A_K := \left\lfloor \frac{\log K + n \cdot (\lambda_1 - \lambda_0)}{\log \frac{\lambda_1}{\lambda_0}} \right\rfloor,\,$$

if for example  $\lambda_1 > \lambda_0$ . If  $\lambda_1 < \lambda_0$  then replace > with < in the argument above

Due to the stability of the Poisson distribution

$$\sum_{j=1}^{n} X_j \sim \text{Poisson}(n\lambda_0),$$

holds under  $H_0$ . Thus choose K as the smallest nonnegative number with

$$P_0\left(\sum_{j=1}^n X_j > A_K\right) \le \alpha,$$

and set

$$\gamma = \frac{\alpha - P_0(\sum_{j=1}^n X_j > A_K)}{P_0(\sum_{j=1}^n X_j = A_K)},$$

where

$$P_0\left(\sum_{j=1}^n X_j > A_K\right) = 1 - \sum_{j=0}^{A_K} e^{-\lambda_0 n} \frac{(\lambda_0 n)^j}{j!},$$

$$P_0\left(\sum_{j=1}^n X_j = A_K\right) = e^{-\lambda_0 n} \frac{(\lambda_0 n)^{A_K}}{A_K!}.$$

Hence, the parameters K and  $\gamma$  have been found and a Neyman-Pearson test  $\varphi_K$  has been constructed.

### 3.3.3 One-sided Neyman-Pearson tests

So far Neyman-Pearson tests for simple hypotheses like  $H_i: \theta = \theta_i, i = 0, 1$  have been considered. This section aims to introduce one-sided Neyman-Pearson tests for hypotheses of the form  $H_0: \theta \leq \theta_0$  vs.  $H_1: \theta > \theta_0$ . In an initial step, a test for the following hypotheses is constructed: Let  $(X_1, \ldots, X_n)$  be a random sample,  $X_i$  i.i.d. with

$$X_i \sim F_\theta \in \Lambda = \{F_\theta : \theta \in \Theta\},\$$

where  $\Theta \subset \mathbb{R}$  is open and  $\Lambda$  uniquely parameterized, i.e.,

$$\theta \neq \theta' \Rightarrow F_{\theta} \neq F_{\theta'}$$
.

Furthermore, assume that  $F_{\theta}$  has a density  $g_{\theta}$  w.r.t. the Lebesgue measure (resp. counting measure)  $\mu$  on  $\mathbb{R}$  (resp.  $\mathbb{N}_0$ ). Then

$$f_{\theta}(x) = \prod_{j=1}^{n} g_{\theta}(x_j), \quad x = (x_1, \dots, x_n)$$

is a density of  $(X_1, \ldots, X_n)$  with respect to  $\mu^n$  on B.

**Definition 3.3.10.** A distribution on B with density  $f_{\theta}$  is a member of the class of distributions with monotone likelihood ratio in T, if for all  $\theta < \theta'$  exist a monotonically increasing function  $h : \mathbb{R} \times \Theta^2 \to \mathbb{R} \cup \infty$  on  $\mathbb{R}$  and a statistic  $T : B \to \mathbb{R}$  with the property

$$\frac{f_{\theta'}(x)}{f_{\theta}(x)} = h(T(x), \theta, \theta'),$$

where

$$h(T(x), \theta, \theta') = \infty$$
, for all  $x \in B : f_{\theta}(x) = 0$  and  $f_{\theta'}(x) > 0$ .

The case  $f_{\theta}(x) = f_{\theta'}(x) = 0$  occurs with probability  $P_{\Theta}$ - resp.  $P_{\Theta'}$  zero.

**Definition 3.3.11.** Let  $Q_{\theta}$  be a distribution on  $(B, \mathcal{B})$  with probability density function  $f_{\theta}$  w.r.t.  $\mu$ .  $Q_{\theta}$  is an element of the *one-parametric exponential* family  $(\theta \in \Theta \subset \mathbb{R} \text{ open})$ , if the density is given by:

$$f_{\theta}(x) = \exp\left\{c(\theta) \cdot T(x) + a(\theta)\right\} \cdot l(x), \quad x = (x_1, \dots, x_n) \in B,$$

where  $c(\theta)$  is a monotonically increasing function and  $\operatorname{Var}_{\theta} T(X_1, \dots, X_n) > 0, \ \theta \in \Theta$ .

**Lemma 3.3.12.** Distributions of the one-parametric exponential family have a monotone likelihood ratio.

**Proof** Let  $Q_{\theta}$  be in the one-parametric exponential family with probability density function

$$f_{\theta}(x) = \exp \{c(\theta) \cdot T(x) + a(\theta)\} \cdot l(x).$$

For  $\theta < \theta'$ 

$$\frac{f_{\theta'}(x)}{f_{\theta}(x)} = \exp\left\{ (c(\theta') - c(\theta)) \cdot T(x) + a(\theta') - a(\theta) \right\}$$

is monotone with respect to T, since  $c(\theta') - c(\theta) > 0$  because of the monotonicity of  $c(\theta)$ . Thus  $f_{\theta}$  has a monotone likelihood ratio.

# Example 3.3.13.

# 1. Normal distributed random samples

Let  $X_i \sim \mathcal{N}(\mu, \sigma_0^2)$ , i = 1, ..., n be i.i.d. random variables, with unknown  $\mu \in \mathbb{R}$  and known variance  $\sigma_0^2 > 0$  (here  $\mu$  denotes the expected value of  $X_i$  and not a measure on  $\mathbb{R}$ ). The probability density function of the vector  $X = (X_1, ..., X_n)^{\top}$  is given by

$$\begin{split} f_{\mu}(x) &= \prod_{i=1}^{n} g_{\mu}(x_{i}) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma_{0}^{2}}} e^{-\frac{(x_{i}-\mu)^{2}}{2\sigma_{0}^{2}}} \\ &= \frac{1}{(2\pi\sigma_{0}^{2})^{n/2}} \exp\left\{-\frac{1}{2\sigma_{0}^{2}} \sum_{i=1}^{n} (x_{i}-\mu)^{2}\right\} \\ &= \frac{1}{(2\pi\sigma_{0}^{2})^{n/2}} \exp\left\{-\frac{1}{2\sigma_{0}^{2}} \left(\sum_{i=1}^{n} x_{i}^{2} - 2\mu \sum_{i=1}^{n} x_{i} + \mu^{2}n\right)\right\} \\ &= \exp\left(\underbrace{\frac{\mu}{\sigma_{0}^{2}} \cdot \sum_{i=1}^{n} x_{i}}_{c(\mu)} \cdot \underbrace{\frac{1}{2\sigma_{0}^{2}}}\right) \cdot \underbrace{\frac{1}{(2\pi\sigma_{0}^{2})^{n/2}} \exp\left(-\frac{\sum_{i=1}^{n} x_{i}^{2}}{2\sigma_{0}^{2}}\right)}_{l(x)}. \end{split}$$

Thus  $\mathcal{N}(\mu, \sigma_0^2)$  is a member of the one-parametric exponential family with  $c(\mu) = \frac{\mu}{\sigma_0^2}$  and  $T(x) = \sum_{i=1}^n x_i$ .

#### 2. Binomial distributed random samples

Let  $X_i \sim \text{Bin}(k, p)$  be i.i.d., i = 1, ..., n. The parameter  $p \in (0, 1)$  is assumed to be unknown. The probability mass function of  $X = (X_1, ..., X_n)^{\top}$  is given by

$$f_{p}(x) = P_{p}(X_{i} = x_{i}, i = 1, ..., n)$$

$$= \prod_{i=1}^{n} {n \choose x_{i}} p^{x_{i}} (1-p)^{k-x_{i}} = p^{\sum_{i=1}^{n} x_{i}} \cdot \frac{(1-p)^{nk}}{(1-p)^{i-1}} \cdot \prod_{i=1}^{n} {k \choose x_{i}}$$

$$= \exp\left\{\left(\sum_{i=1}^{n} x_{i}\right) \cdot \underbrace{\log\left(\frac{p}{1-p}\right)}_{c(p)} + \underbrace{nk \cdot \log(1-p)}_{a(p)}\right\} \cdot \prod_{i=1}^{n} {k \choose x_{i}},$$

thus Bin(n, p) is a member of the one-parametric exponential family with

$$c(p) = \log\left(\frac{p}{1-p}\right)$$

and

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

**Lemma 3.3.14.** If  $\varphi_K$  is the Neyman-Pearson test of the hypotheses  $H_0$ :  $\theta = \theta_0$  vs.  $H_1: \theta = \theta_1$ , then

$$\mu^n(K_0 \cup K_1) = \mu^n(\{x \in B : f_1(x) \neq K f_0(x)\}) > 0.$$

**Proof** Since  $\theta_0 \neq \theta_1$  and because of the unique parametrization it holds that  $f_0 \neq f_1$  on a set with  $\mu$ -measure greater than 0.

Assume that  $\mu(K_0 \cup K_1) = 0$ . Then  $f_1(x) = K \cdot f_0(x)$   $\mu$ -almost surely, which means that

$$1 = \int_{B} f_1(x)dx = K \cdot \int_{B} f_0(x)dx.$$

This yields K = 1 and  $f_1(x) = f_0(x)$   $\mu$ -almost surely, which is a contradiction to the unique parametrization.

Assume that  $(X_1, \ldots, X_n)$  is an i.i.d. random sample, where  $X_i$  have the density  $g_{\theta}$ ,  $i = 1, \ldots, n$  and  $(X_1, \ldots, X_n)$  has the density  $f_{\theta}(x) = \prod_{i=1}^n g_{\theta}(x_i)$  from the class of distributions with monotone likelihood ratio and a statistic  $T(x_1, \ldots, x_n)$ .

Consider the hypotheses  $H_0: \theta \leq \theta_0$  vs.  $H_1: \theta > \theta_0$  and the Neyman-Pearson test:

$$\varphi_{K^*}^*(x) = \begin{cases} 1, & \text{if } T(x) > K^*, \\ \gamma^*, & \text{if } T(x) = K^*, \\ 0, & \text{if } T(x) < K^* \end{cases}$$
(3.6)

for  $K^* \in \mathbb{R}$  and  $\gamma^* \in [0,1]$ . The power function of  $\varphi_{K^*}^*$  at  $\theta_0$  is given by

$$G_n(\theta_0) = \mathbb{E}_0 \, \varphi_{K^*}^* = P_0 \left( T(X_1, \dots, X_n) > K^* \right) + \gamma^* \cdot P_0 \left( T(X_1, \dots, X_n) = K^* \right).$$

#### Theorem 3.3.15.

- 1. If  $\alpha = \mathbb{E}_0 \varphi_{K^*}^*(X_1, \dots, X_n) > 0$ , then the defined test is a most powerful test of the one-sided hypotheses  $H_0$  vs.  $H_1$  with confidence level  $\alpha$ .
- 2. For every confidence level  $\alpha \in (0,1)$  exists a  $K^* \in \mathbb{R}$  and  $\gamma^* \in [0,1]$ , such that  $\varphi_{K^*}^*$  is a most powerful test with scope  $\alpha$ .
- 3. The power function  $G_n(\theta)$  of  $\varphi_{K^*}^*(\theta)$  is monotonically nondecreasing in  $\theta$ . If  $0 < G_n(\theta) < 1$ , then  $G_n$  is even monotonically increasing.

#### Proof

1. Assume  $\theta_1 > \theta_0$  and consider the simple hypotheses  $H_0': \theta = \theta_0$  and  $H_1': \theta = \theta_1$ . Let

$$\varphi_K(x) = \begin{cases} 1, & f_1(x) > K f_0(x), \\ \gamma, & f_1(x) = K f_0(x), \\ 0, & f_1(x) < K f_0(x), \end{cases}$$

be the Neyman-Pearson test for  $H'_0, H'_1$  with K > 0. Since  $f_\theta$  has the monotone likelihood ratio with statistic T, i.e.,

$$\frac{f_1(x)}{f_0(x)} = h(T(x), \theta_0, \theta_1),$$

there exists K > 0, such that

$$\left\{ x : f_1(x)/f_0(x) > K \\ < K \right\} \subset \left\{ T(x) > K^* \\ < K^* \right\} \quad \text{with } K = h(K^*, \theta_0, \theta_1).$$

 $\varphi_K$  is a most powerful Neyman-Pearson test with confidence level  $\alpha = \mathbb{E}_0 \, \varphi_K = \mathbb{E}_0 \, \varphi_{K^*}^*$ .

 $\alpha > 0$  implies that  $K < \infty$ , since  $K = \infty$  would yield

$$0 < \alpha = \mathbb{E}_0 \, \varphi_K \le P_0 \, (T(X_1, \dots, X_n) \ge K^*)$$

$$\le P_0 \, \left( \frac{f_1(X_1, \dots, X_n)}{f_0(X_1, \dots, X_n)} = \infty \right)$$

$$= P_0 \, (f_1(X_1, \dots, X_n) > 0, f_0(X_1, \dots, X_n) = 0)$$

$$= \int_B I \, (f_1(x) > 0, f_0(x) = 0) \cdot f_0(x) \mu(dx)$$

$$= 0$$

For the test  $\varphi_{K^*}^*$  in (3.6) it holds that

$$\varphi_{K^*}^*(x) = \begin{cases} 1, & \text{if } f_1(x)/f_0(x) > K, \\ \gamma^*(x), & \text{if } f_1(x)/f_0(x) = K, \\ 0, & \text{if } f_1(x)/f_0(x) < K, \end{cases}$$

where  $\gamma^*(x) \in \{\gamma^*, 0, 1\}$ . Thus  $\varphi_{K^*}^*$  is a most powerful Neyman-Pearson test for  $H_0'$  vs.  $H_1'$  (cf. Remark 3.3.4, 1. and Remark 3.3.6) for an arbitrary  $\theta_1 > \theta_0$ . That is why  $\varphi_{K^*}^*$  is a most powerful Neyman-Pearson test for  $H_0''$ :  $\theta = \theta_0$  vs.  $H_1''$ :  $\theta > \theta_0$ .

The same assertion is obtained by part 3. of the theorem for  $H_0: \theta \le \theta_0$  vs.  $H_1: \theta > \theta_0$ , since then  $G_n(\theta) \le G_n(\theta_0) = \alpha$  for all  $\theta < \theta_0$ .

- 2. See proof of Theorem 3.3.7, 1.).
- 3. It has to be shown that  $G_n(\theta)$  is monotone. In order to do so, let  $\theta_1 < \theta_2$  and show that  $\alpha_1 = G_n(\theta_1) \le G_n(\theta_2)$ . Consider the new, simple hypotheses  $H_0'': \theta = \theta_1$  vs.  $H_1'': \theta = \theta_2$ . The test  $\varphi_{K^*}^*$  can similarly to 1. be stated as a Neyman-Pearson test (for the hypotheses  $H_0''$  and  $H_1''$ ), which is a most powerful test with confidence level  $\alpha_1$ . Consider another constant test  $\varphi(x) = \alpha_1$ . Then  $\alpha_1 = \mathbb{E}_{\theta_2} \varphi \le \mathbb{E}_{\theta_2} \varphi_{K^*}^* = G_n(\theta_2)$ . This implies that  $G_n(\theta_1) \le G_n(\theta_2)$ .

It is now to be shown that for  $G_n(\theta) \in (0,1)$  it holds that  $G_n(\theta_1) < G_n(\theta_2)$ . Assume that  $\alpha_1 = G_n(\theta_1) = G_n(\theta_2)$  and  $\theta_1 < \theta_2$  for  $\alpha_1 \in (0,1)$ . Then  $\varphi(x) = \alpha_1$  is also a most powerful test for  $H_0''$  and  $H_1''$ . Theorem 3.3.7, 2. implies

$$\mu^{n}(\{x \in K_{0} \cup K_{1} : \underbrace{\varphi(x)}_{=\alpha_{1}} \neq \varphi_{K^{*}}^{*}(x)\}) = 0$$

which is a contradiction to the construction of the test  $\varphi_{K^*}^*$ . This test can not be equal to  $\alpha_1 \in (0,1)$  on  $K_0 \cup K_1$ .

#### Remark 3.3.16.

1. Theorem 3.3.15 can be applied to the Neyman-Pearson tests of the one-sided hypotheses

$$H_0: \theta \geq \theta_0$$
 vs.  $H_1: \theta < \theta_0$ ,

with the corresponding difference

$$\theta \mapsto -\theta$$
$$T \mapsto -T.$$

Thus the most powerful  $\alpha$  test also exists in that case.

2. It can be shown that the power function  $G_n(\varphi_{K^*}^*, \theta)$  of the most powerful Neyman-Pearson tests on  $\Theta_0 = (-\infty, \theta_0)$  attains the following minimization property:

$$G_n(\varphi_{K^*}^*, \theta) \le G_n(\varphi, \theta) \quad \forall \varphi \in \Psi(\alpha), \ \theta \le \theta_0.$$

**Example 3.3.17.** Consider a normally distributed random sample  $(X_1, \ldots, X_n)$  of i.i.d. random variables  $X_i$  with  $X_i \sim \mathcal{N}(\mu, \sigma_0^2)$  and known  $\sigma_0^2 > 0$ . The hypotheses

$$H_0: \mu \leq \mu_0 \text{ vs. } H_1: \mu > \mu_0$$

are tested. Example 3.1.7 provides the test statistic

$$T(X_1,\ldots,X_n)=\sqrt{n}\frac{\overline{X}_n-\mu_0}{\sigma_0},$$

where under  $H_0$  it holds that  $T(X_1, \ldots, X_n) \sim \mathcal{N}(0, 1)$ .  $H_0$  is rejected, if

$$T(X_1,\ldots,X_n) > z_{1-\alpha}, \text{ with } \alpha \in (0,1).$$

It will be shown that this test is the most powerful Neyman-Pearson test with confidence level  $\alpha$ . Example 3.3.13 implies that the probability density function  $f_n$  of  $(X_1, \ldots, X_n)$  is a member of the one-parametric exponential family with

$$\tilde{T}(X_1,\ldots,X_n) = \sum_{i=1}^n X_i.$$

Then  $f_{\mu}$  of  $(x_1, \ldots, x_n)$  is also a member of the one-parametric exponential family with respect to the statistic

$$T(X_1,\ldots,X_n)=\sqrt{n}\frac{\overline{X}_n-\mu}{\sigma_0},$$

since it holds that

$$f_{\mu}(x) = \exp\left(\underbrace{\frac{\mu}{\sigma_0^2}}_{\tilde{c}(\mu)} \cdot \underbrace{\sum_{i=1}^n x_i}_{\tilde{T}} \underbrace{-\frac{\mu^2 n}{2\sigma_0^2}}_{\tilde{a}(\mu)}\right) \cdot l(x)$$
$$= \exp\left(\underbrace{\frac{\mu\sqrt{n}}{\sigma_0}}_{c(\mu)} \cdot \underbrace{\sqrt{n}\frac{\overline{x}_n - \mu}{\sigma_0}}_{T} + \underbrace{\frac{\mu^2 n}{2\sigma_0^2}}_{a(\mu)}\right) \cdot l(x).$$

The statistic T can be used in the construction of the Neyman-Pearson tests (cf. Equation (3.6)):

$$\varphi_{K^*}^*(x) = \begin{cases} 1, & \text{if } T(x) > z_{1-\alpha}, \\ 0, & \text{if } T(x) = z_{1-\alpha}, \\ 0, & \text{if } T(x) < z_{1-\alpha} \end{cases}$$

(with  $K^* = z_{1-\alpha}$  and  $\gamma^* = 0$ ). Theorem 3.3.15 implies that this test is the most powerful Neyman-Pearson test with confidence level  $\alpha$  for our hypotheses:

$$G_n(\varphi_{K^*}, \mu_0) = P_0(T(X_1, \dots, X_n) > z_{1-\alpha}) + 0 \cdot P_0(T(X_1, \dots, X_n) \le z_{1-\alpha})$$
  
=  $1 - \Phi(z_{1-\alpha}) = 1 - (1 - \alpha) = \alpha$ .

#### 3.3.4 Unbiased two-sided tests

Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables with probability density function

$$f_{\theta}(x) = \prod_{i=1}^{n} g_{\theta}(x_i).$$

In the following, a two-sided test of the hypotheses

$$H_0: \theta = \theta_0 \text{ vs. } H_1: \theta \neq \theta_0$$

is considered. There cannot be a most powerful test  $\varphi$  with confidence level  $\alpha$  for all  $\alpha \in (0,1)$ . Assume that  $\varphi$  is the most powerful with confidence level  $\alpha$  for  $H_0$  vs.  $H_1$ , then  $\varphi$  would be the most powerful test for the hypotheses

- 1.  $H'_0: \theta = \theta_0 \text{ vs. } H'_1: \theta > \theta_0$
- 2.  $H_0'': \theta = \theta_0 \text{ vs. } H_1'': \theta < \theta_0.$

By Theorem 3.3.15, 3. the power function would then be given by

- 1.  $G_n(\varphi, \theta) < \alpha \text{ on } \theta < \theta_0, \text{ resp.}$
- 2.  $G_n(\varphi, \theta) > \alpha$  on  $\theta < \theta_0$ ,

which is a contradiction!

That is why the class of all possible tests is reduced to the class of unbiased tests (cf. Definition 3.1.12). The test  $\varphi$  is unbiased if and only if

$$G_n(\varphi, \theta) \leq \alpha \text{ for } \theta \in \Theta_0 \text{ and } G_n(\varphi, \theta) \geq \alpha \text{ for } \theta \in \Theta_1.$$

# Example 3.3.18.

- 1.  $\varphi(x) \equiv \alpha$  is unbiased.
- 2. The two-sided Gauss test is unbiased, (cf. Example 3.1.7):  $G_n(\varphi, \mu) \ge \alpha$  for all  $\mu \in \mathbb{R}$ .

Assume that  $X_i$  are i.i.d. The probability density function  $f_{\theta}$  of  $(X_1, \dots, X_n)$  is assumed to be a member of the one-parametric exponential family

$$f_{\theta}(x) = \exp\left\{c(\theta) \cdot T(x) + a(\theta)\right\} \cdot l(x),\tag{3.7}$$

where  $c(\theta)$  and  $a(\theta)$  are continuously differentiable on  $\Theta$  with

$$c'(\theta) > 0$$
 and  $\operatorname{Var}_{\theta} T(X_1, \dots, X_n) > 0$ 

for all  $\theta \in \Theta$ . Let  $f_{\theta}(x)$  be continuous in  $(x, \theta)$  on  $B \times \Theta$ .

Exercise 3.3.19. Show that the following relation holds:

$$a'(\theta) = -c'(\theta)\mathbb{E}_{\theta} T(X_1, \dots, X_n).$$

**Lemma 3.3.20.** Let  $\varphi$  be an unbiased test with confidence level  $\alpha$  for

$$H_0: \theta = \theta_0 \text{ vs. } H_1: \theta \neq \theta_0.$$

Then

1. 
$$\alpha = \mathbb{E}_0 \varphi(X_1, \dots, X_n) = G_n(\varphi, \theta_0),$$

2. 
$$\mathbb{E}_0 \left[ T(X_1, \dots, X_n) \varphi(X_1, \dots, X_n) \right] = \alpha \cdot \mathbb{E}_0 T(X_1, \dots, X_n),$$

#### Proof

1. The power function of  $\varphi$  is given by

$$G_n(\varphi,\theta) = \int_{\mathbb{R}} \varphi(x) f_{\theta}(x) \mu n(dx).$$

Since  $f_{\theta}$  is in the one-parametric exponential family,  $G_n(\varphi, \theta)$  is differentiable (under the integral) with respect to  $\theta$  and hence continuous in  $\theta$ . Since  $\varphi$  is unbiased, it holds that

$$G_n(\varphi, \theta_0) \le \alpha, \ G_n(\varphi, \theta) \ge \alpha, \quad \theta \ne \theta_0$$

Thus  $G_n(\varphi, \theta_0) = \alpha$  and  $\theta_0$  minimizes  $G_n$ , which proves 1.

2. Since  $\theta_0$  minimizes  $G_n$ , it holds that

$$0 = G'_n(\varphi, \theta_0) = \int_B \varphi(x) (c'(\theta_0) T(x) + a'(\theta_0)) f_0(x) \mu(dx)$$

$$= c'(\theta_0) \cdot \mathbb{E}_0 \left[ \varphi(X_1, \dots, X_n) T(X_1, \dots, X_n) \right] + a'(\theta_0) \cdot G_n(\varphi, \theta_0)$$

$$= c'(\theta_0) \cdot \mathbb{E}_0 \left[ \varphi(X_1, \dots, X_n) T(X_1, \dots, X_n) \right] + \alpha a'(\theta_0)$$
(Exerc. 3.3.19)
$$= c'(\theta_0) \left( \mathbb{E}_0 \left( \varphi \cdot T \right) - \alpha \mathbb{E}_0 T \right)$$

Therefore,  $\mathbb{E}_0(\varphi T) = \alpha \mathbb{E}_0 T$ .

In the following paragraph, a modification of the Neyman-Pearson test for simple hypotheses of the form

$$H_0: \theta = \theta_0 \text{ vs. } H_1': \theta = \theta_1, \quad \theta_1 \neq \theta_0,$$

is introduced. For  $\lambda, K \in \mathbb{R}, \gamma : B \to [0, 1]$ , define

$$\varphi_{K,\lambda}(x) = \begin{cases} 1, & \text{if } f_1(x) > (K + \lambda T(x)) f_0(x), \\ \gamma(x), & \text{if } f_1(x) = (K + \lambda T(x)) f_0(x), \\ 0, & \text{if } f_1(x) < (K + \lambda T(x)) f_0(x), \end{cases}$$
(3.8)

where T(x) is the statistic in (3.7).

Let  $\tilde{\Psi}(\alpha)$  be the class of all tests that satisfy the conditions 1. and 2. of Lemma 3.3.20. Lemma 3.3.20 implies that the set of unbiased tests with confidence level  $\alpha$  is a subset of  $\tilde{\Psi}(\alpha)$ .

**Theorem 3.3.21.** The modified Neyman-Pearson test  $\varphi_{K,\lambda}$  is the most powerful  $\alpha$  test in  $\tilde{\Psi}(\alpha)$  for the hypotheses  $H_0$  vs.  $H'_1$  with confidence level  $\alpha = \mathbb{E}_0 \varphi_{K,\lambda}$ , if  $\varphi_{K,\lambda} \in \tilde{\Psi}(\alpha)$ .

**Proof** It has to be shown that  $\mathbb{E}_1 \varphi_{K,\lambda} \geq \mathbb{E}_1 \varphi$  for all  $\varphi \in \tilde{\Psi}(\alpha)$ , resp.  $\mathbb{E}_1 (\varphi_{K,\lambda} - \varphi) \geq 0$ . It holds that

$$\mathbb{E}_{1}\left(\varphi_{K,\lambda}-\varphi\right) = \int_{B} (\varphi_{K,\lambda}(x)-\varphi(x))f_{1}(x)\mu(dx)$$

$$\stackrel{(\text{Rem. 3.3.6, 2.}))}{\geq} \int_{B} (\varphi_{K,\lambda}(x)-\varphi(x))(K+\lambda T(x))f_{0}(x)\mu(dx)$$

$$= K\left(\underbrace{\mathbb{E}_{0}\,\varphi_{K,\lambda}}_{=\alpha}-\underbrace{\mathbb{E}_{0}\,\varphi}_{=\alpha}\right) + \lambda\left(\underbrace{\mathbb{E}_{0}\,(\varphi_{K,\lambda}\cdot T)}_{\alpha\mathbb{E}_{0}\,T}-\underbrace{\mathbb{E}_{0}\,(\varphi\cdot T)}_{=\alpha\cdot\mathbb{E}_{0}\,T}\right)$$

$$= 0.$$

since  $\varphi, \varphi_{K,\lambda} \in \tilde{\Psi}(\alpha)$ .

Consider the following decision rule, which will later be used in testing twosided hypotheses given by

$$H_0: \theta = \theta_0 \text{ vs. } H_1: \theta \neq \theta_0,$$

$$\varphi_c(x) = \begin{cases}
1, & \text{if } T(x) \notin (c_1, c_2), \\
\gamma_1, & \text{if } T(x) = c_1, \\
\gamma_2, & \text{if } T(x) = c_2, \\
0, & \text{if } T(x) \in (c_1, c_2),
\end{cases}$$
(3.9)

for  $c_1 \leq c_2 \in \mathbb{R}$ ,  $\gamma_1, \gamma_2 \in [0, 1]$  and the statistic T(x),  $x = (x_1, \ldots, x_n) \in B$ , which is in the density (3.7). In the following it is shown that  $\varphi_c$  can be rewritten as a Neyman-Pearson test.

For the density

$$f_{\theta}(x) = \exp\{c(\theta)T(x) + a(\theta)\} \cdot l(x)$$

assume that l(x) > 0, c'(x) > 0, and a'(x) exists for  $\theta \in \Theta$ .

**Lemma 3.3.22.** Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables with probability density function  $f_{\theta}(x), x \in B$ , which is a member of the one-parametric exponential family. Let T(x) be the respective statistic in the exponent of the density  $f_{\theta}$ . For arbitrary real numbers  $c_1 \leq c_2$ ,  $\gamma_1, \gamma_2 \in [0, 1]$  and parameters  $\theta_0, \theta_1 \in \Theta : \theta_0 \neq \theta_1$  the test  $\varphi_c$  in (3.9) can be rewritten as a modified Neyman-Pearson test  $\varphi_{K,\lambda}$  as in (3.8) with  $K, \lambda \in \mathbb{R}$ ,  $\gamma(x) \in [0, 1]$ .

**Proof** If the notation

$$f_{\theta_i}(x) = f_i(x), \quad i = 0, 1,$$

is used then

$$\frac{f_1(x)}{f_0(x)} = \exp\Big\{\underbrace{(c(\theta_1) - c(\theta_0))}_{c} T(x) + \underbrace{a(\theta_1) - a(\theta_0)}_{a}\Big\},\,$$

and therefore

$$\{x \in B : f_1(x) > (K + \lambda T(x)) f_0(x)\} = \{x \in B : \exp(cT(x) + a) > K + \lambda T(x)\}.$$

Can one find such K and  $\lambda$  in  $\mathbb{R}$  for the line  $K + \lambda t$ ,  $t \in \mathbb{R}$ , which intersects or touches the convex curve  $\exp(ct + a)$  exactly in  $c_1$  and  $c_2$  (if  $c_1 \neq c_2$ ) or in  $t = c_1$  (if  $c_1 = c_2$ ) resp.? As it turns out, such K and  $\lambda$  can always be found (cf. Figure 3.6).

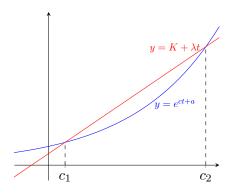


Figure 3.6: Intersection of a line with a convex curve

Let 
$$\gamma(x) = \gamma_i$$
 for  $\{x \in B : T(x) = c_i\}$ . Then

$${x : \exp(cT(x) + a) > K + \lambda T(x)} = {x : T(x) \notin [c_1, c_2]}$$

and

$${x : \exp(cT(x) + a) < K + \lambda T(x)} = {x : T(x) \in (c_1, c_2)}.$$

#### Remark 3.3.23.

- 1. The inversion of Lemma 3.3.22 does not hold, since for given curves  $y = K + \lambda t$  and  $y = \exp(ct + a)$  the intersections  $c_1$  and  $c_2$  do not have to exist. The line can be underneath the curve  $y = \exp(ct + a)$ .
- 2. The test  $\varphi_c$  does not explicitly use the parameters  $\theta_0$  and  $\theta_1$ , which makes it different from  $\varphi_{K,\lambda}$ , since it uses the densities  $f_0$  and  $f_1$ .

In the following, the fundamental theorem for two-sided tests for the hypotheses

$$H_0: \theta = \theta_0 \text{ vs. } H_1: \theta \neq \theta_0$$

will be presented.

#### Theorem 3.3.24. Fundamental theorem for two-sided tests

Under the conditions of Lemma 3.3.22, let  $\varphi_c$  be a test as in (3.9), for which  $\varphi_c \in \tilde{\Psi}(\alpha)$  holds. Then  $\varphi_c$  is the most powerful unbiased test with confidence level  $\alpha$  (and thus most powerful test in  $\tilde{\Psi}(\alpha)$ ) for the hypotheses

$$H_0: \theta = \theta_0 \text{ vs. } H_1: \theta \neq \theta_0.$$

**Proof** Let  $\theta_1 \in \Theta$ ,  $\theta_1 \neq \theta_0$  be arbitrary. By Lemma 3.3.22,  $\varphi_c$  is a modified Neyman-Pearson test  $\varphi_{K,\lambda}$  for a specific choice of K and  $\lambda \in \mathbb{R}$ , but  $\varphi_{K,\lambda}$  is a most powerful test in  $\tilde{\Psi}(\alpha)$  by Theorem 3.3.21 for  $H_0: \theta = \theta_0$  vs.  $H_1': \theta = \theta_1$ . Since  $\varphi_c$  does not depend on  $\theta_1$ , it is the most powerful test in  $\tilde{\Psi}(\alpha)$  for  $H_1: \theta \neq \theta_0$ . Since unbiased tests with confidence level  $\alpha$  are in  $\tilde{\Psi}(\alpha)$  it only has to be shown that  $\varphi_c$  is unbiased.  $\varphi_c$  is the most powerful test and thus not worse than the constant unbiased test  $\varphi = \alpha$ , i.e.

$$G_n(\varphi_c, \theta) \ge G_n(\varphi, \theta) = \alpha, \quad \theta \ne \theta_0.$$

Thus  $\varphi_c$  is also unbiased.

**Remark 3.3.25.** It has been shown that  $\varphi_c$  is the most powerful test within its scope. It should still be shown that for arbitrary  $\alpha \in (0,1)$  constants  $c_1, c_2, \gamma_1, \gamma_2$  can be found, which satisfy  $\mathbb{E}_0 \varphi_c = \alpha$ . The proof is rather technical and will thus be omitted here. The following example shows how  $c_1, c_2, \gamma_1, \gamma_2$  have to be chosen.

## Example 3.3.26. Two-sided-Gauss-test

Example 3.1.7 considers the following test for the expectation of a normally distributed random sample  $X = (X_1, ..., X_n)$  with i.i.d.  $X_i$  and  $X_i \sim \mathcal{N}(\mu, \sigma_0^2)$  where  $\sigma_0^2$  is known. The hypotheses

$$H_0: \mu = \mu_0 \text{ vs. } H_1: \mu \neq \mu_0$$

are tested. The test  $\varphi(x)$  is given by

$$\varphi(x) = I\left(x \in \mathbb{R}^n : |T(x)| > z_{1-\alpha/2}\right),\,$$

where

$$T(x) = \sqrt{n} \frac{\overline{x}_n - \mu_0}{\sigma_0}.$$

It has to be shown that  $\varphi$  is the most powerful test with confidence level  $\alpha$  in  $\tilde{\Psi}(\alpha)$  (and thus the most powerful unbiased test). By Theorem 3.3.24 it has to be shown that  $\varphi$  can be rewritten as  $\varphi_c$  with (3.9), since the n-dimensional Normal distribution with probability density function  $f_{\mu}$  (cf. example 3.3.17) is a member of the one-parametric exponential family with statistic

$$T(x) = \sqrt{n} \frac{\overline{x}_n - \mu}{\sigma_0}.$$

Let  $c_1 = -z_{1-\alpha/2}, c_2 = z_{1-\alpha/2}, \gamma_1 = \gamma_2 = 0$ . Then

$$\varphi(x) = \varphi_c(x) = \begin{cases} 1, & \text{if } |T(x)| > z_{1-\alpha/2}, \\ 0, & \text{if } |T(x)| \le z_{1-\alpha/2}. \end{cases}$$

The assertion is thus proven, since the power function  $G_n(\varphi, \theta)$  of  $\varphi$  as in Example 3.1.7 implies, that  $\varphi$  is an unbiased test with confidence level  $\alpha$  (and therefore  $\varphi \in \tilde{\Psi}(\alpha)$ ).

**Remark 3.3.27.** So far, we only assumed that *one* parameter of the distribution of the random sample  $(X_1, \ldots, X_n)$  is unknown. This has been necessary in order to be able to introduce the above theory of most powerful (Neyman-Pearson) tests for one-parametric exponential families. In order to consider the case with more unknown parameters (as in the example of two-sided tests for the expected value of a normally distributed random sample with unknown variance), a deeper understanding of randomized tests is needed. If one is interested, the theory can be found in [26].

## 3.4 Goodness-of-fit tests

Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables with  $X_i \sim F$  for  $i = 1, \ldots, n$ . Goodness-of-fit testing tests the hypotheses

$$H_0: F = F_0 \text{ vs. } H_1: F \neq F_0,$$

where  $F_0$  is a given distribution function.

A goodness-of-fit test has already been introduced in the lecture "Elementary probability theory". The Kolmogorow-Smirnov test can be found in Remark 7.6.8.

In this section, further non-parametric goodness-of-fit tests are introduced. The first one, namely the  $\chi^2$ -goodness-of-fit test, was introduced by K. Pearson.

## 3.4.1 $\chi^2$ -goodness-of-fit test

The Kolmogorov-Smirnov test is based on the distance

$$D_n = \sup_{x \in \mathbb{R}} | \hat{F}_n(x) - F_0(x) |$$

between the empirical distribution function of the random sample  $(X_1, \ldots, X_n)$  and the distribution function  $F_0$ . In practice this test is usually too sensitive, since irregularities in the random samples might lead to an unjustified rejection of  $H_0$ . A solution to this problem is a test which coarsens the null hypothesis  $H_0$  and is based on the  $\chi^2$ -goodness-of-fit statistic.

Partition the domain of  $X_i$  into r classes  $(a_j, b_j], j = 1, ..., r$  with the property

$$-\infty \le a_1 < b_1 = a_2 < b_2 = \dots = a_r < b_r \le \infty.$$

Instead of  $X_i$ , i = 1, ..., n, consider the so-called *class sizes*  $Z_j$ , j = 1, ..., r, where

$$Z_j = \#\{i : a_j < X_i \le b_j, 1 \le i \le n\}.$$

**Lemma 3.4.1.** The random vector  $Z = (Z_1, \ldots, Z_r)^{\top}$  is multinomial distributed with parameter vector

$$p = (p_1, \dots, p_{r-1})^{\top} \in [0, 1]^{r-1},$$

where

$$p_j = P(a_j < X_1 \le b_j) = F(b_j) - F(a_j), \ j = 1, \dots, r - 1, \quad p_r = 1 - \sum_{j=1}^{r-1} p_j.$$

Notation:

$$Z \sim M_{r-1}(n, p)$$
.

**Proof** We show that for all numbers  $k_1, \ldots k_r \in \mathbb{N}_0$  with  $k_1 + \ldots + k_r = n$ 

$$P(Z_i = k_i, i = 1, \dots, r) = \frac{n!}{k_1! \dots k_r!} p_1^{k_1} \dots p_r^{k_r}$$
 (3.10)

holds. Since  $X_i$  are i.i.d. it holds that

$$P(X_j \in (a_{i_j}, b_{i_j}], j = 1, \dots, n) = \prod_{j=1}^n P(a_{i_j} < X_1 \le b_{i_j}) = p_1^{k_1} \cdot \dots \cdot p_r^{k_r},$$

if the sequence of intervals  $(a_{i_j}, b_{i_j}]_{j=1,\dots,n}$  contains the interval  $(a_i, b_i]$   $k_i$  times,  $i=1,\dots,r$ . The formula (3.10) results from the law of total probability as a sum over all permutations of sequences  $(a_{i_j}, b_{i_j}]_{j=1,\dots,n}$ .

In the sense of Lemma 3.4.1 new hypotheses w.r.t. the nature of F are tested:

$$H_0: p = p_0 \text{ vs. } H_1: p \neq p_0,$$

where  $p = (p_1, \dots, p_{r-1})^{\top}$  is the parameter vector of Z, and  $p_0 = (p_{01}, \dots, p_{0,r-1})^{\top} \in (0,1)^{r-1}$  with  $\sum_{i=1}^{r-1} p_{0i} < 1$ . In this case,

$$\Lambda_0 = \{ F \in \Lambda : F(b_i) - F(a_i) = p_{0i}, \quad j = 1, \dots, r - 1 \}$$

and  $\Lambda_1 = \Lambda \setminus \Lambda_0$  holds, where  $\Lambda$  is the set of all distribution functions. In order to test  $H_0$  vs.  $H_1$ , the *Pearson test statistic* 

$$T_n(x) = \sum_{j=1}^r \frac{(z_j - np_{0j})^2}{np_{0j}},$$

where  $x = (x_1, ..., x_n)$  is an explicit sample and  $z_j$ , j = 1, ..., r the corresponding class sizes. Under  $H_0$ ,

$$\mathbb{E} Z_i = n p_{0i}, \quad j = 1, \dots, r,$$

holds and thus  $H_0$  is rejected, if  $T_n(X)$  attains higher values than expected. The following theorem shows that  $T(X_1, \ldots, X_n)$  is asymptotically (for  $n \to \infty$ )  $\chi^2_{r-1}$ -distributed, which leads to the following goodness-of-fit test ( $\chi^2$  goodness-of-fit test):

$$H_0$$
 is rejected, if  $T_n(x_1,\ldots,x_n) > \chi^2_{r-1,1-\alpha}$ .

This test is named after its inventor Karl Pearson (1857-1936).

**Theorem 3.4.2.** Under  $H_0$ ,

$$\lim_{n \to \infty} P_{p_0} \left( T_n(X_1, \dots, X_n) > \chi^2_{r-1, 1-\alpha} \right) = \alpha, \ \alpha \in (0, 1),$$

holds, which means the  $\chi^2$ -Pearson test is an asymptotic test with confidence level  $\alpha$ .

**Proof** Denote by  $Z_{nj} = Z_j(X_1, ..., X_n)$  the class sizes of the random samples  $(X_1, ..., X_n)$ . By Lemma 3.4.1

$$Z_n = (Z_{n1}, \dots, Z_{nr}) \sim M_{r-1}(n, p_0)$$
 under  $H_0$ 

holds. Moreover,  $\mathbb{E} Z_{nj} = np_{0j}$  and

$$Cov(Z_{ni}, Z_{nj}) = \begin{cases} np_{0j}(1 - p_{0j}), & i = j, \\ -np_{0i}p_{0j}, & i \neq j \end{cases}$$

should hold for all i, j = 1, ..., r. Since

$$Z_{nj} = \sum_{i=1}^{n} I(a_j < X_i \le b_j), \quad j = 1, \dots, r,$$

it holds that  $Z_n = (Z_{n1}, \ldots, Z_{n,r-1})$  is the sum of n i.i.d. random vectors with  $Y_i \in \mathbb{R}^{r-1}$  with coordinates  $Y_{ij} = I(a_j < X_i \le b_j), j = 1, \ldots, r-1$ . Thus, the multivariate limit theorem (which is proven in Lemma 3.4.3) yields

$$Z'_n = \frac{Z_n - \mathbb{E} Z_n}{\sqrt{n}} = \frac{\sum_{i=1}^n Y_i - n\mathbb{E} Y_1}{\sqrt{n}} \xrightarrow[n \to \infty]{d} Y \sim \mathcal{N}(0, K),$$

with  $\mathcal{N}(0, K)$  a (r-1) dimensional multivariate normal distribution (cf. [33, Example 3.4.5.3.] with expectation vector 0 and covariance matrix  $K = (\sigma_{ij}^2)$ , where

$$\sigma_{ij}^2 = \begin{cases} -p_{0i}p_{0j}, & i \neq j, \\ p_{0i}(1 - p_{0j}), & i = j \end{cases}$$

for i, j = 1, ..., r - 1. This matrix K is invertible with  $K^{-1} = A = (a_{ij})$ ,

$$a_{ij} = \begin{cases} \frac{1}{p_{0r}}, & i \neq j, \\ \frac{1}{p_{0i}} + \frac{1}{p_{0r}}, & i = j. \end{cases}$$

Moreover, K (as a covariance matrix) is symmetric and positive semi-definite. Results from Linear Algebra ensure the existence of an invertible  $(r-1) \times (r-1)$  matrix  $A^{1/2}$ , with  $A = A^{1/2}(A^{1/2})^{\top}$ . Thus,

$$K = A^{-1} = ((A^{1/2})^{\top})^{-1} \cdot (A^{1/2})^{-1}.$$

If  $(A^{1/2})^{\top}$  is applied to  $Z'_n$ , we get

$$(A^{1/2})^{\top} \cdot Z'_n \xrightarrow[n \to \infty]{d} (A^{1/2})^{\top} \cdot Y,$$

where

$$(A^{1/2})^{\top} \cdot Y \sim \mathcal{N}\left(0, (A^{1/2})^{\top} \cdot K \cdot A^{1/2}\right) = \mathcal{N}\left(0, \mathcal{I}_{r-1}\right)$$

by the properties of the multivariate normal distribution. Furthermore, the continuous mapping theorem, which has been introduced in [32, Theorem 3.4.4.], implies that

$$Y_n \xrightarrow[n \to \infty]{d} Y \Longrightarrow \varphi(Y_n) \xrightarrow[n \to \infty]{d} \varphi(Y)$$

for random variables  $\{Y_n\}$ ,  $Y \in \mathbb{R}^m$ , and continuous mappings  $\varphi : \mathbb{R} \to \mathbb{R}$ . Repeatedly applying the continuous mapping theorem implies that

$$\left|(A^{1/2})^\top Z_n'\right|^2 \overset{d}{\underset{n \to \infty}{\longrightarrow}} \left|(A^{1/2})^\top Y\right|^2 = R \sim \chi_{r-1}^2.$$

It needs to be shown that

$$T_n(X_1,\ldots,X_n) = \left| (A^{1/2})^{\top} Z_n' \right|^2.$$

Now,

$$\begin{aligned} \left| (A^{1/2})^{\top} Z_n' \right|^2 &= ((A^{1/2})^{\top} Z_n')^{\top} ((A^{1/2})^{\top} Z_n') \\ &= Z_n'^{\top} \cdot \underbrace{A^{1/2} \cdot (A^{1/2})^{\top}}_{A} Z_n' = Z_n'^{\top} A Z_n' \\ &= n \sum_{j=1}^{r-1} \frac{1}{p_{0j}} \left( \frac{Z_{nj}}{n} - p_{0j} \right)^2 + \frac{n}{p_{0r}} \sum_{i=1}^{r-1} \sum_{j=1}^{r-1} \left( \frac{Z_{ni}}{n} - p_{0i} \right) \left( \frac{Z_{nj}}{n} - p_{0j} \right) \\ &= \sum_{j=1}^{r-1} \frac{(Z_{nj} - np_{0j})^2}{np_{0j}} + \frac{n}{p_{0r}} \left( \sum_{j=1}^{r-1} \left( \frac{Z_{nj}}{n} - p_{0j} \right) \right)^2 \\ &= \sum_{j=1}^{r-1} \frac{(Z_{nj-np_{0j}})^2}{np_{0j}} + \frac{n}{p_{0r}} \left( \frac{Z_{nr}}{n} - p_{0r} \right)^2 \\ &= \sum_{j=1}^{r} \frac{(Z_{nj} - np_{0j})^2}{np_{0j}} = T_n(X_1, \dots, X_n), \end{aligned}$$

since

$$\sum_{j=1}^{r-1} Z_{nj} = n - Z_{nr},$$

$$\sum_{j=1}^{r-1} p_{0j} = 1 - p_{0r}.$$

## Lemma 3.4.3. Multivariate central limit theorem

Let  $\{Y_n\}_{n\in\mathbb{N}}$  be a sequence of i.i.d. random vectors with  $\mathbb{E}Y_1 = \mu \in \mathbb{R}^m$  and covariance matrix  $K \in \mathbb{R}^{m \times m}$ . Then

$$\frac{\sum_{i=1}^{n} Y_i - n\mu}{\sqrt{n}} \xrightarrow[n \to \infty]{d} Y \sim N(0, K).$$
 (3.11)

**Proof** Let  $Y_j = (Y_{j1}, \dots, Y_{jm})^{\top}$ . By the continuous mapping theorem for characteristic functions the convergence in (3.11) is equivalent to

$$\varphi_n(t) \underset{n \to \infty}{\longrightarrow} \varphi(t), \quad t \in \mathbb{R}^m,$$
 (3.12)

where

$$\varphi_n(t) = \mathbb{E} e^{itS_n} = \mathbb{E} \exp \left\{ i \sum_{j=1}^m t_j \frac{Y_{1j} + \ldots + Y_{nj} - n\mu_j}{\sqrt{n}} \right\},$$

is the characteristic function of the random vector

$$S_n = \frac{\sum_{i=1}^n Y_i - n\mu}{\sqrt{n}},$$

and

$$\varphi(t) = e^{-t^{\top} K t / 2}$$

is the characteristic function of the  $\mathcal{N}(0,K)$  distribution. The function  $\varphi_n(t)$  can be rewritten as

$$\varphi_n(t) = \mathbb{E} \exp \left\{ i \sum_{i=1}^n \frac{\sum_{i=1}^m t_j (Y_{ij} - \mu_j)}{\sqrt{n}} \right\}, \quad t = (t_1, \dots, t_m)^\top \in \mathbb{R}^m,$$

where

$$L_i := \sum_{j=1}^m t_j (Y_{ij} - \mu_j)$$

is a random variable with

$$\mathbb{E} L_i = 0,$$

$$\operatorname{Var} L_i = \mathbb{E} \left[ \sum_{k,j=1}^m t_j (Y_{ij} - \mu_j) (Y_{ik-\mu_k}) t_k \right] = t^\top K t, \quad i \in \mathbb{N}.$$

If  $t^{\top}Kt = 0$ , then  $L_i = 0$  almost surely for all  $i \in \mathbb{N}$ , which implies  $\varphi_n(t) = \varphi(t) = 1$ . Thus the convergence in (3.11) holds. If  $t^{\top}Kt > 0$ , then  $\varphi_n(t)$  is the characteristic function of a random variable

$$\sum_{i=1}^{n} L_i / \sqrt{n}$$

evaluated at 1, and  $\varphi(t)$  is the characteristic function of a one-dimensional normal distribution  $\mathcal{N}(0, t^{\top}Kt)$  evaluated at 1. The central limit theorem for one-dimensional random variables then implies (cf. [33, Theorem 5.2.2.])

$$\sum_{i=1}^{n} \frac{L_i}{\sqrt{n}} \xrightarrow[n \to \infty]{d} L \sim \mathcal{N}(0, t^{\top}Kt)$$

and thus

$$\varphi_n(t) = \varphi_{\left(\sum_{i=1}^n L_i/\sqrt{n}\right)}(1) \underset{n \to \infty}{\longrightarrow} \varphi_L(1) = \varphi(t),$$

which proves the convergence in (3.11).

#### Remark 3.4.4.

- 1. The method of reducing a multidimensional convergence to a onedimensional convergence, using linear combinations of random variables, as in the proof above is called *Cramér-Wold* device.
- 2. The  $\chi^2$ -Pearson test works asymptotically for large random samples. Naturally the question of how big n is arises. In this case, the "rule of thumb" is given by:  $np_{0j}$  should be larger or equal to a, with  $a \in (2,\infty)$ . For a larger class number, i.e.,  $r \geq 10$ , even a=1 is sufficient. In the following, it is shown that the  $\chi^2$  goodness-of-fit test is consistent.

**Lemma 3.4.5.** The  $\chi^2$ -Pearson test is consistent, i.e., for all  $p \in [0,1]^{r-1}, p \neq p_0$ 

$$\lim_{n \to \infty} P_p \left( T_n(X_1, \dots, X_n) > \chi_{r-1, 1-\alpha}^2 \right) = 1$$

holds.

**Proof** Under  $H_1$ , the strong law of large numbers implies

$$\frac{Z_{nj}}{n} = \frac{\sum_{i=1}^{n} I(a_j < X_i \le b_j)}{n} \xrightarrow[n \to \infty]{a.s.} \underbrace{\mathbb{E} I(a_j < X_1 \le b_j)}_{=p_j}.$$

Choose j such that  $p_j \neq p_{0j}$ . Then

$$T_n(X_1,\ldots,X_n) \ge \frac{(Z_{nj} - np_{0j})^2}{np_{0j}} \ge \underbrace{n\left(\frac{Z_{nj}}{n} - p_{0j}\right)^2}_{\sim n(p_j - p_{0j})^2} \xrightarrow[n \to \infty]{a.s.} \infty,$$

and thus

$$P_p\left(T_n(X_1,\ldots,X_n)>\chi^2_{r-1,1-\alpha}\right) \xrightarrow[n\to\infty]{f.s.} 1.$$

## 3.4.2 $\chi^2$ -goodness-of-fit test of Pearson-Fisher

Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables  $X_i$ ,  $i = 1, \ldots, n$ . The goal is to test whether the distribution function F of  $X_i$  is an element of a given parametric family

$$\Lambda_0 = \{ F_\theta : \theta \in \Theta \}, \quad \Theta \subset \mathbb{R}^m$$

Let  $a_i, b_i, i = 1, ..., r$  be given with m < r,

$$-\infty \le a_1 < b_1 = a_2 < b_2 = \dots = a_r < b_r \le \infty$$

and

$$Z_{nj} = \#\{X_i, i = 1, \dots, n : a_j < X_i \le b_j\}, \quad j = 1, \dots, r,$$
  
 $Z_n = (Z_{n1}, \dots, Z_{nr})^{\top}.$ 

Lemma 3.4.1 implies  $Z \sim M_{r-1}(n,p)$ ,  $p = (p_0, \ldots, p_{r-1})^{\top} \in [0,1]^{r-1}$ . Under  $H_0: F \in \Lambda_0, p = p(\theta), \theta \in \Theta \subset \mathbb{R}^m$  holds. Presume  $p \in C(\Theta)$ . By coarsening the hypothesis  $H_0$  the new hypotheses:

$$H_0: p \in \{p(\theta): \theta \in \Theta\} \text{ vs. } H_1: p \notin \{p(\theta): \theta \in \Theta\}$$

are to be tested. In order to test these hypotheses, the  $\chi^2$ -Pearson-Fisher test is constructed as follows:

- 1. Find a maximum-likelihood estimator  $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n)$  (weakly consistent) for  $\theta$ , such that  $\hat{\theta}_n \xrightarrow{P}_{n \to \infty} \theta$ . Here,  $\{\hat{\theta}_n\}_{n \in \mathbb{N}}$  is asymptotically normal distributed.
- 2. Construct the plug-in estimator  $p(\hat{\theta}_n)$  for  $p(\theta)$ .
- 3. For the test statistic

$$\hat{T}_n(X_1, \dots, X_n) = \sum_{i=1}^r \frac{\left(Z_{nj} - np_j(\hat{\theta})\right)^2}{np_j(\hat{\theta})} \xrightarrow[n \to \infty]{d} \eta \sim \chi_{r-m-1}^2$$

holds under  $H_0$  and certain assumptions.

4.  $H_0$  is rejected, if  $\hat{T}_n(X_1, \ldots, X_n) > \chi^2_{r-m-1,1-\alpha}$ . This is an asymptotic test with confidence level  $\alpha$ .

#### Remark 3.4.6.

- 1. The  $\chi^2$ -Pearson-Fisher test assumes that the function  $p(\theta)$  can be stated explicitly, but  $\theta$  is unknown. That means for every class of distributions  $\Lambda_0$ , the function  $p(\cdot)$  has to be calculated.
- 2. Why is  $\hat{T}_n$  able to discriminate between the hypotheses  $H_0$  and  $H_1$ ? The strong law of large numbers implies

$$\frac{1}{n}Z_{nj} - p_j(\hat{\theta}_n) = \underbrace{\frac{1}{n}Z_{nj} - p_j(\theta)}_{\stackrel{P}{\to}0} - \underbrace{(p_j(\hat{\theta}_n) - p_j(\theta))}_{\stackrel{P}{\to}0} \xrightarrow{\stackrel{P}{\to}0} 0,$$

if  $\hat{\theta}_n$  is weakly consistent and  $p_j(\cdot)$  a continuous function for all  $j = 1, \ldots, r$ .

Thus, under  $H_0$   $\hat{T}_n(X_1, ..., X_n)$  is supposed to take relatively small values. A significant deviation of this behavior is supposed to lead to the rejection of  $H_0$ .

For the distribution  $F_{\theta}$  of  $X_i$  the following regularity properties are assumed to hold (cf. Theorem 1.2.22).

- 1. The distribution function  $F_{\theta}$  is either absolutely continuous or discrete for all  $\theta \in \Theta$ .
- 2. The parametrization is unique, i.e.  $\theta \neq \theta_1 \Leftrightarrow F_{\theta} \neq F_{\theta_1}$ .
- 3. The support supp  $L(x,\theta) = \{x \in \mathbb{R} : L(x,\theta) > 0\}$  of the likelihood function given by

$$L(x,\theta) = \begin{cases} P_{\theta}(X_1 = x), & \text{in case of discrete } F_{\theta}, \\ f_{\theta}(x), & \text{in the absolutely continuous case,} \end{cases}$$

does not depend on  $\theta$ .

4.  $L(x,\theta)$  is assumed to be three times continuously differentiable and for  $k = 1, \ldots, 3$  and  $i_1, \ldots, i_k \in \{1, \ldots, m\}$ ,

$$\sum_{k=1}^{3} \int \frac{\partial^{k} L(x,\theta)}{\partial \theta_{i_{1}} \cdot \ldots \cdot \partial \theta_{i_{k}}} dx = \frac{\partial^{k}}{\partial \theta_{i_{1}} \cdot \ldots \cdot \partial \theta_{i_{k}}} \sum_{k=1}^{3} \int L(x,\theta) dx = 0$$

holds.

5. For all  $\theta_0 \in \Theta$  there exist a constant  $c_{\theta_0}$  and a measurable function  $g_{\theta_0} : \text{supp } L \to \mathbb{R}_+$ , such that

$$\left| \frac{\partial^3 \log L(x,\theta)}{\partial \theta_{i_1} \partial \theta_{i_2} \partial \theta_{i_3}} \right| \le g_{\theta_0}(x), \quad |\theta - \theta_0| < c_{\theta_0},$$

and

$$\mathbb{E}_{\theta_0} g_{\theta_0}(X_1) < \infty.$$

Define the Fisher information matrix by

$$I(\theta) = \left( \mathbb{E} \left[ \frac{\partial \log L(X_1, \theta)}{\partial \theta_i} \frac{\partial \log L(X_1, \theta)}{\partial \theta_j} \right] \right)_{i, j = 1, \dots, m}.$$
 (3.13)

# Theorem 3.4.7. Asymptotical normal distribution of consistent maximum likelihood estimator $\hat{\theta}_n$ , multivariate case m > 1

Let  $X_1, \ldots, X_n$  be i.i.d. with likelihood function L, which satisfies the regularity assumptions 1.-5. Let  $I(\theta)$  be positive definite for all  $\theta \in \Theta \subset \mathbb{R}^m$  and  $\hat{\theta}_n = \hat{\theta}(X_1, \ldots, X_n)$  be a sequence of weakly consistent maximum likelihood estimators for  $\theta$ . Then

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow[n \to \infty]{d} N(0, I^{-1}(\theta)).$$

Without proof (cf. proof of Theorem 1.2.22).

For the coarsed hypothesis  $H_0: p \in \{p(\theta), \theta \in \Theta\}$  construct the piecewise constant likelihood function

$$L(x,\theta) = p_j(\theta), \text{ if } x \in (a_j,b_j].$$

Then, the likelihood function of the random sample  $(x_1, \ldots, x_n)$  is given by

$$L(x_1, ..., x_n, \theta) = \prod_{j=1}^r p_j(\theta)^{Z_j(x_1, ..., x_n)}$$

$$\Rightarrow \log L(x_1, \dots, x_n, \theta) = \sum_{j=1}^r Z_j(x_1, \dots, x_n) \cdot \log p_j(\theta).$$

For the maximum likelihood estimator, we get

$$\hat{\theta}_n = \hat{\theta}(x_1, \dots, x_n) = \operatorname*{argmax}_{\theta \in \Theta} \log L(x_1, \dots, x_n, \theta)$$

$$\Rightarrow \sum_{j=1}^{r} Z_j(x_1, \dots, x_n) \frac{\partial p_j(\theta)}{\partial \theta_i} \cdot \frac{1}{p_j(\theta)} = 0, \quad i = 1, \dots, m.$$

Furthermore, the property  $\sum_{j=1}^{r} p_j(\theta) = 1$  implies

$$\sum_{j=1}^{r} \frac{\partial p_j(\theta)}{\partial \theta_i} = 0 \Rightarrow \sum_{j=1}^{r} \frac{Z_j(x_1, \dots, x_n) - np_j(\theta)}{p_j(\theta)} \cdot \frac{\partial p_j(\theta)}{\partial \theta_i} = 0, \quad i = 1, \dots, m.$$

**Lemma 3.4.8.** In the case above,  $I(\theta) = C^{\top}(\theta) \cdot C(\theta)$  holds, where  $C(\theta)$  is a  $(r \times m)$ -matrix with elements

$$c_{ij}(\theta) = \frac{\partial p_i(\theta)}{\partial \theta_j} \cdot \frac{1}{\sqrt{p_i(\theta)}}.$$

Proof

$$\mathbb{E}_{0} \left[ \frac{\partial \log L(X_{1}, \theta)}{\partial \theta_{i}} \cdot \frac{\partial \log L(X_{1}, \theta)}{\partial \theta_{j}} \right] = \sum_{k=1}^{r} \frac{\partial \log p_{k}(\theta)}{\partial \theta_{i}} \cdot \frac{\partial \log p_{k}(\theta)}{\partial \theta_{j}} \cdot p_{k}(\theta)$$

$$= \sum_{k=1}^{r} \frac{\partial p_{k}(\theta)}{\partial \theta_{i}} \frac{1}{p_{k}(\theta)} \cdot \frac{\partial p_{k}(\theta)}{\partial \theta_{j}} \cdot \frac{1}{p_{k}(\theta)} \cdot p_{k}(\theta)$$

$$= \left( C^{\top}(\theta) \cdot C(\theta) \right)_{ij},$$

since

$$\log L(X_1, \theta) = \sum_{i=1}^{r} \log p_j(\theta) \cdot I \left( x \in (a_j, b_j] \right).$$

Theorem 3.4.7 implies

Corollary 3.4.9. Let  $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n)$  be a weakly consistent maximum likelihood estimator of  $\theta$  in the coarsened model, which satisfies the regularity assumptions 1.-5. Assume that the Fisher information matrix  $I(\theta) = C^{\top}(\theta) \cdot C(\theta)$  is positive definite for all  $\theta \in \Theta$ . Then,  $\hat{\theta}$  is asymptotically normal distributed

$$\sqrt{n}\left(\hat{\theta}_n - \theta\right) \xrightarrow[n \to \infty]{d} Y \sim \mathcal{N}\left(0, I^{-1}(\theta)\right).$$

**Theorem 3.4.10.** Let  $\hat{\theta}_n$  be a maximum likelihood estimator in the coarsed model for  $\theta$ , which satisfies all assumptions of Corollary 3.4.9. The test statistic

$$\hat{T}_n(X_1, \dots, X_n) = \sum_{i=1}^r \frac{(Z_j(X_1, \dots, X_n) - np_j(\hat{\theta}_n))^2}{np_j(\hat{\theta}_n)}$$

is asymptotically  $\chi^2_{r-m-1}$ -distributed under  $H_0$ :

$$\lim_{n \to \infty} P_{\theta} \left( \hat{T}_n(X_1, \dots, X_n) > \chi^2_{r-m-1, 1-\alpha} \right) = \alpha.$$

Without proof (cf. [27]).

This theorem implies that the  $\chi^2$ -Pearson-Fisher test is an asymptotic test with confidence level  $\alpha$ .

#### Example 3.4.11.

1.  $\chi^2$ -Pearson-Fisher test of the normal distribution Let  $(X_1, \ldots, X_n)$  be a random sample. We test, whether  $X_i \sim \mathcal{N}(\mu, \sigma^2)$ . Define

$$\theta := (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_+.$$

Let  $\{(a_j,b_j]\}_{j=1,\dots,r}$  be an arbitrary partition of  $\mathbb{R}$  in r disjoint intervals. Recall that density of the one-dimensional  $\mathcal{N}(\mu,\sigma^2)$ -distribution is given by

$$f_{\theta}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

and define

$$p_j(\theta) := P_0(X_1 \in (a_j, b_j]) = \int_{a_j}^{b_j} f_{\theta}(x) dx, \quad j = 1, \dots, r$$

with class sizes

$$Z_j = \# \{i : X_i \in (a_j, b_j]\}.$$

The goal is to find the maximum-likelihood estimator in the coarsed model

$$\begin{split} \frac{\partial p_j(\theta)}{\partial \mu} &= \int_{a_j}^{b_j} \frac{\partial}{\partial \mu} f_{\theta}(x) dx = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \int_{a_j}^{b_j} \frac{x - \mu}{\sigma^2} \cdot e^{-\frac{1}{2} \left(\frac{x - \mu}{\sigma}\right)^2} dx, \\ \frac{\partial p_j(\theta)}{\partial \sigma^2} &= \int_{a_j}^{b_j} \frac{\partial}{\partial \sigma^2} f_{\theta}(x) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{a_j}^{b_j} \left[ -\frac{1}{2} \cdot \frac{1}{(\sigma^2)^{3/2}} e^{-\frac{1}{2} \left(\frac{x - \mu}{\sigma}\right)^2} \right. \\ &\left. + \frac{1}{\sqrt{\sigma^2}} e^{-\frac{1}{2} \left(\frac{x - \mu}{\sigma}\right)^2} \cdot \left(\frac{(x - \mu)^2}{2(\sigma^2)^2}\right) \right] dx \\ &= -\frac{1}{2} \frac{1}{\sigma^2} \int_{a_j}^{b_j} f_{\theta}(x) dx + \frac{1}{2(\sigma^2)^2} \int_{a_j}^{b_j} (x - \mu)^2 f_{\theta}(x) dx. \end{split}$$

The necessary conditions for a maximum are

$$\sum_{i=1}^{r} Z_{j} \frac{\int_{a_{j}}^{b_{j}} x f_{\theta}(x) dx}{\int_{a_{j}}^{s} f_{\theta}(x) dx} - \mu \sum_{j=1}^{r} Z_{j} = 0,$$

$$\frac{1}{\sigma^{2}} \sum_{j=1}^{r} Z_{j} \frac{\int_{a_{j}}^{b_{j}} (x - \mu)^{2} f_{\theta}(x) dx}{\int_{a_{j}}^{s} f_{\theta}(x) dx} - \sum_{j=1}^{r} Z_{j} = 0,$$

which results in the maximum likelihood estimators  $\hat{\mu}$  and  $\hat{\sigma}^2$  for  $\mu$  and  $\sigma^2$ 

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^{r} Z_{j} \frac{\int_{a_{j}}^{b_{j}} x f_{\theta}(x) dx}{\int_{a_{j}}^{b_{j}} f_{\theta}(x) dx}, \quad \hat{\sigma}^{2} = \frac{1}{n} \sum_{j=1}^{r} Z_{j} \frac{\int_{a_{j}}^{b_{j}} (x - \mu)^{2} f_{\theta}(x) dx}{\int_{a_{j}}^{b_{j}} f_{\theta}(x) dx}.$$

Construct an approximation of  $\hat{\mu}$  and  $\hat{\sigma}^2$  for  $r \to \infty$  as follows: If  $r \to \infty$  (and thus  $n \to \infty$ ), then  $b_j - a_j$  is small and by the simple quadratic rule

$$\int_{a_j}^{b_j} x f_{\theta}(x) dx \approx (b_j - a_j) y_j f_{\theta}(y_j),$$
$$\int_{a_j}^{b_j} f_{\theta}(x) dx \approx (b_j - a_j) f_{\theta}(y_j),$$

holds, where  $y_1 = b_1$ ,  $y_r = b_{r-1} = a_r$ , and

$$y_j = (b_{j+1} + b_j)/2, \quad j = 2, \dots, r - 1.$$

Thus, for the maximum likelihood estimators  $\hat{\mu}$  and  $\hat{\sigma}^2$ 

$$\hat{\mu} \approx \frac{1}{n} \sum_{j=1}^{r} y_j \cdot Z_j = \tilde{\mu}$$

$$\hat{\sigma}^2 \approx \frac{1}{n} \sum_{j=1}^{r} (y_j - \tilde{\mu})^2 Z_j = \tilde{\sigma}^2$$

and

$$\tilde{\theta} = \left(\tilde{\mu}, \tilde{\sigma}^2\right)$$

holds. In the  $\chi^2$ -Pearson-Fisher test  $H_0$  is rejected if

$$\hat{T}_n = \frac{\sum_{j=1}^r \left( Z_j - np_j(\tilde{\theta}) \right)^2}{np_j(\tilde{\theta})} > \chi_{r-3,1-\alpha}^2.$$

2.  $\chi^2$ -Pearson-Fisher test for the Poisson distribution

Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables. We aim to test, whether  $X_i \sim \text{Poisson}(\lambda)$ ,  $\lambda > 0$ . Set  $\theta = \lambda$  and  $\Theta = (0, +\infty)$ . Coarsing  $\Theta$  leads to

$$-\infty = a_1 < \underbrace{b_1}_{=0} = a_2 < \underbrace{b_2}_{=1} = a_3 < \dots < \underbrace{b_{r-1}}_{=r-2} = a_r < b_r = +\infty.$$

Then,

$$p_{j}(\lambda) = P_{\lambda} \left( X_{1} = j - 1 \right) = e^{-\lambda} \frac{\lambda^{j-1}}{(j-1)!}, \quad j = 1, \dots, r-1,$$

$$p_{r}(\lambda) = \sum_{i=r-1}^{\infty} e^{-\lambda} \frac{\lambda^{i}}{i!},$$

$$\frac{dp_{j}(\lambda)}{d\lambda} = -e^{-\lambda} \frac{\lambda^{j-1}}{(j-1)!} + (j-1) \frac{\lambda^{j-2}}{(j-1)!} e^{-\lambda} = e^{-\lambda} \frac{\lambda^{j-1}}{(j-1)!} \left( \frac{j-1}{\lambda} - 1 \right)$$

$$= p_{j}(\lambda) \cdot \left( \frac{j-1}{\lambda} - 1 \right), \quad j = 1, \dots, r-1,$$

$$\frac{dp_{r}(\lambda)}{d\lambda} = \sum_{i \geq r-1} p_{i}(\lambda) \left( \frac{i-1}{\lambda} - 1 \right).$$

Next, we have to solve the maximum likelihood equation

$$0 = \sum_{j=1}^{r-1} Z_j \cdot \left(\frac{j-1}{\lambda} - 1\right) + Z_r \frac{\sum_{i \ge r-1} p_i(\lambda) \left(\frac{i-1}{\lambda} - 1\right)}{p_r(\lambda)}.$$

If  $r \to \infty$ , then r(n) exists for every n with  $Z_{r(n)} = 0$ . Thus, for r > r(n)

$$\sum_{j=1}^{r-1} (j-1)Z_j - \lambda \sum_{j=1}^{r} Z_j = 0$$

holds, which yields the maximum likelihood estimator

$$\frac{1}{n}\sum_{j=1}^{r-1}(j-1)Z_j = \frac{1}{n}\sum_{j=1}^n X_j = \overline{X}_n.$$

Hence, the  $\chi^2$ -Pearson-Fisher test rejects  $H_0$ , if

$$\hat{T}_n = \sum_{j=1}^r \frac{\left(Z_j - np_\lambda(\overline{X}_n)\right)^2}{\left(np_j(\overline{X}_n)\right)^2} > \chi_{r-2,1-\alpha}^2.$$

## 3.4.3 Shapiros goodness-of-fit test

Let  $(X_1, \ldots, X_n)$  be a random sample of i.i.d. random variables  $X_i \sim F$ . The hypotheses

$$H_0: F \in \{N(\mu, \sigma^2) : \mu \in \mathbb{R}, \ \sigma^2 > 0\} \text{ vs.}$$
  
 $H_1: F \notin \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \ \sigma^2 > 0\}$ 

are to be tested. The  $\chi^2$ -tests of sections 3.4.1 - 3.4.2 are asymptotic, which makes them impractical for small sample sizes.

The following test is more suitable for testing  $H_0$ , if only a small sample is available.

Consider the order statistic  $X_{(1)}, \ldots, X_{(n)}$ , i.e.,  $X_{(1)} \leq X_{(2)} \leq \ldots \leq X_{(n)}$  and compare their correlation to the mean of the corresponding order statistic of a  $\mathcal{N}(0,1)$ -distribution. Let  $(Y_1,\ldots,Y_n)$  be a random sample of i.i.d. random variables with  $Y_1 \sim \mathcal{N}(0,1)$ . Define  $a_i \coloneqq \mathbb{E} Y_{(i)}, i = 1,\ldots,n$ . If the empirical correlation coefficient  $\rho_{aX}$  between  $(a_1,\ldots,a_n)$  and  $(X_{(1)},\ldots,X_{(n)})$  is close to 1, the random sample is normally distributed. In the following, the approach above will be formalized.

Let  $b_i$  be the expected value of the *i*-th order statistic of a  $\mathcal{N}(\mu, \sigma^2)$ -distributed random sample of i.i.d. random variables  $Z_i$ , with  $b_i = \mathbb{E} Z_{(i)}$ ,  $i = 1, \ldots, n$ . It holds  $b_i = \mu + \sigma a_i$ ,  $i = 1, \ldots, n$  and considering the correlation coefficient yields

$$\rho_{bX} = \frac{\sum_{i=1}^{n} \left(b_{i} - \overline{b}_{n}\right) \left(X_{(i)} - \overline{X}_{n}\right)}{\sqrt{\sum_{i=1}^{n} \left(b_{i} - \overline{b}_{n}\right)^{2} \sum_{i=1}^{n} \left(X_{(i)} - \overline{X}_{n}\right)^{2}}}.$$
(3.14)

Since  $\rho$  is invariant with respect to linear transformations and

$$\sum_{i=1}^{n} a_{i} = \sum_{i=1}^{n} \mathbb{E} Y_{i} = \mathbb{E} \left( \sum_{i=1}^{n} Y_{i} \right) = 0,$$

$$\rho_{bX} = \rho_{aX} = \frac{\sum_{i=1}^{n} a_{i} \left( X_{(i)} - \overline{X}_{n} \right)}{\sqrt{\sum_{i=1}^{n} a_{i}^{2} \sum_{i=1}^{n} \left( X_{i} - \overline{X}_{n} \right)^{2}}} = \frac{\sum_{i=1}^{n} a_{i} X_{(i)} - \overline{X}_{n} \sum_{i=1}^{n} a_{i}}{\sqrt{\sum_{i=1}^{n} a_{i}^{2} \sum_{i=1}^{n} \left( X_{i} - \overline{X}_{n} \right)^{2}}}$$

$$= \frac{\sum_{i=1}^{n} a_{i} X_{(i)}}{\sqrt{\sum_{i=1}^{n} a_{i}^{2} \cdot \sum_{i=1}^{n} \left( X_{i} - \overline{X}_{n} \right)^{2}}}$$

holds.

The test statistic is then given by

$$T_n = \frac{\sum_{i=1}^n a_i X_{(i)}}{\sqrt{\sum_{i=1}^n a_i^2 \sum_{i=1}^n \left(X_i - \overline{X}_n\right)^2}}$$
 (Shapiro-Francia test)

The values  $a_i$  are known and can be found in tables or by using statistic software. Note that  $|T_n| \leq 1$ .

 $H_0$  is rejected if  $T_n \leq q_{n,\alpha}$ , where  $q_{n,\alpha}$  is the  $\alpha$ -quantile of the distribution of  $T_n$ . Those quantiles can also be found in tables or by using Monte-Carlo-Simulations.

**Remark 3.4.12.** Another famous test of this kind is obtained by replacing the linear transformation  $b_i = \mu + \sigma a_i$  with another linear transformation given by

$$(a'_1,\ldots,a'_n)^{\top} = K^{-1} \cdot (a_1,\ldots,a_n),$$

where  $K = (k_{ij})_{j=1}^n$  is the covariance matrix of  $(Y_{(1)}, \dots, Y_{(n)})$  with

$$k_{ij} = \mathbb{E}\left(Y_{(i)} - a_i\right)\left(Y_{(j)} - a_j\right), \quad i, j = 1, \dots, n.$$

The constructed test is called Shapiro-Wilk test.

## 3.5 More nonparametric tests

#### 3.5.1 Binomial test

Let  $(X_1, ..., X_n)$  be a random sample of i.i.d random variables with  $X_i \sim \text{Bernoulli}(p)$ , where  $p \in [0, 1]$ . We want to test the hypotheses

$$H_0: p = p_0 \text{ vs. } H_1: p \neq p_0$$

The test statistic is given by

$$T_n = \sum_{i=1}^n X_i \underset{H_0}{\sim} \operatorname{Bin}(n, p_0),$$

and  $H_0$  is rejected if

$$T_n \notin [Bin(n, p_0)_{\alpha/2}, Bin(n, p_0)_{1-\alpha/2}],$$

where  $Bin(n,p)_{\alpha}$  is the  $\alpha$  quantile of the Bin(n,p) distribution For different  $H_0$ , like  $p \leq p_0$  ( $p \geq p_0$ ) the rejection region has to be adjusted. The quantiles  $Bin(n,p)_{\alpha}$  can also be found in tables or by using Monte-Carlo simulations. If n is sufficiently large, the quantiles can be approximated using the central limit theorem of DeMoivre-Laplace:

$$P\left(T_n \le x\right) = P\left(\frac{T_n - np_0}{\sqrt{np_0(1 - p_0)}} \le \frac{x - np_0}{\sqrt{np_0(1 - p_0)}}\right) \underset{n \to \infty}{\approx} \Phi\left(\frac{x - np_0}{\sqrt{np_0(1 - p_0)}}\right).$$

This yields

$$\begin{split} z_{\alpha} &\approx \frac{\mathrm{Bin}(n,p_0)_{\alpha} - np_0}{\sqrt{np_0(1-p_0)}} \\ \Rightarrow \mathrm{Bin}(n,p_0)_{\alpha} &\approx \sqrt{np_0(1-p_0)} \cdot z_{\alpha} + np_0 \end{split}$$

Using Poisson approximation (for  $n \to \infty, np_0 \to \lambda_0$ )

$$\operatorname{Bin}(n, p_0)_{\alpha/2} \approx \operatorname{Poisson}(\lambda_0)_{\alpha/2},$$
  
 $\operatorname{Bin}(n, p_0)_{1-\alpha/2} \approx \operatorname{Poisson}(\lambda_0)_{1-\alpha/2},$ 

holds if  $\lambda_0 = np_0$ .

**Question:** Can the symmetry of a distribution be tested by using the binomial test?

Let  $(Y_1, \ldots, Y_n)$  be a random sample of i.i.d. random variables with distribution function F. The hypotheses are

 $H_0: F$  is symmetric vs.  $H_1: F$  is not symmetric.

A symmetric distributions median is around 0. Thus the hypothesis  $H_0$  is coarsed, and

$$H'_0: F^{-1}(0,5) = 0 \text{ vs. } H'_1: F^{-1}(0,5) \neq 0$$

is tested instead. More generally, for  $\beta \in [0,1]$  we consider

$$H_0'': F^{-1}(\beta) = \gamma_\beta \text{ vs. } H_1'': F^{-1}(\beta) \neq \gamma_\beta.$$

 $H_0''$  vs.  $H_1''$  is tested by using the binomial test: Define  $X_i := I(Y_i \le \gamma_\beta)$ . Under  $H_0''$ 

$$X_i \sim \text{Bernoulli}(F(\gamma_\beta)) = \text{Bernoulli}(\beta).$$

holds. For  $a_1 = -\infty$ ,  $b_1 = \gamma_{\alpha}$ ,  $a_2 = b_1$ ,  $b_2 = +\infty$  define two disjoint classes  $(a_1, b_1]$ ,  $(a_2, b_2]$  in the sense of the  $\chi^2$ -Pearson test. The test statistic is given by

$$T_n = \sum_{i=1}^n X_i = \# \{Y_i : Y_i \le \gamma_\beta\} \sim Bin(n, \beta), \quad p = F(\gamma_\beta),$$

and the hypothesis  $F^{-1}(\beta) = \gamma_{\beta}$  is equivalent to  $H_0''': p = \beta$ . In this case, the decision rule states that  $H_0'''$  is rejected, if

$$T_n \notin \left[ \operatorname{Bin}(n,\beta)_{\alpha/2}, \operatorname{Bin}(n,\beta)_{1-\alpha/2} \right].$$

This is a test with confidence level  $\alpha$ .

#### 3.5.2 Randomness iteration tests

Sometimes in biological research a sequence of 0s and 1s is tested for its randomness or the existence of bigger clusters within those numbers. This hypothesis can be tested statistically by using the so-called *iteration tests*. Let  $(X_1,...,X_n)$  be a random sample,  $X_i \in \{0,1\}$ ,  $\sum\limits_{i=1}^n X_i = n_1$  the total number of ones,  $n_2 = n - n_1$  the total number of zeroes and  $n_1, n_2$  predetermined. An exemplary realization of  $(X_1,...,X_n)$  with  $n=18, n_1=12$  could be

$$x = (0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1).$$

The following hypotheses are to be tested:

 $H_0$ : every sequence x is equally likely vs.

 $H_1$ : There are preferred sequences (clustering).

Let

$$\Omega = \left\{ x = (x_1, \dots, x_n) : \quad x_i \in \{0, 1\}, \ i = 1, \dots, n, \sum_{i=1}^n x_i = n_1 \right\}$$

be the sample space. Then, the space  $(\Omega, \mathcal{F}, P)$  with  $\mathcal{F} = \mathcal{P}(\Omega)$ ,

$$P(x) = \frac{1}{|\Omega|} = \frac{1}{\binom{n}{n_1}}$$

is a Laplace space.

Let

$$T_n(X) = \#\{\text{Iterations in } X\} = \#\{\text{Subsequences of zeros or ones}\}\$$
  
=  $\#\{\text{Change spots from 0 to 1 or from 1 to 0}\} + 1.$ 

For x = (0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0),  $T_n(x) = 7 = 6 + 1$  holds.  $T_n(X)$  is used as a test statistic for  $H_0$  vs.  $H_1$  as follows.  $H_0$  is rejected, if T(x) is small, i.e.  $T_n(x) < F_{T_n}^{-1}(\alpha)$ . This is a test with confidence level  $\alpha$ . The question arises, how the quantiles  $F_{T_n}^{-1}$  can be calculated?

**Theorem 3.5.1.** Under  $H_0$ 

1.

$$P(T_n = k) = \begin{cases} \frac{2\binom{n_1 - 1}{i - 1}\binom{n_2 - 1}{i - 1}}{\binom{n}{n_1}}, & \text{if } k = 2i, \\ \frac{\binom{n_1 - 1}{i}\binom{n_2 - 1}{i - 1} + \binom{n_1 - 1}{i - 1}\binom{n_2 - 1}{i}}{\binom{n}{n_1}}, & \text{if } k = 2i + 1, \end{cases}$$

2.

$$\mathbb{E}\,T_n = 1 + \frac{2n_1n_2}{n},$$

3.

$$Var(T_n) = \frac{2n_1n_2(2n_1n_2 - n)}{n^2(n-1)}.$$

holds.

#### **Proof**

1. Assume that k = 2i (the uneven case works analogously). How can i clusters of ones be selected? The number of those possibilities is equal to the number of ways, how  $n_1$  particles can be distributed to i classes.

$$0|00|\dots|0|(n_1).$$

This is the number of possibilities, how i-1 partitions can be distributed on  $n_1-1$  positions, which is equal to  $\binom{n_1-1}{i-1}$ . The same holds for the zeroes.

2. Let  $Y_j = I \{X_{j-1} \neq X_j\}_{j=2,...,n}$ . Then,

$$\mathbb{E} T_n(X) = 1 + \sum_{j=2}^n \mathbb{E} Y_j = 1 + \sum_{j=2}^n P(X_{j-1} \neq X_j)$$

and the probabilities can be rewritten as

$$P(X_{j-1} \neq X_j) = \frac{2\binom{n-2}{n_1-1}}{\binom{n}{n_1}} = 2 \cdot \frac{\frac{(n-2)!}{(n-2-(n_1-1))!(n_1-1)!}}{\frac{n!}{(n-n_1)!n_1!}}$$
$$= \frac{2n_1(n-n_1)}{(n-1)n}$$
$$= \frac{2n_1n_2}{n(n-1)}.$$

Hence,

$$\mathbb{E} T_n = 1 + (n-1) \frac{2n_1 n_2}{n(n-1)} = 1 + 2 \frac{n_1 n_2}{n}.$$

Exercise 3.5.2. Proof the third assertion.

**Example 3.5.3** (Wald-Wolfowitz test). Let  $Y = (Y_1, \ldots, Y_{n_1}), Z = (Z_1, \ldots, Z_{n_2})$  be two independent random samples of i.i.d. random variables,  $Y_i \sim F$ ,  $Z_i \sim G$ .

$$H_0: F = G \text{ vs. } H_1: F \neq G.$$

is to be tested. Define  $(Y, Z) := (Y_1, \ldots, Y_{n_1}, Z_1, \ldots, Z_{n_2})$  and let  $X_i'$  be the sample variables of (Y, Z),  $i = 1, \ldots, n$ ,  $n = n_1 + n_2$ . Consider the order statistic  $X_{(i)}'$ ,  $i = 1, \ldots, n$  and set

$$X_i = \begin{cases} 1, & \text{if } X'_{(i)} = Y_j \text{ for a } j = 1, \dots, n_1, \\ 0, & \text{if } X'_{(i)} = Z_j \text{ for a } j = 1, \dots, n_2. \end{cases}$$

Under  $H_0$ , the sample values in (Y, Z) are well distributed, i.e., every combination of 0 and 1 in  $(X_1, \ldots, X_n)$  is equally likely. Thus, the randomness iteration test can be applied to test  $H_0$  vs.  $H_1$ .  $H_0$  is rejected if  $T_n(x) \leq F_{T_n}^{-1}(\alpha)$ ,  $x = (x_1, \ldots, x_n)$ .

The quantiles of  $F_{T_n}$  can be calculated directly if n is sufficiently large, since

$$\frac{n_1}{n_1+n_2}\underset{n\to\infty}{\longrightarrow} p\in(0,1)$$

implies that  $T_n$  is asymptotically normal distributed.

**Theorem 3.5.4.** Under the assumptions above

$$\lim_{n \to \infty} \frac{\mathbb{E} T_n}{n} = 2p(1-p),$$

$$\lim_{n \to \infty} \frac{1}{n} \operatorname{Var} T_n = 4p^2 (1-p)^2,$$

$$\frac{T_n - 2p(1-p)}{2\sqrt{n}p(1-p)} \xrightarrow[n \to \infty]{d} Y \sim \mathcal{N}(0,1), \quad \text{if } \frac{n_1}{n_1 + n_2} \longrightarrow p \in (0,1)$$

holds. Thus, the quantiles for  $T_n$  can be approximated for large n by

$$\alpha = P\left(T_n \le F_{T_n}^{-1}(\alpha)\right) = P\left(\frac{T_n - 2np(1-p)}{2\sqrt{n}p(1-p)} \le \frac{x - 2np(1-p)}{2\sqrt{n}p(1-p)}\right)\Big|_{x = F_{T_n}^{-1}(\alpha)}$$

$$\approx \Phi\left(\frac{F_{T_n}^{-1}(\alpha) - 2np(1-p)}{2\sqrt{n}p(1-p)}\right)$$

$$\Rightarrow z_\alpha \approx \frac{F_{T_n}^{-1}(\alpha) - 2np(1-p)}{2\sqrt{n}p(1-p)},$$

which implies

$$F_{T_{-}}^{-1}(\alpha) \approx 2np(1-p) + 2\sqrt{n}p(1-p) \cdot z_{\alpha}$$

In practice  $\hat{p} = \frac{n_1}{n_1 + n_2}$  is used for p.

## Chapter 4

## Linear Regression

In Section 6.7.3 of the lecture "Elementary probability theory and statistics", (cf. [33]) a simple form of linear regression was introduced via

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n.$$

Using matrix notation, this can be rewritten as  $Y = X\beta + \varepsilon$ , where  $Y = (Y_1, \ldots, Y_n)^{\top}$  is a random vector and

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

is a  $n \times 2$  matrix, which contains the so-called *predictor variables*  $x_i, i = 1, \ldots, n$  and is called *design matrix*. Further,  $\beta = (\beta_0, \beta_1)^{\top}$  resp.  $\varepsilon = (\varepsilon_1, \ldots, \varepsilon_n)^{\top}$  is the so-called parameter resp. error vector. With respect to the error vector, we assume that  $\varepsilon \sim \mathcal{N}(0, \mathcal{I} \cdot \sigma^2)$  is multivariate normally distributed.

In multivariate linear regression, i.e. not simple simple linear regression, an arbitrary  $(n \times m)$  design matrix

$$X = (x_{ij})_{\substack{i=1,...,n\\j=1,...,m}}$$

and a *m*-dimensional parameter vector  $\beta = (\beta_1, \dots, \beta_m)^{\top}$  are permissible for  $m \geq 2$ . That means we consider

$$Y = X\beta + \varepsilon, \tag{4.1}$$

where  $\varepsilon \sim \mathcal{N}(0, K)$  and K an arbitrary covariance matrix. In general, this choice of K can result in the errors not being independent which means that  $K \neq \text{diag } (\sigma_1^2, \dots, \sigma_n^2)$ .

The goal of this chapter is to construct estimators and tests for  $\beta$ . But before getting into detail about this, properties of the multivaritate normal distribution need to be discussed.

### 4.1 Multivariate normal distribution

In the lecture notes of "Elementary probability theory and statistics" (cf. [33]) the multivariate normal distribution was introduced in example 3.4.5 as follows:

**Definition 4.1.1.** Let  $X = (X_1, \ldots, X_n)^{\top}$  be a *n*-dimensional random vector,  $\mu \in \mathbb{R}^n$ , K a symmetric positive definite  $(n \times n)$  matrix. X is multivariate normal distribution with parameters  $\mu$  and K  $(X \sim \mathcal{N}(\mu, K))$ , if X is absolutely continuous distributed with probability density function

$$f_X(x) = \frac{1}{(2\pi)^{n/2}} \frac{1}{\sqrt{\det(K)}} \exp\left\{-\frac{1}{2} (x - \mu)^\top K^{-1} (x - \mu)\right\},$$

where  $x = (x_1, \dots, x_n)^{\top} \in \mathbb{R}^n$ .

However, this is not the only way to define the multivariate normal distribution. Thus, let us discuss three more definitions of  $\mathcal{N}(\mu, K)$ .

**Definition 4.1.2.** The random vector  $X = (X_1, \dots, X_n)^{\top}$  is multivariate normally distributed  $(X \sim \mathcal{N}(\mu, K))$  with mean vector  $\mu \in \mathbb{R}^n$  and covariance matrix K, if the characteristic function  $\varphi_X(t) = \mathbb{E} e^{i(t,X)}$ ,  $t \in \mathbb{R}^n$ , is given by

$$\varphi_X(t) = \exp\left\{it^\top \mu - \frac{1}{2}t^\top Kt\right\}, \quad t \in \mathbb{R}^n.$$

**Definition 4.1.3.** The random vector  $X = (X_1, \ldots, X_n)^{\top}$  is multivariate normally distributed  $(X \sim \mathcal{N}(\mu, K))$  with mean vector  $\mu \in \mathbb{R}^n$  and covariance matrix K, if

for all  $a \in \mathbb{R}^n$ : the random variable  $(a, X) = a^\top X \sim \mathcal{N}(a^\top \mu, a^\top K a)$ 

is a one-dimensional normally distributed random variable.

**Definition 4.1.4.** Let  $\mu \in \mathbb{R}^n$  and K be a covariance matrix. A random vector  $X = (X_1, \dots, X_n)^{\top}$  is multivariate normally distributed with mean vector  $\mu$  and covariance matrix K  $(X \sim \mathcal{N}(\mu, K))$ , if

$$X \stackrel{d}{=} \mu + C \cdot Y,$$

where C is a  $n \times m$  matrix with rank(C) = m,  $K = C \cdot C^{\top}$  and  $Y \sim \mathcal{N}(0, \mathcal{I})$  is an m-dimensional random vector with i.i.d.  $\sim \mathcal{N}(0, 1)$  coordinates.

**Remark 4.1.5.** Definition 4.1.4 is an analogue to the one-dimensional case where we know that  $Y \sim \mathcal{N}(\mu, \sigma^2)$  if and only if  $Y \stackrel{d}{=} \mu + \sigma X$  with  $X \sim \mathcal{N}(0,1)$ .

Exercise 4.1.6. Show that the function

$$f_X(x) = \frac{1}{(2\pi)^{n/2}} \frac{1}{\sqrt{\det(K)}} \exp\left\{-\frac{1}{2} (x - \mu)^\top K^{-1} (x - \mu)\right\}, x \in \mathbb{R}^n$$

from Definition 4.1.1 is indeed a probability density function.

**Lemma 4.1.7.** Let X and Y be two n-dimensional random vectors with characteristic functions

$$\varphi_X(t) = \mathbb{E} e^{i(t,X)} = \mathbb{E} e^{it^\top X}$$
$$\varphi_Y(t) = \mathbb{E} e^{i(t,Y)} = \mathbb{E} e^{it^\top Y}$$

for  $t \in \mathbb{R}^n$ . Then, it holds

1. Uniqueness theorem:

$$X \stackrel{d}{=} Y \Leftrightarrow \varphi_X(t) = \varphi_Y(t), \quad t \in \mathbb{R}^n$$

2. If X and Y are independent, then:

$$\varphi_{X+Y}(t) = \varphi_X(t) \cdot \varphi_Y(t), \quad t \in \mathbb{R}^n.$$

without proof (cf. proof of Theorem 2.1.4 (5), [32, Corollary 2.1.10].

#### Theorem 4.1.8.

- 1. The definitions 4.1.2 4.1.4 of the multivariate normal distribution are equivalent.
- 2. The definition 4.1.1 and 4.1.4 are equivalent for n = m.

#### Remark 4.1.9.

- 1. If the matrix K in Definition 4.1.4 has full rank n, then X's probability density function is given as in Definition 4.1.1. In this case it is called regular.
- 2. If  $\operatorname{rank}(K) = m < n$ , then the distribution  $\mathcal{N}(\mu, K)$  is concentrated on the m-dimensional subspace

$$\{y \in \mathbb{R}^n : y = \mu + Cx, x \in \mathbb{R}^m\}$$

by Definition 4.1.4. In this case,  $\mathcal{N}(\mu, K)$  is obviously not absolutely continuous distributed and is thus called *singular*.

**Proof** In an initial step, Definition  $4.1.3 \Leftrightarrow 4.1.2 \Leftrightarrow 4.1.4$  is proven.

1. (a) First, we show that Definitions 4.1.2 and 4.1.3 are equivalent, i.e., for the random variable X with characteristic function  $\varphi_X$  it holds that

$$\varphi_X(t) = \exp\{it^\top \mu - \frac{1}{2}t^\top K t\}$$
  
\$\implies \text{ for all } a \in \mathbb{R}^n : a^\tau X \simeq \mathcal{N}(a^\tau \mu, a^\tau K a).

Simple calculations yield

$$\varphi_{t^{\top}X}(1) = \mathbb{E} e^{it^{\top}X \cdot 1} \stackrel{\varphi_{\mathcal{N}(\mu,\sigma^2)}}{=} \exp\{it^{\top}\mu - \frac{1}{2}t^{\top}Kt\} = \varphi_X(t),$$

for all  $t \in \mathbb{R}$ . (This is called the *Procedure of Cramér-Wold*, cf. multivariate central limit theorem).

(b) Next, we show that Definitions 4.1.3 and 4.1.4 are equivalent. Using the notation  $y = C^{\top}t$  we compute

$$\varphi_{\mu+CY}(t) = \mathbb{E} e^{i(t,\mu+CY)} = \mathbb{E} e^{it^{\top}\mu+it^{\top}CY} = e^{it^{\top}\mu} \cdot \mathbb{E} e^{i(C^{\top}t,Y)}$$

$$\stackrel{Y \sim \mathcal{N}(0,\mathcal{I})}{=} e^{it^{\top}\mu} \cdot \exp\left(-\frac{1}{2}y^{\top} \cdot y\right)$$

$$= \exp\left\{it^{\top}\mu - \frac{1}{2}t^{\top}C \cdot C^{\top}t\right\}$$

$$= \exp\left\{it^{\top}\mu - \frac{1}{2}t^{\top}Kt\right\}, t \in \mathbb{R}^{n}.$$

2. It needs to be shown that for  $X \sim \mathcal{N}(\mu, K)$  in the sense of Definition 4.1.4 and  $Y \sim \mathcal{N}(\mu, K)$  in the sense of Definition 4.1.1 the relation  $\operatorname{rank}(K) = n$  implies that  $\varphi_X = \varphi_Y$ .

Definition 4.1.2 (which is equivalent to Definition 4.1.4) implies, that

$$\varphi_X(t) = \exp\left\{it^\top \mu - \frac{1}{2}t^\top Kt\right\}, \quad t \in \mathbb{R}^n,$$

$$\varphi_Y(t) = \mathbb{E} e^{it^\top Y}$$

$$= \int_{\mathbb{R}^n} e^{it^\top y} \frac{1}{(2\pi)^{n/2} \sqrt{\det K}} \cdot \exp\left\{-\frac{1}{2} \underbrace{(y-\mu)}^x {}^\top K^{-1} \underbrace{(y-\mu)}^x\right\} dy$$

$$= e^{it^\top \mu} \cdot \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{n/2} \sqrt{\det K}} \cdot \exp\left\{it^\top x - \frac{1}{2} x^\top K^{-1} x\right\} dx$$

Diagonalising  $K : \exists$  orthogonal  $(n \times n)$  matrix  $V : V^{\top} = V^{-1}$  and  $V^{\top}KV = \text{diag } (\lambda_1, \dots, \lambda_n)$ , where  $\lambda_i > 0, i = 1, \dots, n$ . By applying

the substitution x = Vz, t = Vs it holds that

$$\varphi_{Y}(t) = \frac{e^{it^{\top}\mu}}{(2\pi)^{n/2}\sqrt{\det K}} \cdot \int_{\mathbb{R}^{n}} \exp\left\{is^{\top}V^{\top}Vz - \frac{1}{2}z^{\top}V^{\top}K^{-1}Vz\right\} dz$$

$$= \frac{e^{it^{\top}\mu}}{\sqrt{(2\pi)^{n}\lambda_{1}\cdot\ldots\lambda_{n}}} \cdot \int_{\mathbb{R}^{n}} \ldots \int_{\mathbb{R}^{n}} \exp\left\{is^{\top}z - \frac{1}{2}\sum_{i=1}^{n} \frac{z_{i}^{2}}{\lambda_{i}}\right\} dz_{1} \ldots dz_{n}$$

$$= e^{it^{\top}\mu} \prod_{i=1}^{n} \int_{\mathbb{R}^{n}} \frac{1}{\sqrt{2\pi\lambda_{i}}} e^{is_{i}z_{i} - \frac{z_{i}^{2}}{(2\lambda_{i})}} dz_{i}$$

$$= e^{it^{\top}\mu} \cdot \prod_{i=1}^{n} \varphi_{\mathcal{N}(0,\lambda_{i})}(s_{i}) = e^{it^{\top}\mu} \prod_{i=1}^{n} e^{\frac{-s_{i}^{2}\lambda_{i}}{2}}$$

$$= \exp\left\{it^{\top}\mu - \frac{1}{2}s^{\top}\operatorname{diag}(\lambda_{1},\ldots,\lambda_{n})s\right\}$$

$$= \exp\left\{it^{\top}\mu - \frac{1}{2}(V^{\top}t)^{\top}V^{\top}KVV^{\top}t\right\}$$

$$= \exp\left\{it^{\top}\mu - \frac{1}{2}t^{\top}\underbrace{VV^{\top}}_{\mathcal{I}}K\underbrace{VV^{\top}}_{\mathcal{I}}t\right\}$$

$$= \exp\left\{it^{\top}\mu - \frac{1}{2}t^{\top}\underbrace{VV^{\top}}_{\mathcal{I}}K\underbrace{VV^{\top}}_{\mathcal{I}}t\right\}$$

$$= \exp\left\{it^{\top}\mu - \frac{1}{2}t^{\top}Kt\right\}, t \in \mathbb{R}^{n}.$$

## 4.1.1 Properties of the multivariate normal distribution

**Theorem 4.1.10.** Let  $X = (X_1, ..., X_n) \sim \mathcal{N}(\mu, K), \mu \in \mathbb{R}^n$ , K symmetric and positive semidefinite. Then the following properties hold:

1.  $\mu$  is the vector of expectations of X:

$$\mathbb{E} X = \mu$$
, that means  $\mathbb{E} X_i = \mu_i$ ,  $i = 1, \dots, n$ .

K is the *covariance matrix* of X:

$$K = (k_{ij})$$
, with  $k_{ij} = \text{Cov } (X_i, X_j)$ .

- 2. Every partial vector  $X' = (X_{i_1}, \ldots, X_{i_k})^{\top}$   $(1 \leq i_1 < \ldots < i_k \leq n)$  of X is also multivariate normally distributed,  $X' \sim \mathcal{N}(\mu', K')$ , where  $\mu' = (\mu_{i_1}, \ldots, \mu_{i_k})^{\top}$ ,  $K' = (k'_{jl}) = (\text{Cov}(X_{i_j}, X_{i_l}))$ ,  $j, l = 1, \ldots, k$ . In particular it holds that  $X_i \sim \mathcal{N}(\mu_i, k_{ii})$ , where  $k_{ii} = \text{Var } X_i$ ,  $i = 1, \ldots, n$ .
- 3. Two partial vectors of X are independent if and only if the corresponding elements  $k_{ij}$  of K, which represent the cross covariances, are zero, i.e.  $X' = (X_1, \ldots, X_k)^{\top}, X'' = (X_{k+1}, \ldots, X_n)$  are independent

(where the given order is chosen for the sake of simplicity) if and only if  $k_{ij} = 0$  for  $1 \le i \le k$ , j > k or i > k,  $1 \le j \le k$ , i.e.

$$K = \left(\begin{array}{c|c} K' & 0 \\ \hline 0 & K'' \end{array}\right)$$

where K' and K'' are covariance matrices of X' resp. X''.

4. Conclusion stability: If X and Y are independent, n-dimensional random vectors with  $X \sim \mathcal{N}(\mu_1, K_1)$  and  $Y \sim \mathcal{N}(\mu_2, K_2)$ , then

$$X + Y \sim \mathcal{N}(\mu_1 + \mu_2, K_1 + K_2).$$

Exercise 4.1.11. Prove Theorem 4.1.10.

**Theorem 4.1.12** (Linear transformation of  $\mathcal{N}(\mu, K)$ ). Let  $X \sim \mathcal{N}(\mu, K)$  be an n-dimensional random vector and A an  $(m \times n)$  matrix with rank $(A) = m \leq n, b \in \mathbb{R}^m$ . Then the random vector Y = AX + b is multivariate normally distributed with

$$Y \sim \mathcal{N}(A\mu + b, AKA^{\top}).$$

**Proof** Without loss of generality assume  $\mu = 0$  and b = 0, since  $\varphi_{Y-a}(t) = e^{-it^{\top}a} \cdot \varphi_Y(t)$ , for  $a = A\mu + b$ . It has to be shown that:

$$Y = AX, X \sim \mathcal{N}(0, K) \Rightarrow Y \sim \mathcal{N}(0, AKA^{\top}).$$

This can be done by calculation as follows.

$$\varphi_Y(t) = \varphi_{AX}(t) = \mathbb{E} e^{it^\top AX} = \mathbb{E} e^{i(X, A^\top t)}$$

$$\stackrel{\text{(Def. 4.1.2)}}{=} \exp\left\{-\frac{1}{2}s^\top Ks\right\} = \exp\left\{-\frac{1}{2}t^\top AKA^\top t\right\}, t \in \mathbb{R}^n$$

$$\Rightarrow Y \sim N\left(0, AKA^\top\right).$$

# 4.1.2 Linear and quadratic forms of normally distributed random variables

**Definition 4.1.13.** Let  $X = (X_1, \ldots, X_n)^{\top}$ ,  $Y = (Y_1, \ldots, Y_n)^{\top}$  be two random vectors on  $(\Omega, \mathcal{F}, P)$  and A be a symmetric, real-valued  $(n \times n)$  matrix.

1. Z = AX is called *linear form* of X with matrix A.

2.  $Z = Y^{\top}AX$  is called *bilinear form* of X and Y with matrix A. Furthermore, rewriting the vector notation yields

$$Z = \sum_{i=1}^{n} \sum_{j=1}^{n} a_{ij} X_{j} Y_{i}.$$

3. The random variable  $Z = X^{\top}AX$  (which is a bilinear form X with itself) is called *quadratic form* of X with matrix A.

**Theorem 4.1.14.** Let  $Z = Y^{\top}AX$  be a bilinear form of random vectors  $X, Y \in \mathbb{R}^n$  with respect to the symmetric matrix A. If  $\mu_X = \mathbb{E} X$ ,  $\mu_Y = \mathbb{E} Y$  and  $K_{XY} = (\text{Cov}(X_i, Y_j))_{i,j=1,...,n}$  the cross covariance matrix of X and Y, then

$$\mathbb{E} Z = \mu_Y^{\top} A \mu_X + \operatorname{trace}(AK_{XY}).$$

Proof

$$\mathbb{E} Z = \mathbb{E} \operatorname{trace}(Z) = \mathbb{E} \operatorname{trace}(Y^{\top}AX) \quad (\operatorname{since} \operatorname{trace}(AB) = \operatorname{trace}(BA))$$

$$= \mathbb{E} \operatorname{trace}(AXY^{\top}) = \operatorname{trace}(A\mathbb{E} (XY^{\top}))$$

$$= \operatorname{trace} \left( A\mathbb{E} \left( (X - \mu_X) \cdot (Y - \mu_Y)^{\top} + \mu_X Y^{\top} + X \mu_Y^{\top} - \mu_X \mu_Y^{\top} \right) \right)$$

$$= \operatorname{trace} \left( A(K_{XY} + \mu_X \mu_Y^{\top} + \mu_X \mu_Y^{\top} - \mu_X \mu_Y^{\top}) \right)$$

$$= \operatorname{trace} \left( AK_{XY} + A\mu_X \mu_Y^{\top} \right)$$

$$= \operatorname{trace} \left( AK_{XY} + \mu_X \mu_Y^{\top} \right)$$

Corollary 4.1.15. For quadratic forms it holds that

$$\mathbb{E}\left(X^{\top}AX\right) = \mu_X^{\top}A\mu_X + \operatorname{trace}(A \cdot K),$$

where  $\mu_X = \mathbb{E} X$  and K is the covariance matrix of X.

**Theorem 4.1.16** (Covariance of quadratic forms). Let  $X \sim \mathcal{N}(\mu, K)$  be an n-dimensional random vector and  $A, B \in \mathbb{R}^{n \times n}$  two symmetric matrices. Then

$$\operatorname{Cov}\left(X^{\top}AX, X^{\top}BX\right) = 4\mu^{\top}AKB\mu + 2 \cdot \operatorname{trace}(AKBK).$$

**Lemma 4.1.17** (mixed moments). Let  $Y = (Y_1, \dots, Y_n)^{\top} \sim \mathcal{N}(0, K)$  be a random vector. Then

$$\mathbb{E}(Y_i Y_j Y_k) = 0,$$

$$\mathbb{E}(Y_i Y_j Y_k Y_l) = k_{ij} \cdot k_{kl} + k_{ik} \cdot k_{jl} + k_{jk} \cdot k_{il}, \quad 1 \le i, j, k, l \le n,$$

where  $K = (k_{ij})_{i,j=1,...,n}$  is the covariance matrix of Y.

Exercise 4.1.18. Prove the Lemma.

Proof of Theorem 4.1.16.

$$\operatorname{Cov}\left(X^{\top}AX,X^{\top}BX\right) = \mathbb{E}\left(X^{\top}AX \cdot X^{\top}BX\right) \\ - \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}BX\right) \\ = Y \\ = Y \\ - \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}BX\right) \\ = Y \\ - \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}BX\right) \\ - \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}BX\right) \\ - \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}BX\right) \\ - \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) + \mathbb{E}\left(X^{\top}AX\right) \\ - \mathbb{E}\left[\left(Y^{\top}AX\right) + 2\mu^{\top}AX\right) + \mu^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \\ - \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) + \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) + \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) + \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) + \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) + \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) + \mathbb{E}\left(X^{\top}AX\right) + \mathbb{E}\left(X^{\top}AX\right) \cdot \mathbb{E}\left(X^{\top}AX\right) + \mathbb$$

Since

$$\mathbb{E}\left(Y^{\top}AY \cdot Y^{\top}BY\right) = \mathbb{E}\left(\sum_{i,j=1}^{n} a_{ij}Y_{i}Y_{j} \cdot \sum_{k,l=1}^{n} b_{kl}Y_{k}Y_{l}\right)$$

$$= \sum_{i,j,k,l=1}^{n} a_{ij}b_{kl}\mathbb{E}\left(Y_{i}Y_{j}Y_{k}Y_{l}\right)$$

$$\stackrel{\text{(Lemma 4.1.17)}}{=} \sum_{i,j,k,l=1}^{n} a_{ij}b_{kl}\left(k_{ij} \cdot k_{kl} + k_{ik} \cdot k_{jl} + k_{jk} \cdot k_{il}\right)$$

$$= \sum_{i,j=1}^{n} a_{ij} k_{ij} \cdot \sum_{k,l=1}^{n} b_{kl} \cdot k_{kl} + 2 \sum_{i,j,k,l=1}^{n} a_{ij} \cdot k_{jl} \cdot b_{lk} \cdot k_{ki}$$
$$= 2 \cdot \operatorname{trace} (AKBK) + \operatorname{trace} (AK) \cdot \operatorname{trace} (BK)$$

it holds that

Cov 
$$(X^{\top}AX, X^{\top}BX)$$
  
=  $2 \cdot \operatorname{trace}(AKBK) + \operatorname{trace}(AK) \cdot \operatorname{trace}(BK) + 4\mu^{\top}AKB\mu$   
-  $\operatorname{trace}(AK) \cdot \operatorname{trace}(BK) = 4\mu^{\top}AKB\mu + 2 \cdot \operatorname{trace}(AKBK).$ 

Corollary 4.1.19.

$$\operatorname{Var}\left(\boldsymbol{X}^{\top} \boldsymbol{A} \boldsymbol{X}\right) = 4 \boldsymbol{\mu}^{\top} \boldsymbol{A} \boldsymbol{K} \boldsymbol{A} \boldsymbol{\mu} + 2 \cdot \operatorname{trace}\left((\boldsymbol{A} \boldsymbol{K})^{2}\right)$$

**Theorem 4.1.20.** Let  $X \sim \mathcal{N}(\mu, K)$  and  $A, B \in \mathbb{R}^{n \times n}$  be two symmetric matrices. Then

$$Cov (BX, X^{\top}AX) = 2BKA\mu$$

**Proof** 

$$\operatorname{Cov} (BX, X^{\top}AX) =$$

$$\stackrel{\text{(Folgerung 4.1.15)}}{=} \mathbb{E} \left[ (BX - B\mu)(X^{\top}AX - \mu^{\top}A\mu - \operatorname{trace}(AK)) \right]$$

$$= \mathbb{E} \left[ B(X - \mu) \left( (X - \mu)^{\top}A(X - \mu) + 2\mu^{\top}AX - 2\mu^{\top}A\mu - \operatorname{trace}(AK) \right) \right],$$

Since

$$(X - \mu)^{\top} A (X - \mu) = X^{\top} A X - \mu^{\top} A X - X^{\top} A \mu + \mu^{\top} A \mu$$

and by substituting  $Z=X-\mu$  (which implies  $\mathbb{E}\,Z=0)$ 

$$\operatorname{Cov} (BX, X^{\top}AX) = \mathbb{E} \left[ BZ(Z^{\top}AZ + 2\mu^{\top}AZ - \operatorname{trace}(AK)) \right]$$

$$= \mathbb{E} (BZ \cdot Z^{\top}AZ) + 2\mathbb{E} (BZ \cdot \mu^{\top}AZ)$$

$$= B\mathbb{E} Z = 0$$

$$- \operatorname{trace}(AK) \cdot \overline{\mathbb{E}} (BZ)$$

$$= 2\mathbb{E} (BZ \cdot Z^{\top}A\mu) + \mathbb{E} (BZZ^{\top}AZ)$$

$$= 2B \underbrace{\mathbb{E} (ZZ^{\top})}_{\operatorname{Cov} X = K} A\mu + B \cdot \underbrace{\mathbb{E} (ZZ^{\top}AZ)}_{=0}$$

$$= 2BKA\mu,$$

since  $Z \sim \mathcal{N}(0, K)$  and Lemma 4.1.17 and the proof of Theorem 4.1.16.  $\square$ 

**Definition 4.1.21.** Let  $X_i \sim \mathcal{N}(\mu_i, 1), i = 1, ..., n$  be independent. Then the random variable

$$Y = X_1^2 + \ldots + X_n^2$$

is non-centered  $\chi^2_{n,\mu}$  distributed with n degrees of freedom and the non-centrality parameter

$$\mu = \sum_{i=1}^{n} \mu_i^2.$$

In Remark 2.2.6, WT&SP (cf. [32]), the moment generating function of random variables were introduced. For the proof of Theorem 4.1.23 the following uniqueness theorem will be used:

**Lemma 4.1.22** (Uniqueness theorem for moment generating functions). Let  $X_1$  and  $X_2$  be two absolutely continuous random variables with moment generating functions

$$M_{X_i}(t) = \mathbb{E} e^{tX_i}, \quad i = 1, 2,$$

which are defined on the interval (a,b). If  $f_1$  and  $f_2$  are the probability density functions of the distributions of  $X_1$  and  $X_2$ , then

$$f_1(x) = f_2(x)$$
 for almost all  $x \in \mathbb{R} \Leftrightarrow M_{X_1}(t) = M_{X_2}(t), t \in (a, b)$ .

Without proof.

**Theorem 4.1.23.** The probability density function of a  $\chi_{n,\mu}^2$  distributed random variable X (with  $n \in \mathbb{N}$  and  $\mu > 0$ ) is given by the mixture function of the density of a  $\chi_{n+2J}^2$  distribution with mixture variable  $J \sim \text{Poisson}(\mu/2)$ :

$$f_X(x) = \begin{cases} \sum_{j=0}^{\infty} e^{-\mu/2} \frac{(\mu/2)^j}{j!} \cdot \frac{e^{-x/2} x^{\frac{n+2j}{2}-1}}{\Gamma(\frac{n+2j}{2}) \cdot 2^{\frac{n+2j}{2}}}, & x \ge 0, \\ 0, & x < 0. \end{cases}$$
(4.2)

#### Proof

1. First, calculate  $M_X(t)$ ,  $X \sim \chi^2_{n,\mu}$ :

$$M_X(t) = \mathbb{E}\left(e^{tX}\right) = \mathbb{E}\exp\left\{t\sum_{i=1}^n X_i^2\right\}$$
$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \cdot \int_{-\infty}^{\infty} e^{tx_i^2} \cdot e^{-\frac{(x_i - \mu_i)^2}{2}} dx_i \quad \left(t < \frac{1}{2}, X_i \sim \mathcal{N}(\mu_i, 1)\right)$$

It holds that

$$tx_i^2 - \frac{(x_i - \mu_i)^2}{2} = \frac{1}{2}(2tx_i^2 - x_i^2 + 2x_i\mu_i - \mu_i^2)$$

$$= -\frac{1}{2}\left(x_i^2(1 - 2t) - 2x_i\mu_i + \frac{\mu_i^2}{(1 - 2t)} - \frac{\mu_i^2}{(1 - 2t)} + \mu_i^2\right)$$

$$= -\frac{1}{2}\left(\left(x_i \cdot \sqrt{1 - 2t} - \frac{\mu_i}{\sqrt{1 - 2t}}\right)^2 + \mu_i^2\left(1 - \frac{1}{1 - 2t}\right)\right)$$

$$= -\frac{1}{2}\left(\frac{(x_i(1 - 2t) - \mu_i)^2}{1 - 2t} - \mu_i^2 \cdot \frac{2t}{1 - 2t}\right)$$

Substituting

$$y_i = \frac{(x_i \cdot (1 - 2t) - \mu_i)}{\sqrt{1 - 2t}}$$

yields

$$M_X(t) = (1 - 2t)^{-\frac{n}{2}} \prod_{i=1}^n \exp\left\{\mu_i^2 \cdot \left(\frac{t}{1 - 2t}\right)\right\} \cdot \underbrace{\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{y_i^2}{2}} dy_i}_{=1}$$
$$= (1 - 2t)^{-\frac{n}{2}} \cdot \exp\left\{\frac{t}{1 - 2t} \cdot \sum_{i=1}^n \mu_i^2\right\}$$
$$= \frac{1}{(1 - 2t)^{n/2}} \cdot \exp\left\{\frac{\mu t}{1 - 2t}\right\}, \quad t < \frac{1}{2}.$$

2. Let Y be a random variable with probability density function (4.2). Calculating  $M_Y(t)$  yields

$$M_{Y}(t) = \sum_{j=0}^{\infty} e^{-\frac{\mu}{2}} \frac{(\mu/2)^{j}}{j!} \int_{0}^{\infty} e^{xt} \cdot \frac{e^{-\frac{x}{2}} \cdot x^{\frac{n+2j}{2}-1}}{\Gamma\left(\frac{n+2j}{2}\right) \cdot \frac{n+2j}{2}} dx$$

$$= M_{\chi_{n+2j}^{2}}(t) = \frac{1}{(1-2t)^{(n+2j)/2}} \operatorname{Satz} 1.1.4$$

$$= \frac{e^{-\frac{\mu}{2}}}{(1-2t)^{\frac{n}{2}}} \cdot \sum_{j=1}^{\infty} \left(\frac{\mu}{2(1-2t)}\right)^{j} \cdot \frac{1}{j!}$$

$$= \frac{1}{(1-2t)^{\frac{n}{2}}} \cdot \exp\left\{-\frac{\mu}{2} + \frac{\mu}{2(1-2t)}\right\}$$

$$= \frac{1}{(1-2t)^{\frac{n}{2}}} \cdot \exp\left\{\frac{\mu \cdot (1-(1-2t))}{2 \cdot (1-2t)}\right\}$$

$$= (1-2t)^{-\frac{n}{2}} \cdot \exp\left\{\frac{\mu t}{1-2t}\right\}$$

$$\implies M_{X}(t) = M_{Y}(t), \quad t < \frac{1}{2}$$

Using Lemma 4.1.22 implies  $f_X(x) = f_Y(x)$  for almost all  $x \in \mathbb{R}$ .

#### Remark 4.1.24.

1. Definition 4.1.21 can be rewritten as:

If  $X \sim \mathcal{N}(\vec{\mu}, \mathcal{I})$ ,  $\vec{\mu} = (\mu_1, \dots, \mu_n)^{\top}$ , then  $|X|^2 = X^{\top}X \sim \chi_{n,\mu}^2$ , where  $\mu = |\vec{\mu}|^2$ .

2. The property above can be generalized for  $X \sim \mathcal{N}(\vec{\mu}, K)$ , with a symmetric, positive definite  $(n \times n)$  matrix K:

$$X^{\top}K^{-1}X \sim \chi^2_{n,\tilde{\mu}}, \quad \text{where } \tilde{\mu} = \vec{\mu}^{\top}K^{-1}\vec{\mu},$$

and since K is positive definite, there exists a  $K^{\frac{1}{2}}$ , such that  $K=K^{\frac{1}{2}}K^{\frac{1}{2}\top}$ . Then

$$Y = K^{-\frac{1}{2}}X \sim \mathcal{N}(K^{-\frac{1}{2}}\mu, \mathcal{I}),$$

since

$$K^{-\frac{1}{2}}KK^{-\frac{1}{2}\top} = K^{-\frac{1}{2}} \cdot K^{\frac{1}{2}} \cdot K^{\frac{1}{2}\top} \cdot K^{-\frac{1}{2}\top} = \mathcal{T}$$

and thus  $Y^{\top}Y \stackrel{1}{\sim} \chi^2_{n,\tilde{\mu}}$ , with

$$\tilde{\mu} = \left(K^{-\frac{1}{2}}\vec{\mu}\right)^{\top}K^{-\frac{1}{2}}\vec{\mu} = \vec{\mu}^{\top}K^{-\frac{1}{2}\top}K^{-\frac{1}{2}}\vec{\mu} = \vec{\mu}^{\top}K^{-1}\vec{\mu}.$$

**Theorem 4.1.25.** Let  $X \sim \mathcal{N}(\mu, K)$ , where K is a symmetric, positive definite  $(n \times n)$  matrix and let A be another symmetric  $(n \times n)$  matrix with the property  $AK = (AK)^2$  (idempotence) and  $\operatorname{rank}(A) = r \leq n$ . Then:

$$X^{\top}AX \sim \chi^2_{r,\tilde{\mu}}$$
, where  $\tilde{\mu} = \mu^{\top}A\mu$ .

**Proof** A is positive semidefinite since

$$\begin{split} AK &= (AK)^2 = AK \cdot AK \quad \mid K^{-1} \\ &\Longrightarrow A = AKA \Rightarrow \ \forall x \in \mathbb{R}^n: \ x^\top Ax = x^\top AKAx \\ &= \underbrace{(Ax)}^\top K\underbrace{(Ax)}_{=y} \geq 0 \text{ because of the positive definiteness of } K. \\ &\Longrightarrow \exists H: \ \text{a} \ (n \times r) \text{ matrix with rank } (H) = r: \ A = HH^\top \end{split}$$

Thus it holds that

$$\boldsymbol{X}^{\top} \boldsymbol{A} \boldsymbol{X} = \boldsymbol{X}^{\top} \boldsymbol{H} \cdot \boldsymbol{H}^{\top} \boldsymbol{X} = (\underbrace{\boldsymbol{H}^{\top} \boldsymbol{X}}_{=\boldsymbol{Y}})^{\top} \cdot \boldsymbol{H}^{\top} \boldsymbol{X} = \boldsymbol{Y}^{\top} \boldsymbol{Y}$$

Further,  $Y \sim \mathcal{N}(H^{\top}\mu, \mathcal{I}_r)$ , since by Theorem 4.1.12  $Y \sim \mathcal{N}(H^{\top}\mu, H^{\top}KH)$  and rank (H) = r. Consequently,  $H^{\top}H$  is a regular  $(r \times r)$  matrix and

$$H^{\top}KH = (H^{\top}H)^{-1}(H^{\top}\underbrace{H \cdot H^{\top}KH \cdot (H^{\top}H)(H^{\top}H)^{-1}}_{=AKA=A}$$
$$= (H^{\top}H)^{-1}H^{\top} \cdot \underbrace{A}_{=HH^{T}} \cdot H(H^{\top}H)^{-1}$$
$$= \mathcal{I}_{r}$$

Then

$$X^{\top}AX = |Y|^2 \sim \chi_{r,\tilde{\mu}}^2 \text{ with } \tilde{\mu} = (H^{\top}\mu)^2 = \mu^{\top}H \cdot H^{\top}\mu = \mu^{\top}A\mu.$$

**Theorem 4.1.26** (Independence). Let  $X \sim \mathcal{N}(\mu, K)$  and K be a symmetric, positive semidefinite  $(n \times n)$  matrix.

- 1. Let A, B be  $(r_1 \times n)$  resp.  $(r_2 \times n)$  matrices,  $r_1, r_2 \leq n$  with  $AKB^{\top} = 0$ . Then the vectors AX and BX are independent.
- 2. Furthermore, let C be a symmetric, positive semidefinite  $(n \times n)$  matrix with the property AKC = 0. Then AX and  $X^{\top}CX$  are independent.

### Proof

1. By theorem 4.1.10, 3) it holds that AX and BX are independent, if and only if  $\varphi_{(AX,BX)}(t) = \varphi_{AX}(t) \cdot \varphi_{BX}(t)$ ,  $t = (t_1, t_2)^{\top} \in \mathbb{R}^{r_1+r_2}$ ,  $t_1 \in \mathbb{R}^{r_1}$ ,  $t_2 \in \mathbb{R}^{r_2}$ . It has to be shown that:

$$\varphi_{(AX,BX)}(t) = \mathbb{E} e^{\left(it_1^\top A + t_2^\top B\right) \cdot X} \stackrel{!}{=} \mathbb{E} e^{it_1^\top AX} \cdot \mathbb{E} e^{it_2^\top BX}.$$

It holds that

$$\begin{split} \varphi_{(AX,BX)}(t) &= \mathbb{E}e^{i\left(t_1^\top A + t_2^\top B\right) \cdot X} \\ &\stackrel{(Def.4.1.2)}{=} e^{i\left(t_1^\top A + t_2^\top B\right) \cdot \mu - \frac{1}{2} \cdot \left(t_1^\top A + t_2^\top B\right) \cdot K \cdot \left(t_1^\top A + t_2^\top B\right)^\top}. \end{split}$$

and with

$$\begin{split} \left(t_{1}^{\top}A + t_{2}^{\top}B\right) \cdot K \cdot \left(t_{1}^{\top}A + t_{2}^{\top}B\right)^{\top} \\ &= \left(t_{1}^{\top}A\right) K \left(t_{1}^{\top}A\right)^{\top} + \left(t_{1}^{\top}A\right)^{\top} K \left(t_{2}^{\top}B\right) \\ &+ \left(t_{2}^{\top}B\right) K \left(t_{1}^{\top}A\right)^{\top} + \left(t_{2}^{\top}B\right) K \left(t_{2}^{\top}B\right)^{\top} \\ &= t_{1}^{\top}AKA^{\top}t_{1} + t_{1}^{\top} \cdot \underbrace{AKB^{\top}}_{=0} \cdot t_{2} + t_{2}^{\top} \cdot \underbrace{BKA^{\top}}_{=(AKB^{\top})^{\top}=0} \cdot t_{1} + t_{2}^{\top}BKB^{\top}t_{2} \end{split}$$

we get

$$\varphi_{(AX,BX)}(t) = e^{it_1^{\top} A - \frac{1}{2} t_1^{\top} A K A^{\top} t_1} \cdot e^{it_2^{\top} B - \frac{1}{2} t_2^{\top} B K B^{\top} t_2}$$
$$= \varphi_{AX}(t_1) \cdot \varphi_{BX}(t_2), \quad t_1 \in \mathbb{R}^{r_1}, t_2 \in \mathbb{R}^{r_2}$$

2. C is symmetric, positive semidefinite  $\Longrightarrow$  There exists a  $(n \times r)$  matrix H with rank  $(H) = r \le n$  and  $C = HH^{\top}, \Longrightarrow H^{\top}H$  has rank r and is thus invertible. Then

$$X^{\top}CX = X^{\top}HH^{\top}X = (H^{\top}X)^{\top} \cdot H^{\top}X = |H^{\top}X|^2.$$

If AX and  $H^{\top}X$  are independent, then AX and  $X^{\top}CX = |H^{\top}X|^2$  are independent by the transformation theorem for random vectors. By 1) AX and  $H^{\top}X$  are independent, if  $AK(H^{\top})^{\top} = AKH = 0$ . By assumption

$$AKC = AKH \cdot H^{\top} = 0 \Longrightarrow AKH \cdot H^{\top}H = 0,$$

since  $\exists (H^{\top}H)^{-1}$ , it holds that

$$0 = AKH \cdot H^{\top}H \cdot (H^{\top}H)^{-1} = AKH \Longrightarrow AKH = 0$$
  
\Rightarrow AX and  $H^{\top}X$  are independent  
\Rightarrow AX and  $X^{\top}CX$  are independent.

# 4.2 Multivariate linear regression models with full rank

Multivariate linear regression has the form

$$Y = X\beta + \varepsilon$$
,

where  $Y = (Y_1, \dots, Y_n)^{\top}$  is the random vector of the so-called *response* variables, the design matrix

$$X = (x_{ij})_{\substack{i=1,\dots,n\\j=1,\dots,m}}$$

is deterministic and has full rank, i.e., rank  $(X) = r = m \leq n$ ,  $\beta = (\beta_1, \dots, \beta_m)^{\top}$  is the parameter vector and  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^{\top}$  is the random vector of the error terms. In our setting, the error terms fulfill  $\mathbb{E} \, \varepsilon_i = 0$ ,  $\operatorname{Var} \varepsilon_i = \sigma^2 > 0$ ,  $i \in \{1, \dots, n\}$ . The goal of this section is to estimate  $\beta$  and  $\sigma^2$ .

# 4.2.1 Method of least squares

Let the design matrix  $X = (X_1, \ldots, X_m)$  be defined by deterministic vectors  $X_j = (x_{1j}, x_{2j}, \ldots, x_{nj})^{\top}$ ,  $j = 1, \ldots, m$ , which generate the *m*-dimensional linear subspace  $L_X = \langle X_1, \ldots, X_m \rangle$ . Further, define the mean squared deviation between Y and  $X\beta$  via

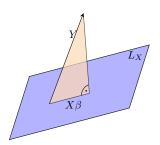
$$e(\beta) = \frac{1}{n}|Y - X\beta|^2 = \frac{1}{n}\sum_{i=1}^n (Y_i - x_{i1}\beta_1 - \dots - x_{im}\beta_m)^2.$$

Then, the ordinary least squares estimator, or *OLS estimator* for short,  $\hat{\beta}$  of  $\beta$  is defined as

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}}(e(\beta)). \tag{4.3}$$

Why does a solution  $\beta \in \mathbb{R}^m$  of the quadratic optimisation (4.3) exist? Geometrically,  $X\hat{\beta}$  can be interpreted as the orthogonal projection of the data vector Y on the linear subvector  $L_X$  as depicted in Figure 4.1. Formally, the existence of the solution will be shown by using the following theorem.

Figure 4.1: Projection on the linear subspace  $L_X$ 



**Theorem 4.2.1.** Under the above conditions, there exists an unique OLS estimator  $\hat{\beta}$ , which solves the so-called *normal equation* 

$$X^{\top}X\beta = X^{\top}Y. \tag{4.4}$$

Thus, it holds that

$$\hat{\beta} = \left( X^{\top} X \right)^{-1} X^{\top} Y.$$

**Proof** The necessary condition for the existence of the minimum is  $e'(\beta) = 0$ , that means

$$e'(\beta) = \left(\frac{\partial e(\beta)}{\partial \beta_1}, \dots, \frac{\partial e(\beta)}{\partial \beta_m}\right)^{\top} = 0.$$

It holds that

$$e'(\beta) = \frac{2}{n} \left( X^{\top} X \beta - X^{\top} Y \right)$$

 $\implies \hat{\beta}$  is a solution of the normal equation  $X^{\top}X\beta = X^{\top}Y$ . Sufficient conditions for a minimum are given, since

$$e''(\beta) = \left(\frac{\partial^2 e(\beta)}{\partial \beta_i \partial \beta_j}\right)_{i,j=1,\dots,m} = \frac{2}{n} X^{\top} X.$$

 $X^{\top}X$  is symmetric and positive definite, since X has full rank:

$$\forall y \neq 0, y \in \mathbb{R}^m : y^{\top} X^{\top} X y = (Xy)^{\top} X y = |Xy|^2 > 0$$

and  $y \neq 0 \Longrightarrow Xy \neq 0$  implies that  $e''(\beta)$  is positive definite. Thus  $X^{\top}X$  is invertible. That means,  $\hat{\beta}$  minimizes  $e(\beta)$ . The estimator  $\hat{\beta} = \left(X^{\top}X\right)^{-1}X^{\top}Y$  can be obtained, by multiplying  $\left(X^{\top}X\right)^{-1}$  to the left of the normal equation  $X^{\top}X\beta = X^{\top}Y$ .

# Example 4.2.2.

### 1. Ordinary least squares

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad m = 2, \ \beta = (\beta_1, \beta_2)^\top, \ Y = X\beta + \varepsilon$$

 $\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2)$  yields the OLS estimator from [33].

$$\hat{\beta}_2 = \frac{S_{XY}^2}{S_{YY}^2}, \quad \hat{\beta}_1 = \overline{Y}_n - \overline{X}_n \hat{\beta}_2,$$

where

$$\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i, \quad \overline{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$S_{XY}^2 = \frac{1}{n-1} \sum_{i=1}^n \left( X_i - \overline{X}_n \right) \left( Y_i - \overline{Y}_n \right)$$

$$S_{XX}^2 = \frac{1}{n-1} \sum_{i=1}^n \left( X_i - \overline{X}_n \right)^2$$

## Exercise 4.2.3. Prove that!

2. Multiple linear regression

 $Y = X\beta + \varepsilon$  with design matrix

$$X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1m} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & \cdots & x_{nm} \end{pmatrix} \text{ for } \beta = (\beta_0, \beta_1, \dots, \beta_m)^{\top}.$$

The OLS estimator  $\hat{\beta} = (X^{\top}X)^{-1}X^{\top}Y$  is obviously a linear estimator with respect to Y.

Next, let us show that  $\hat{\beta}$  is the best linear, unbiased estimator of  $\beta$  (BLUE) in the class

$$\mathcal{L} = \left\{ \tilde{\beta} = AY + b : \mathbb{E}\,\tilde{\beta} = \beta \right\}$$

of all linear unbiased estimators.

**Theorem 4.2.4** (Properties of the OLS estimator  $\hat{\beta}$ ). Let  $Y = X\beta + \varepsilon$  be a multivariate linear regression model with full rank m and error terms  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^{\top}$ , which satisfy the following conditions:

$$\mathbb{E} \varepsilon = 0$$
, Cov  $(\varepsilon_i, \varepsilon_j) = \sigma^2 \delta_{ij}$ ,  $i, j = 1, \dots, n$  for a  $\sigma^2 \in (0, \infty)$ .

Then,

- 1. the OLS estimator  $\hat{\beta} = (X^{\top}X)^{-1} X^{\top}Y$  is unbiased:  $\mathbb{E}\,\hat{\beta} = \beta$ .
- 2. Cov  $(\hat{\beta}) = \sigma^2 \left( X^{\top} X \right)^{-1}$
- 3.  $\hat{\beta}$  has minimal variance among the estimators from  $\mathcal{L}$ , i.e.,

$$\forall \tilde{\beta} \in \mathcal{L} : \operatorname{Var} \tilde{\beta}_j \geq \operatorname{Var} \hat{\beta}_j, \quad j = 1, \dots, m.$$

**Proof** 1. Let us compute

$$\begin{split} \mathbb{E}\,\hat{\beta} &= \mathbb{E}\,\left[\left(X^\top X\right)^{-1} X^\top \left(X\beta + \varepsilon\right)\right] \\ &= \left(X^\top X\right)^{-1} \cdot X^\top X \cdot \beta + \left(X^\top X\right)^{-1} X^\top \cdot \underbrace{\mathbb{E}\,\varepsilon}_{=0} \\ &= \beta \quad \forall \beta \in \mathbb{R}^m. \end{split}$$

2. For all 
$$\tilde{\beta} = AY + b \in \mathcal{L}$$
 it holds that
$$\beta = \mathbb{E} \, \tilde{\beta} = A\mathbb{E} \, Y + b = AX\beta + b \quad \forall \beta \in \mathbb{R}^m.$$

$$\Longrightarrow b = 0, \quad AX = \mathcal{I}.$$

$$\Longrightarrow \tilde{\beta} = AY = A \left( X\beta + \varepsilon \right) = AX\beta + A\varepsilon$$

$$= \beta + A\varepsilon.$$

For

$$\hat{\beta} = \underbrace{\left(X^{\top}X\right)^{-1}X^{\top}}_{-A}Y$$

it holds that

$$\operatorname{Cov} \hat{\beta} = \left( \mathbb{E} \left( \left( \hat{\beta}_{i} - \beta_{i} \right) \left( \hat{\beta}_{j} - \beta_{j} \right) \right) \right)_{i,j=1,\dots,m}$$

$$= \mathbb{E} \left( A \varepsilon \cdot \left( A \varepsilon \right)^{\top} \right) = \mathbb{E} \left( A \varepsilon \varepsilon^{\top} A^{\top} \right) = A \mathbb{E} \left( \varepsilon \varepsilon^{\top} \right) \cdot A^{\top}$$

$$= A \cdot \sigma^{2} \mathcal{I} A^{\top} = \sigma^{2} A A^{\top} = \sigma^{2} \left( X^{\top} X \right)^{-1} X^{\top} \left( \left( X^{\top} X \right)^{-1} X^{\top} \right)^{\top}$$

$$= \sigma^{2} \left( X^{\top} X \right)^{-1} X^{\top} X \left( X^{\top} X \right)^{-1} = \sigma^{2} \left( X^{\top} X \right)^{-1}.$$

3. Let  $\tilde{\beta} \in \mathcal{L}$ ,  $\tilde{\beta} = \beta + A\varepsilon$ . It has to be shown that

$$\left(\operatorname{Cov}(\tilde{\beta})\right)_{ii} = \sigma^2(AA^\top)_{ii} \ge \left(\operatorname{Cov}(\hat{\beta})\right)_{ii} = \sigma^2(X^\top X)_{ii}^{-1},$$

for i = 1, ..., m.

Let 
$$D = A - (X^{\top}X)^{-1}X^{\top}$$
, then  $A = D + (X^{\top}X)^{-1}X^{\top}$ ,  

$$AA^{\top} = \left(D + \left(X^{\top}X\right)^{-1}X^{\top}\right)\left(D^{\top} + X\left(X^{\top}X\right)^{-1\top}\right)$$

$$= DD^{\top} + \left(X^{\top}X\right)^{-1}$$
,

since

$$DX \left( X^{\top} X \right)^{-1} = \left( \underbrace{AX}_{=\mathcal{I}} - \underbrace{\left( X^{\top} X \right)^{-1} X^{\top} X}_{=\mathcal{I}} \right) \left( X^{\top} X \right)^{-1}$$
$$= 0$$

and

$$\begin{pmatrix} X^{\top}X \end{pmatrix}^{-1} X^{\top}D^{\top} = \begin{pmatrix} X^{\top}X \end{pmatrix}^{-1} X^{\top} \begin{pmatrix} A^{\top} - X \begin{pmatrix} X^{\top}X \end{pmatrix}^{-1\top} \end{pmatrix}$$

$$= \begin{pmatrix} X^{\top}X \end{pmatrix}^{-1} \left( \underbrace{(AX)^{\top}}_{=\mathcal{I}} - \underbrace{X^{\top}X \begin{pmatrix} X^{\top}X \end{pmatrix}^{-1}}_{=\mathcal{I}} \right)$$

$$= 0.$$

$$\Longrightarrow \left(AA^{\top}\right)_{ii} = \underbrace{\left(DD^{\top}\right)_{ii}}_{\geq 0} + \left(X^{\top}X\right)_{ii}^{-1} \geq \left(X^{\top}X\right)_{ii}^{-1}$$
$$\Longrightarrow \operatorname{Var} \hat{\beta}_{i} \leq \operatorname{Var} \tilde{\beta}_{i}, \quad i = 1, \dots, m.$$

**Theorem 4.2.5.** Let  $\hat{\beta}_n$  be the OLS estimator of the linear regression model from above and  $\{a_n\}_{n\in\mathbb{N}}$  be a sequence with  $a_n\neq 0, n\in\mathbb{N}, a_n\to 0 \ (n\to\infty)$ . Additionally, assume that there exists a regular  $(m\times m)$  matrix Q with

$$Q = \lim_{n \to \infty} a_n \left( X_n^\top X_n \right).$$

Then,  $\hat{\beta}_n$  is weakly consistent, i.e.,

$$\hat{\beta}_n \xrightarrow[n \to \infty]{p} \beta.$$

Proof

$$\hat{\beta}_n \xrightarrow[n \to \infty]{p} \beta \Longleftrightarrow P\left(\left|\hat{\beta}_n - \beta\right| > \varepsilon\right) \xrightarrow[n \to \infty]{0} \quad \forall \varepsilon > 0.$$

$$P\left(\left|\hat{\beta}_{n} - \beta\right| > \varepsilon\right) = P\left(\left|\hat{\beta}_{n} - \beta\right|^{2} > \varepsilon^{2}\right)$$

$$= P\left(\sum_{i=1}^{m} \left|\hat{\beta}_{in} - \beta_{i}\right|^{2} > \varepsilon^{2}\right)$$

$$\leq P\left(\bigcup_{i=1}^{m} \left\{\left|\hat{\beta}_{in} - \beta_{i}\right|^{2} > \frac{\varepsilon^{2}}{m}\right\}\right)$$

$$\leq \sum_{i=1}^{m} P\left(\left|\hat{\beta}_{in} - \beta_{i}\right| > \frac{\varepsilon}{\sqrt{m}}\right)$$

$$\stackrel{\text{Tschebyschew}}{\leq} m \sum_{i=1}^{m} \frac{\operatorname{Var} \hat{\beta}_{in}}{\varepsilon^{2}} \xrightarrow[n \to \infty]{} 0$$

$$\text{if } \operatorname{Var} \hat{\beta}_{in} \xrightarrow[n \to \infty]{} 0, \quad i = 1, \dots, m.$$

 $\operatorname{Var} \hat{\beta}_{in}$  is a diagonal element of the matrix

$$\operatorname{Cov} \, \hat{\beta}_n \overset{(Satz4.2.4)}{=} \sigma^2 \left( X_n^\top X_n \right)^{-1}.$$

If Cov  $\hat{\beta}_n \xrightarrow[n \to \infty]{} 0$  is true, then so is the theorem. Thus, let us show this convergence.

Since there exists a matrix

$$Q^{-1} = \lim_{n \to \infty} \frac{1}{a_n} \left( X_n^\top X_n \right)^{-1}$$

we can show that

$$\lim_{n \to \infty} \operatorname{Cov} \, \hat{\beta}_n = \sigma^2 \lim_{n \to \infty} \left( X_n^\top X_n \right)^{-1} = \sigma^2 \lim_{n \to \infty} a_n \cdot \frac{1}{a_n} \left( X_n^\top X_n \right)^{-1}$$
$$= 0 \cdot Q^{-1} \cdot \sigma^2 = 0.$$

# **4.2.2** Estimator of the variance $\sigma^2$

Introduce the estimator  $\hat{\sigma}^2$  for the variance  $\sigma^2$  of the error terms  $\varepsilon_i$  as follows:

$$\hat{\sigma}^2 = \frac{1}{n-m} \left| Y - X \hat{\beta} \right|^2. \tag{4.5}$$

This is a generalized version of the variance estimator from the simple linear regression, which has already been introduced in [33].

**Theorem 4.2.6** (Expectation). The variance estimator

$$\hat{\sigma}^2 = \frac{1}{n-m} \left| Y - X \hat{\beta} \right|^2$$

is unbiased. That means,

$$\mathbb{E}\,\hat{\sigma}^2 = \sigma^2$$

Proof

$$\hat{\sigma}^2 = \frac{1}{n-m} \left( Y - X \hat{\beta} \right)^\top \left( Y - X \hat{\beta} \right)$$

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

$$= \frac{1}{n-m} (DY)^\top DY$$

where  $D = \mathcal{I} - X(X^{\top}X)^{-1}X^{\top}$  is a  $(n \times n)$  matrix. Then,

$$\hat{\sigma}^2 = \frac{1}{n-m} Y^\top D^\top DY = \frac{1}{n-m} Y^\top D^2 Y = \frac{1}{n-m} Y^\top DY,$$

if  $D^{\top}=D$  and  $D^2=D,$  i.e. D is symmetric and idempotent. Indeed, it holds that:

$$D^{\top} = \mathcal{I} - \left(X^{\top}\right)^{\top} \left(X^{\top}X\right)^{\top^{-1}} X^{\top} = \mathcal{I} - X \left(X^{\top}X\right)^{-1} X^{\top} = D.$$

$$D^{2} = \left(\mathcal{I} - X(X^{\top}X)^{-1}X^{T}\right) \left(\mathcal{I} - X \left(X^{\top}X\right)^{-1}X^{\top}\right)$$

$$= \mathcal{I} - 2X \left(X^{\top}X\right)^{-1} X^{\top} + X \left(X^{\top}X\right)^{-1} X^{\top}X \left(X^{\top}X\right)^{-1} X^{\top}$$

$$= \mathcal{I} - X \left(X^{\top}X\right)^{-1} X^{\top} = D.$$

Furthermore it holds that

$$\hat{\sigma}^2 = \frac{1}{n-m} \cdot \operatorname{trace}\left(Y^\top D Y\right) = \frac{1}{n-m} \cdot \operatorname{trace}\left(D Y Y^\top\right)$$

$$\Longrightarrow \mathbb{E}\,\hat{\sigma}^2 = \frac{1}{n-m} \cdot \operatorname{trace}\left(D \mathbb{E}\left(Y Y^\top\right)\right) = \frac{\sigma^2}{n-m} \cdot \operatorname{trace}\left(D\right),$$

since

$$\operatorname{trace}\left(D \cdot \mathbb{E}\left(YY^{\top}\right)\right) = \\ = \operatorname{trace}\left(D(X\beta)(X\beta)^{\top} + DX\beta \underbrace{\mathbb{E}\,\varepsilon^{\top}}_{=0} + D\underbrace{\mathbb{E}\,\varepsilon}_{=0}(X\beta)^{\top} + D \cdot \underbrace{\mathbb{E}\,\varepsilon\varepsilon^{\top}}_{=0}\right) \\ = \operatorname{Cov}\,\varepsilon = \sigma^{2} \cdot \mathcal{I}$$

and

$$DX = \left(\mathcal{I} - X \left(X^{\top} X\right)^{-1} X^{T}\right) X$$
$$= X - X \left(X^{\top} X\right)^{-1} X^{\top} X = X - X = 0.$$

Now it needs to be shown that trace(D) = n - m:

$$\operatorname{trace}(D) = \operatorname{trace}\left(\mathcal{I} - X\left(X^{\top}X\right)^{-1}X^{\top}\right)$$
$$= \operatorname{trace}\left(\mathcal{I}\right) - \operatorname{trace}\left(X\left(X^{\top}X\right)^{-1}X^{\top}\right)$$
$$= n - \operatorname{trace}\left(\underbrace{X^{\top}X \cdot \left(X^{\top}X\right)^{-1}}_{\text{eine } (m \times m)\text{-Matrix}}\right) = n - m.$$

# **4.2.3** Maximum likelihood estimator for $\beta$ and $\sigma^2$

In order to construct a maximum likelihood estimator for  $\beta$  and  $\sigma^2$  resp. the distributional properties of the OLS estimators  $\hat{\beta}$  and  $\hat{\sigma}^2$ , the distribution of  $\varepsilon$  resp. Y has to be specified. In the following, normally distributed i.i.d error terms are assumed, i.e.,

$$\varepsilon \sim N\left(0, \sigma^2 \mathcal{I}\right), \quad \sigma^2 > 0.$$

Clearly, this implies

$$Y \sim N\left(X\beta, \sigma^2 \mathcal{I}\right)$$
.

What do the distributions of the OLS estimators  $\hat{\beta}$  and  $\hat{\sigma}^2$  look like? Since  $\hat{\beta}$  is linearly dependent of Y, unbiased and Cov  $\hat{\beta} = \hat{\sigma}^2 \left( X^\top X \right)^{-1}$ , simple calculations yield

$$\hat{\beta} \sim N\left(\beta, \sigma^2 \left(X^{\top} X\right)^{-1}\right).$$

In the following, the maximum likelihood estimator for  $\beta$  and  $\sigma^2$ , namely  $\tilde{\beta}$  and  $\tilde{\sigma}^2$  are calculated. Once they have been calculated, it can be seen that they are closely related to the OLE estimator.

$$\tilde{\beta} = \hat{\beta},$$

$$\tilde{\sigma}^2 = \frac{n-m}{n}\hat{\sigma}^2.$$

Consider the Likelihood function of Y:

$$L(y, \beta, \sigma^2) = f_Y(y) = \frac{1}{\left(\sqrt{2\pi}\sigma\right)^n} \cdot \exp\left\{-\frac{1}{2\sigma^2} \left(y - X\beta\right)^\top \left(y - X\beta\right)\right\}$$

and the Log likelihood function

$$\log L(y,\beta,\sigma^2) = -\frac{n}{2}\log\left(2\pi\right) \underbrace{-\frac{n}{2}\log\left(\sigma^2\right) - \frac{1}{2\sigma^2}\left|y - X\beta\right|^2}_{:=q}.$$

The maximum likelihood estimators are then given by

$$(\tilde{\beta}, \tilde{\sigma}^2) = \underset{\beta \in \mathbb{R}^m, \, \sigma^2 > 0}{\operatorname{argmax}} \log L(y, \beta, \sigma^2),$$

if they exist.

**Theorem 4.2.7** (Maximum likelihood estimation of  $\tilde{\beta}$  and  $\tilde{\sigma}^2$ ). There exist unique maximum likelihood estimators for  $\beta$  and  $\sigma^2$  which are given by

$$\tilde{\beta} = \hat{\beta} = \left(X^{\top}X\right)^{-1}X^{\top}Y$$

$$\tilde{\sigma}^2 = \frac{n-m}{n}\hat{\sigma}^2 = \frac{1}{n}\left|Y - X\tilde{\beta}\right|^2.$$

**Proof** Fix  $\sigma^2 > 0$  and find

$$\tilde{\beta} = \underset{\beta \in \mathbb{R}^m}{\operatorname{argmax}} \log L(Y, \beta, \sigma^2) = \underset{\beta \in \mathbb{R}^m}{\operatorname{argmin}} |Y - X\beta|^2,$$

which implies that  $\tilde{\beta}$  coincides with the known OLS estimator  $\hat{\beta} = (X^{\top}X)^{-1}X^{\top}Y$  and does not depend on  $\sigma^2$ . Therefore, we can compute

$$\tilde{\sigma}^2 = \operatorname*{argmax}_{\sigma^2 > 0} \log L\left(Y, \tilde{\beta}, \sigma^2\right) = \operatorname*{argmax}_{\sigma^2 > 0} g(\sigma^2).$$

It holds that

$$g\left(\sigma^{2}\right) \underset{\sigma^{2} \to +\infty}{\longrightarrow} -\infty, \quad g\left(\sigma^{2}\right) \underset{\sigma^{2} \to 0}{\longrightarrow} -\infty,$$

since  $|Y-X\beta|^2\neq 0$  because  $Y\sim N\left(X\beta,\sigma^2\mathcal{I}\right)\in\{Xy:y\in\mathbb{R}^m\}$  with probability zero. Since

$$g'(\sigma^2) = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{|Y - X\beta|}{2(\sigma^2)^2} = 0, \quad \tilde{\sigma}^2 = \frac{1}{n} |Y - X\tilde{\beta}|^2$$

maximizes  $g(\sigma^2)$ , that means  $\tilde{\sigma}^2$  is a maximum likelihood estimator for  $\sigma^2$ .

Theorem 4.2.8. Under the assumptions above it holds that

- 1.  $\mathbb{E}\,\tilde{\sigma}^2 = \frac{n-m}{n}\sigma^2$ , that means  $\tilde{\sigma}^2$  is biased; but it is asymptotically unbiased.
- 2.  $\frac{n}{\sigma^2}\tilde{\sigma}^2 \sim \chi_{n-m}^2$ ,  $\frac{n-m}{\sigma^2}\hat{\sigma}^2 \sim \chi_{n-m}^2$ .

#### Proof

- 1. Trivial (similar to the proof of Theorem 4.2.6)
- 2. Only the assertion for  $\hat{\sigma}^2$  is shown:

$$\begin{split} \frac{n-m}{\sigma^2} \hat{\sigma}^2 &= \frac{1}{\sigma^2} \left| Y - X \hat{\beta} \right|^2 \\ &= \frac{1}{\sigma^2} Y^\top \underbrace{\mathcal{D}}_{=D^2} Y \quad \text{(by the proof of Theorem 4.2.6)} \\ &= \frac{1}{\sigma^2} (DY)^\top DY = \frac{1}{\sigma^2} (\underbrace{\mathcal{D}(X \, \beta + \varepsilon)})^\top \cdot \underbrace{\mathcal{D}(X \, \beta + \varepsilon)}_{=0} \\ &= \frac{1}{\sigma^2} (D\varepsilon)^\top D\varepsilon = \left(\frac{\varepsilon}{\sigma}\right) D\left(\frac{\varepsilon}{\sigma}\right), \end{split}$$

where

$$\left(\frac{\varepsilon}{\sigma}\right) \sim N\left(0, \mathcal{I}\right).$$

By Theorem 4.1.25 it holds that

$$\frac{\varepsilon^{\top}}{\sigma} D \frac{\varepsilon}{\sigma} \sim \chi_r^2,$$

where  $r = \operatorname{rank}(D)$ , since  $D\mathcal{I} = D$  is idempotent. If r = n - m, then  $(n-m)\hat{\sigma}^2 \sim \chi^2_{n-m}$ . It needs to be shown that  $\operatorname{rank}(D) = r = n - m$ . From linear algebra it is known that  $\operatorname{rank}(D) = n - \dim(\operatorname{Kern}(D))$ . Now  $\operatorname{Kern}(D) = \{Xx : x \in \mathbb{R}^m\}$  and thus  $\dim(\operatorname{Kern}(D)) = m$ , since  $\operatorname{rank}(X) = m$ . It holds that  $\{Xx : x \in \mathbb{R}^n\} \subseteq \operatorname{Kern}(D)$ , since

$$DX = (\mathcal{I} - X(X^{\top}X)^{-1}X^{\top})X = X - (X^{\top}X)^{-1}X^{\top}X = 0.$$

and Kern 
$$(D) \subseteq \{Xx : x \in \mathbb{R}^m\}$$
, since

$$\forall y \in \text{Kern } (D): \quad Dy = 0 \iff (\mathcal{I} - X(X^{\top}X)^{-1}X^{\top})y = 0$$
$$\iff y = X \cdot \underbrace{(X^{\top}X)^{-1}X^{\top}y}_{x} = Xx \in \{Xx : x \in \mathbb{R}^{m}\}.$$

**Theorem 4.2.9.** Let  $Y = X\beta + \varepsilon$  be a multivariate regression model with  $Y = (Y_1, \dots, Y_n)^\top$ , design matrix X with rank  $(X) = m, \beta = (\beta_1, \dots, \beta_m)^\top$ ,  $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathcal{I})$ . Then, the estimators  $\hat{\beta} = (X^\top X)^{-1} X^\top Y$  for  $\beta$  resp.  $\hat{\sigma}^2 = \frac{1}{n-m}|Y - X\hat{\beta}|^2$  for  $\sigma^2$  are independent.

**Proof** In this proof Theorem 4.1.26 is used. In order to do so,  $\hat{\beta}$  has to be expressed as a linear and  $\hat{\sigma}^2$  as quadratic form of  $\varepsilon$ . It has been shown in the proofs of Theorem 4.2.4 and 4.2.8 that

$$\hat{\beta} = \beta + \underbrace{(X^{\top}X)^{-1}X^{\top}}_{=A} \varepsilon,$$

$$\hat{\sigma}^2 = \frac{1}{n-m} \varepsilon^{\top} D \varepsilon, \text{ where } D = \mathcal{I} - X(X^{\top}X)^{-1}X^{\top}.$$

Moreover it holds that AD = 0, by the proof of Theorem 4.2.6

$$(AD)^{\top} = D^{\top}A^{\top} = \underbrace{D \cdot X}_{=0} ((X^{\top}X)^{-1})^{\top} = 0.$$

Since  $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathcal{I})$ , it holds that

$$A\sigma^2 \mathcal{T}D = 0.$$

Thus, the assumptions of Theorem 4.1.26 are satisfied and  $\hat{\beta}$  and  $\hat{\sigma}^2$  are independent.

#### 4.2.4 Tests for regression parameters

In this section the hypotheses

$$H_0: \beta = \beta_0 \text{ vs. } H_1: \beta \neq \beta_0$$

are tested for a  $\beta_0 \in \mathbb{R}^m$ . In order to do so, define the test statistic

$$T = \frac{\left(\hat{\beta} - \beta_0\right)^{\top} X^{\top} X \left(\hat{\beta} - \beta_0\right)}{m\hat{\sigma}^2}.$$

Theorem 4.2.11 implies that under  $H_0$ 

$$T \sim F_{m,n-m}$$
.

Thus,  $H_0$  is rejected, if  $T > F_{m,n-m,1-\alpha}$ , where  $F_{m,n-m,1-\alpha}$  is the  $(1-\alpha)$  quantile of the  $F_{m,n-m}$  distribution. This is a test with confidence level  $\alpha \in (0,1)$ .

Special case: The case  $\beta_0 = 0$  describes test for connectivity; that means it is tested, whether  $\beta_1, \ldots, \beta_m$  are relevant for describing the data Y.

#### Remark 4.2.10.

1. How can we test that the test statistic T can actually distinguish  $H_0$  from  $H_1$ ? Introduce

$$\tilde{Y} = Y - X\hat{\beta} =: Y - \hat{Y}.$$

Then,

$$\hat{\sigma}^2 = \frac{1}{n-m} \left| \tilde{Y} \right|^2$$

and  $\tilde{Y}$  is the vector of *residuals*.

Without loss of generality assume  $\beta_0 = 0$ . If  $H_0$  is false, then  $\beta \neq 0$ , and

$$|X\beta|^2 = (X\beta)^\top X\beta = \beta^\top X^\top X\beta > 0,$$

since X has full rank. This implies that  $H_0$  has to be rejected, if

$$\left| \hat{Y} \right|^2 = \left| X \hat{\beta} \right|^2 = \hat{\beta}^\top X^\top X \hat{\beta} \gg 0.$$

In the test statistic  $|X\hat{\beta}|^2$  the variation of the estimation of  $\beta$  is not considered. In order to do so, a new test statistic T can be defined by dividing  $|X\hat{\beta}|^2$  by  $\hat{\sigma}^2$ , i.e.

$$T = \frac{\hat{\beta}^{\top} X^{\top} X \hat{\beta}}{m \cdot \hat{\sigma}^2} = \frac{\left| \hat{Y} \right|^2}{\frac{m}{n-m} \left| Y - \hat{Y} \right|^2}.$$

The Pythagorean theorem implies

$$|Y|^2 = \left|\tilde{Y}\right|^2 + \left|\hat{Y}\right|^2.$$

Then, under  $H_0$ 

$$\mathbb{E}\,|\hat{Y}|^2 = \mathbb{E}\,|Y|^2 - \mathbb{E}\,|Y - \hat{Y}|^2 = n\sigma^2 - \mathbb{E}\,|\tilde{Y}|^2$$

holds and thus

$$\frac{\mathbb{E}\,|\hat{Y}|^2}{\mathbb{E}\,\left(\frac{m}{n-m}\left|\tilde{Y}\right|^2\right)}\stackrel{(H_0)}{=}\frac{n\sigma^2-\mathbb{E}\,|\tilde{Y}|^2}{\frac{m}{n-m}\mathbb{E}\,|\tilde{Y}|^2}=\frac{n-m}{m}\left(\frac{n\sigma^2}{\mathbb{E}\,|\tilde{Y}|^2}-1\right),$$

where we have used that  $\mathbb{E}|Y|^2 = \mathbb{E}\left(Y^\top Y\right) = \sigma^2 n$  and  $Y \sim \mathcal{N}(0, \sigma^2 \mathcal{I})$ . Consequently, the test statistic T is sensible with respect to variations of  $H_0$ .

#### 2. The term

$$\left| \tilde{Y} \right|^2 = \left| Y - \hat{Y} \right|^2$$

is called residual distribution. Now the coefficient of determination  $R^2$  as introduced in [33] can be generalized as

$$R^{2} = 1 - \frac{|\tilde{Y}|^{2}}{\left|Y - \overline{Y}_{n} \cdot e\right|^{2}},$$

where 
$$e = (1, ..., 1)^{\top}, \overline{Y}_n = \frac{1}{n} \sum_{i=1}^{n} Y_i$$
.

**Theorem 4.2.11.** Under  $H_0: \beta = \beta_0$  it holds that

$$T = \frac{\left(\hat{\beta} - \beta_0\right)^{\top} X^{\top} X \left(\hat{\beta} - \beta_0\right)}{m\hat{\sigma}^2} \sim F_{m,n-m}.$$

**Proof** It holds that

$$\hat{\beta} \sim N \left( \beta_0, \sigma^2 \left( X^\top X \right)^{-1} \right)$$

$$\Longrightarrow \hat{\beta} - \beta_0 \sim N \left( 0, \underbrace{\sigma^2 (X^\top X)^{-1}}_{:=K} \right).$$

If  $A = \frac{X^{\top}X}{\sigma^2}$ , then  $AK = \mathcal{I}$  is idempotent. Then by Theorem 4.1.25

$$(\hat{\beta} - \beta_0)^{\top} A (\hat{\beta} - \beta_0) \stackrel{H_0}{\sim} \chi_m^2$$

Note that under  $H_1$  the distribution of  $(\hat{\beta} - \beta_0)^{\top} A(\hat{\beta} - \beta_0)$  does not follow a centered  $\chi^2$  distribution.

Furthermore, it holds that

$$\frac{n-m}{\sigma^2}\hat{\sigma}^2 \sim \chi_{n-m}^2.$$

Also, Theorem 4.2.9 implies the independence of  $(\hat{\beta} - \beta_0)^{\top} A(\hat{\beta} - \beta_0)$  and  $\frac{n-m}{\sigma^2} \hat{\sigma}^2$ . Therefore,

$$T = \frac{(\hat{\beta} - \beta_0)^{\top} (X^{\top} X)(\hat{\beta} - \beta_0)/m}{(n - m)\hat{\sigma}^2/(n - m)} \sim F_{m, n - m}$$

by definition of the F distribution.

Now, let us test the relevance of the parameters  $\beta_j$ , i.e., we test

$$H_0: \beta_j = \beta_{0j} \text{ vs. } H_1: \beta_j \neq \beta_{0j}.$$

**Theorem 4.2.12.** Under  $H_0: \beta_j = \beta_{0j}$  it holds that

$$T_j = \frac{\hat{\beta}_j - \beta_{0j}}{\hat{\sigma}\sqrt{x^{jj}}} \sim t_{n-m},$$

where  $x^{jj}$  is the j-th diagonal entry of the matrix  $(X^{\top}X)^{-1}$ .

**Proof**  $\hat{\beta} \stackrel{H_0}{\sim} \mathcal{N}(\beta_0, \sigma^2(X^\top X)^{-1})$  implies  $\hat{\beta}_j \stackrel{H_0}{\sim} \mathcal{N}(\beta_{0j}, \sigma^2 x^{jj})$  and thus  $\hat{\beta}_j - \beta_{0j} \sim \mathcal{N}(0, \sigma^2 x^{jj})$ . Consequently,  $A := \frac{\hat{\beta}_j - \beta_{0j}}{\sigma \sqrt{x^{jj}}} \sim \mathcal{N}(0, 1)$ . Furthermore, it holds that  $B := \frac{(n-m)\hat{\sigma}^2}{\sigma^2} \stackrel{H_0}{\sim} \chi^2_{n-m}$ , and by Theorem 4.2.9 the statistics A and B are independent which implies that

$$T_j = \frac{\frac{\hat{\beta}_j - \beta_{0j}}{\sigma \sqrt{x^{jj}}}}{\sqrt{\frac{(n-m)\hat{\sigma}^2}{(n-m)\sigma^2}}} \sim t_{n-m}.$$

Thus, a test for  $H_0: \beta_j = \beta_{j0}$  vs.  $H_1:$  not  $H_0$ , with confidence level  $\alpha$  can be constructed using test statistic T by rejected the null hypothesis, if  $|T| > t_{n-m,1-\alpha/2}$ .

Next, let us test the hypothesis

$$H_0: \beta_{i_1} = \beta_{0i_1}, \dots, \beta_{i_l} = \beta_{0i_l} \text{ vs. } H_1: \exists i \in \{1, \dots, l\}: \beta_{i_i} \neq \beta_{0i_l}.$$

**Exercise 4.2.13.** Show that under  $H_0$  the following assertion holds:

$$T = \frac{(\hat{\beta}' - \beta_0')^\top K'(\hat{\beta}' - \beta_0')}{l\hat{\sigma}^2} \sim F_{l,n-m},$$

where

$$\hat{\beta}' = (\hat{\beta}_{j_1}, \dots, \hat{\beta}_{j_l}),$$

$$\beta'_0 = (\beta_{0j_1}, \dots, \beta_{0j_l}),$$

$$K' = \begin{pmatrix} x^{j_1j_1} & \cdots & x^{j_1j_l} \\ \vdots & \vdots & \vdots \\ x^{j_lj_1} & \cdots & x^{j_lj_l} \end{pmatrix}^{-1}.$$

Construct the corresponding F test!

#### Test for linear combinations of parameters

Let us consider

$$H_0: H\beta = c \text{ vs. } H_1: H\beta \neq c,$$

where H is a  $(r \times m)$  matrix and  $c \in \mathbb{R}^r$ .

**Theorem 4.2.14.** Under  $H_0$  it holds that

$$T = \frac{(H\hat{\beta} - c)^{\top} (H(X^{\top}X)^{-1}H^{\top})^{-1} (H\hat{\beta} - c)}{r\hat{\sigma^2}} \sim F_{r,n-m}.$$

Thus  $H_0: H\beta = c$  is rejected, if  $T > F_{r,n-m,1-\alpha}$ .

Exercise 4.2.15. Prove Theorem 4.2.14!

#### 4.2.5Confidence region

1. Confidence interval for  $\beta_i$ 

In Theorem 4.2.12 it has been shown that

$$\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma} \cdot \sqrt{x^{jj}}} \sim t_{n-m},$$

where  $(X^{\top}X)^{-1} = (x^{ij})_{i,j=1,\dots,m}$ . By using the standard methodology, the following confidence interval for  $\beta_i$  with confidence level  $1-\alpha$  can be constructed as follows

$$P\left(\hat{\beta}_j - t_{n-m,1-\alpha/2} \cdot \hat{\sigma} \sqrt{x^{jj}} \le \beta_j \le \hat{\beta}_j + t_{n-m,1-\alpha/2} \cdot \hat{\sigma} \sqrt{x^{jj}}\right) = 1 - \alpha.$$

2. Simultaneous confidence region for  $\beta = (\beta_1, \dots, \beta_m)^{\top}$ 

From the Bonferroni inequality it is known that

$$P\left(\bigcap_{j=1}^{m} A_j\right) \ge \sum_{j=1}^{m} P(A_j) - (m-1),$$

for sets  $A_1, \ldots, A_m$ . Now, using the sets

$$A_j := \left\{ \beta_j \in \left[ \hat{\beta}_j - t_{n-m,1-\alpha/(2m)} \cdot \hat{\sigma} \sqrt{x^{jj}}, \hat{\beta}_j + t_{n-m,1-\alpha/(2m)} \cdot \hat{\sigma} \sqrt{x^{jj}} \right] \right\}$$

yields

$$P(A_j, j = 1, ..., m)$$
  
 $\geq \sum_{j=1}^{m} P(A_j) - (m-1) = m \cdot \left(1 - \frac{\alpha}{m}\right) - m + 1 = 1 - \alpha.$ 

This implies that

$$\left\{ \beta = (\beta_1, \dots, \beta_m)^\top : \beta_j \in A_j \right\}$$

is a simultaneous confidence region for  $\beta$  with confidence level  $1 - \alpha$ .

3. Confidence ellipsoid for  $\beta$ .

In Theorem 4.2.11 it has been shown that

$$T = \frac{(\hat{\beta} - \beta)^{\top} (X^{\top} X)(\hat{\beta} - \beta)}{m \hat{\sigma}^2} \sim F_{m, n-m}.$$

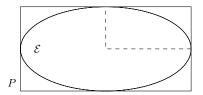
This implies that

$$P\left(T \leq F_{m,n-m,1-\alpha}\right) = 1 - \alpha \quad \text{and}$$

$$\mathcal{E} = \left\{ \beta \in \mathbb{R}^m : \frac{(\hat{\beta} - \beta)^\top (X^\top X)(\hat{\beta} - \beta)}{m\hat{\sigma}^2} \leq F_{m,n-m,1-\alpha} \right\}$$

is a confidence ellipsoid with confidence level  $1 - \alpha$  (see Figure 4.2).

Figure 4.2: Confidence ellipsoid



Since an ellipsoid can be embedded in the minimal parallelepiped P, such that the length of the sides of P are  $2 \times$  length of the half-axes of  $\mathcal{E}$ , the following simultaneous confidence region for  $\beta = (\beta_1, \dots, \beta_m)^{\top}$  can be constructed:

$$P = \left\{ \beta : \hat{\beta}_j - \hat{\sigma} \sqrt{mx^{jj}} F_{m,n-m,1-\alpha} \le \beta_j \le \hat{\beta}_j + \hat{\sigma} \sqrt{mx^{jj}} F_{m,n-m,1-\alpha} \right\}$$

$$j = 1, \dots, m$$

4. Confidence interval for the expected target value  $x_{01}\beta_1 + \ldots + x_{0m}\beta_m$ . Let  $Y_0 = x_{01}\beta_1 + \ldots + x_{0m}\beta_m + \varepsilon_0$  be a new target value with  $\mathbb{E} \, \varepsilon_0 = 0$ . Then

$$\mathbb{E} Y_0 = \sum_{i=1}^n x_{0i} \beta_i.$$

In the following, a confidence interval for  $\mathbb{E} Y_0$  is constructed. In oder to do so, the proof idea of Theorem 4.2.12 combined with Theorem 4.2.14 with  $H = (x_{01}, \dots, x_{0m}) = x_0^{\top}, r = 1$  is used. Then

$$T = \frac{\sum_{i=1}^{m} \hat{\beta}_{i} x_{0i} - \sum_{i=1}^{m} \beta_{i} x_{0i}}{\hat{\sigma} \sqrt{x_{0}^{\top} (X^{\top} X)^{-1} x_{0}}} \sim t_{n-m}.$$

Thus

$$\left\{ \beta = (\beta_1, \dots, \beta_m)^\top : \sum_{i=1}^m x_{0i} \hat{\beta}_i - \hat{\sigma} \sqrt{x_0^\top (X^\top X)^{-1} x_0} \cdot t_{n-m,1-\alpha/2} \right. \\
\leq \sum_{i=1}^m x_{0i} \beta_i \leq \sum_{i=1}^m x_{0i} \hat{\beta}_i + \hat{\sigma} \sqrt{x_0^\top (X^\top X)^{-1} x_0} \cdot t_{n-m,1-\alpha/2} \right\}$$

is a confidence interval  $\sum_{i=1}^{m} x_{0i}\beta_i$  with confidence level  $1-\alpha$ .

5. Forecast interval for the target variable  $Y_0$ .

For  $Y_0 = \sum_{i=1}^m x_{0i}\beta_i + \varepsilon_0$  with  $\varepsilon_0 \sim \mathcal{N}(0, \sigma^2)$  independent of  $\varepsilon_1, \dots, \varepsilon_n$ , it holds that

$$x_0^{\top} \hat{\beta} - Y_0 \sim \mathcal{N}(0, \sigma^2 (1 + x_0^{\top} (X^{\top} X)^{-1} x_0))$$

$$\Longrightarrow \frac{x_0^{\top} \hat{\beta} - Y_0}{\sigma \sqrt{1 + x_0^{\top} (X^{\top} X)^{-1} x_0}} \sim \mathcal{N}(0, 1)$$

$$\Longrightarrow \frac{x_0^{\top} \hat{\beta} - Y_0}{\hat{\sigma} \sqrt{1 + x_0^{\top} (X^{\top} X)^{-1} x_0}} \sim t_{n-m},$$

Thus,

$$\left(x_0^\top \hat{\beta} - c, \ x_0^\top \hat{\beta} + c\right)$$
 with  $c = \hat{\sigma} \sqrt{1 + x_0^\top (X^\top X)^{-1} \cdot x_0} \cdot t_{n-m,1-\alpha/2}$ 

is a forecast interval for the target variable  $Y_0$  with confidence level  $1-\alpha$ .

6. Confidence band for the regression plane  $y = \beta_1 + \sum_{i=2}^{m} x_i \beta_i$  in the multiple regression model.

Let  $Y = X\beta + \varepsilon$ , where

$$X = \begin{pmatrix} 1 & x_{12} & \cdots & x_{1m} \\ 1 & x_{22} & \cdots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n2} & \cdots & x_{nm} \end{pmatrix} \quad \text{and } \varepsilon \sim \mathcal{N}(0, \sigma^2 \cdot \mathcal{I}).$$

The goal is to construct a confidence band B(x) for y. It holds that

$$P\left(y = \beta_1 + \sum_{i=2}^{m} \beta_i x_i \in B(x)\right) = 1 - \alpha \quad \forall x \in \mathbb{R}_1^{m-1},$$

where  $R_1^{m-1} = \{(1, x_2, \dots, x_m)^\top \in \mathbb{R}^m\}.$ 

Theorem 4.2.16. Furthermore,

$$P\left(\max_{x \in \mathbb{R}_1^{m-1}} \frac{\left(x^T \hat{\beta} - \left(\beta_1 + \sum_{i=2}^m \beta_i x_i\right)\right)^2}{\hat{\sigma}^2 x^\top (X^\top X)^{-1} x} \le m \cdot F_{m,n-m,1-\alpha}\right) = 1 - \alpha$$

holds. Without proof.

# 4.3 Multivariate linear regression with rank(X) < m

Let  $Y = X\beta + \varepsilon$ ,  $Y \in \mathbb{R}^n$ , where X is a  $(n \times m)$  matrix with rank (X) = r < m,  $\beta = (\beta_1, \dots, \beta_m)^\top$ ,  $\varepsilon \in \mathbb{R}^n$ ,  $\mathbb{E} \varepsilon = 0$ ,  $\mathbb{E} (\varepsilon_i \varepsilon_j) = \delta_{ij} \sigma^2$ ,  $i, j = 1, \dots, n$ ,  $\sigma^2 > 0$ .

Even though the rank of the matrix is not full anymore, the OLS estimator  $\hat{\beta}$  is still a solution to the normal equation

$$(X^{\top}X)\beta = X^{\top}Y.$$

However,  $X^{\top}X$  is not invertible, since

$$\operatorname{rank} \ (\boldsymbol{X}^{\top} \boldsymbol{X}) \leq \min \left\{ \operatorname{rank} \ (\boldsymbol{X}), \operatorname{rank} \ (\boldsymbol{X}^{\top}) \right\} = r < m.$$

Consequently, in order to obtain  $\hat{\beta}$  from the normal equation, both sides of the equations are multiplied with the *generalized inverse* of  $X^{\top}X$ .

# 4.3.1 Generalized inverse

**Definition 4.3.1.** Let A be a  $(n \times m)$  matrix. A  $(m \times n)$  matrix  $A^-$  is called *generalized inverse* of A, if

$$AA^{-}A = A.$$

The matrix  $A^-$  is not unique, which is shown in the following lemmas.

**Lemma 4.3.2.** Let A be a  $(n \times m)$  matrix,  $m \le n$  with rank  $(A) = r \le m$ . There exist invertible matrices  $P(n \times n)$  and  $Q(m \times m)$ , such that

$$PAQ = \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{pmatrix}$$
, where  $I_r = \operatorname{diag}(\underbrace{1, \dots, 1}_{r \text{ times}})$ . (4.6)

Corollary 4.3.3. For an arbitrary  $(n \times m)$  matrix A with  $n \geq m$ , rank (A)  $r \leq m$  it holds that

$$A^{-} = Q \begin{pmatrix} \mathcal{I}_r & A_2 \\ A_1 & A_3 \end{pmatrix} P, \tag{4.7}$$

where P and Q are matrices as in (4.6),  $\mathcal{I}_r = \text{diag}(1,\ldots,1)$ , and  $A_1$ ,  $A_2$  resp.  $A_3$  are arbitrary  $((m-r)\times r)$ ,  $(r\times (n-r))$  resp.  $((m-r)\times (n-r))$  matrices.

In particular

$$A_1 = 0,$$

$$A_2 = 0,$$

$$A_3 = \operatorname{diag} \left(\underbrace{1, \dots, 1}_{s-r \text{ times}}, 0, \dots, 0\right),$$

$$s \in \{r, \dots, m\}$$

can be chosen, which means rank  $(A^{-}) = s \in \{r, ..., m\}$  for

$$A^{-} = Q \begin{pmatrix} \mathcal{I}_{s} & 0 \\ 0 & 0 \end{pmatrix} P.$$

**Proof** In the following it is shown that for  $A^-$  as in (4.7), it holds that  $AA^-A = A$ . Lemma 4.3.2 implies that

$$A = P^{-1} \cdot \operatorname{diag} (1, \dots, 1, 0, \dots, 0) \cdot Q^{-1} \quad \text{and thus}$$

$$AA^{-}A = P^{-1} \begin{pmatrix} \mathcal{I}_{r} & 0 \\ 0 & 0 \end{pmatrix} Q^{-1}Q \cdot \begin{pmatrix} \mathcal{I}_{r} & A_{2} \\ A_{1} & A_{3} \end{pmatrix} PP^{-1} \begin{pmatrix} \mathcal{I}_{r} & 0 \\ 0 & 0 \end{pmatrix} Q^{-1}$$

$$= P^{-1} \begin{pmatrix} \mathcal{I}_{r} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathcal{I}_{r} & A_{2} \\ A_{1} & A_{3} \end{pmatrix} \begin{pmatrix} \mathcal{I}_{r} & 0 \\ 0 & 0 \end{pmatrix} Q^{-1}$$

$$= P^{-1} \begin{pmatrix} \mathcal{I}_{r} & 0 \\ 0 & 0 \end{pmatrix} Q^{-1} = A.$$

**Lemma 4.3.4.** Let A be an arbitrary  $(n \times m)$  matrix with rank  $(A) = r \le m, m \le n$ .

1. If  $(A^{\top}A)^{-}$  is a generalized inverse of an  $(m \times m)$  matrix  $A^{\top}A$ , then  $\left((A^{\top}A)^{-}\right)^{\top}$  is also a generalized inverse of  $A^{\top}A$ .

2. It holds that

$$(A^{\top}A)(A^{\top}A)^{-}A^{\top} = A^{\top}$$
 resp.  
 $A(A^{\top}A)^{-}(A^{\top}A) = A$ .

# Proof

1.  $A^{\top}A$  is symmetric, i.e.

$$\underbrace{\left(A^{\top}A(A^{\top}A)^{-}A^{\top}A\right)^{\top}}_{=A^{\top}A\left((A^{\top}A)^{-}\right)^{\top}A^{\top}A} = \left(A^{\top}A\right)^{\top} = A^{\top}A.$$

Thus  $((A^{\top}A)^{-})^{\top}$  is a generalized inverse of  $A^{\top}A$ .

2. Let  $B = (A^{\top}A)(A^{\top}A)^{-}A^{\top} - A^{\top}$ . In the following it is shown that B = 0 by proving that  $BB^{\top} = 0$ .

$$\begin{split} BB^\top &= \left( (A^\top A)(A^\top A)^- A^\top - A^\top \right) \left( A \left( (A^\top A)^- \right)^\top A^\top A - A \right) \\ &= A^\top A (A^\top A)^- A^\top A \left( (A^\top A)^- \right)^\top A^\top A - \underbrace{A^\top A (A^\top A)^- A^\top A}_{=A^\top A} \\ &- \underbrace{A^\top A \left( (A^\top A)^- \right)^\top \cdot A^\top A}_{=A^\top A} + A^\top A \\ &= A^\top A - 2A^\top A + A^\top A = 0. \end{split}$$

The assertion  $A(A^{\top}A)^{-}A^{\top}A = A$  can be shown, by transposing the matrices on both sides of the equation  $A^{\top}A(A^{\top}A)^{-}A^{\top} = A^{\top}$ .

# 4.3.2 OLS estimator for $\beta$

**Theorem 4.3.5.** Let X be a  $(n \times m)$  design matrix with rank (X) = r < m in the linear regression model  $Y = X\beta + \varepsilon$ . The generalized solution of the normal equation

$$(X^{\top}X)\beta = X^{\top}Y$$

is given by

$$\beta = \left(X^{\top}X\right)^{-}X^{\top}Y + \left(\mathcal{I}_{m} - \left(X^{\top}X\right)^{-}X^{\top}X\right)z, \quad z \in \mathbb{R}^{m}. \tag{4.8}$$

Proof

1. In the following it is shown that  $\beta$  as in (4.8) is a solution of the normal equation.

$$X^{\top}X\beta = \underbrace{(X^{\top}X)(X^{\top}X)^{-}X^{\top}}_{=X^{\top}(\text{Lemma 4.3.4, 2.}))}Y + \left(X^{\top}X - \underbrace{X^{\top}X(X^{\top}X)^{-}X^{\top}X}_{=X^{\top}X}\right)z$$
$$= X^{\top}Y$$

2. Let us show that an arbitrary solution  $\beta'$  of the normal equation can be written as (4.8). Let  $\beta$  be the solution (4.8). Calculating the difference of the equations yields

$$(X^{\top}X)\beta' = X^{\top}Y$$

$$- (X^{\top}X)\beta = X^{\top}Y$$

$$(X^{\top}X)(\beta' - \beta) = 0$$

$$\beta' = (\beta' - \beta) + \beta$$

$$= \beta' - \beta + (X^{\top}X)^{-}X^{\top}Y + (\mathcal{I}_m - (X^{\top}X)^{-}X^{\top}X)z$$

$$= (X^{\top}X)^{-}X^{\top}Y + (\mathcal{I}_m - (X^{\top}X)^{-}X^{\top}X)z + (\beta' - \beta)$$

$$- \underbrace{(X^{\top}X)^{-}X^{\top}X(\beta' - \beta)}_{=0}$$

$$= (X^{\top}X)^{-}X^{\top}Y + (\mathcal{I}_m - (X^{\top}X)^{-}X^{\top}X)(\underbrace{z + \beta' - \beta}_{=z_0})$$

$$\Longrightarrow \beta' \text{ can be rewritten as } (4.8).$$

**Remark 4.3.6.** Theorem 4.3.5 yields the set of all extreme points of the OLS minimization problem

$$e(\beta) = \frac{1}{n} |Y - X\beta|^2 \longrightarrow \min_{\beta}.$$

Thus, the set of all OLS estimators of  $\beta$  in (4.8) should satisfy additional conditions.

# Theorem 4.3.7.

1. All OLS estimators of  $\beta$  can be written as

$$\overline{\beta} = \left( X^{\top} X \right)^{-} X^{\top} Y,$$

where  $(X^{\top}X)^{-}$  is an arbitrary generalized inverse of  $X^{\top}X$ .

2.  $\overline{\beta}$  is not unbiased, since

$$\mathbb{E}\,\overline{\beta} = \left(X^{\top}X\right)^{-}X^{\top}X\beta.$$

3. It holds that

$$\operatorname{Cov}\, \overline{\beta} = \sigma^2 \left( X^\top X \right)^- \left( X^\top X \right) \left( \ (X^\top X)^- \right)^\top.$$

#### **Proof**

1. We show that  $e(\beta) \ge e(\overline{\beta}) \quad \forall \beta \in \mathbb{R}^m$ .

$$n \cdot e(\beta) = |Y - X\beta|^2 = (Y - X\overline{\beta} + X(\overline{\beta} - \beta))^\top (Y - X\overline{\beta} + X(\overline{\beta} - \beta))$$

$$= (Y - X\overline{\beta})^\top (Y - X\overline{\beta}) + \left(X(\overline{\beta} - \beta)\right)^\top \left(X(\overline{\beta} - \beta)\right)$$

$$+ 2(\overline{\beta} - \beta)^\top X^\top (Y - X\overline{\beta})$$

$$= n \cdot e(\overline{\beta}) + \underbrace{2 \cdot (\overline{\beta} - \beta)^\top (X^\top Y - (X^\top X\overline{\beta}))}_{=0} + \left|X(\overline{\beta} - \beta)\right|^2$$

$$\geq n \cdot e(\overline{\beta}) + 0 = n \cdot e(\overline{\beta}), \quad \text{since}$$

 $\overline{\beta}$  can be rewritten as in (4.8) with z=0 and is thus a solution to the normal equation.

2. It holds that

$$\mathbb{E}\,\overline{\beta} = \mathbb{E}\,\left((X^\top X)^- X^\top Y\right) = \left(X^\top X\right)^- X^\top \mathbb{E}\,Y$$
$$= (X^\top X)^- X^\top X \beta,$$
$$Y = X\beta + \varepsilon, \quad \mathbb{E}\,\varepsilon = 0.$$

Note that this implies  $\mathbb{E} Y = X\beta$ . Why is  $\overline{\beta}$  not unbiased, i.e.  $(X^{\top}X)^{-}X^{\top}X\beta \neq \beta, \ \beta \in \mathbb{R}^{m}$ ?

Since rank (X) = r < m, rank  $(X^{\top}X) < m$  and rank  $((X^{\top}X)^{-}X^{\top}X) < m$ . Thus, there exists a  $\beta \neq 0$ , for which it holds that

$$(X^{\top}X)^{-}X^{\top}X\beta = 0 \neq \beta,$$

so  $\overline{\beta}$  is not unbiased. Furthermore it holds that all solutions of (4.8) are not unbiased estimators. Applying the expectation on (4.8) yields the following in the case of unbiasedness

$$\forall \beta \in \mathbb{R}^m : \quad \beta = (X^\top X)^- X^\top X \beta + \left( \mathcal{I}_m - (X^\top X)^- (X^\top X) \right) z, \quad z \in \mathbb{R}^m.$$

$$\Longrightarrow \left( \mathcal{I}_m - (X^\top X)^- (X^\top X) \right) (z - \beta) = 0 \quad \forall z, \beta \in \mathbb{R}^m.$$

$$\Longrightarrow (X^\top X)^- (X^\top X) (\beta - z) = \beta - z, \quad \forall z, \beta \in \mathbb{R}^m.$$

Since this equation can't hold for all  $\beta \in \mathbb{R}^m$ , the assumption leads to a contradiction.

3. It holds that

$$\begin{aligned} &\operatorname{Cov} \ \left(\overline{\beta}_{i}, \overline{\beta}_{j}\right) = \operatorname{Cov} \left(\left(\underbrace{(X^{\top}X)^{-}X^{\top}}_{:=A=(a_{kl})}Y\right)_{i}, \left((X^{\top}X)^{-}X^{\top}Y\right)_{j}\right) \\ &= \operatorname{Cov} \ \left(\sum_{k=1}^{n} a_{ik}Y_{k}, \sum_{l=1}^{n} a_{jl}Y_{l}\right) \\ &= \sum_{k,l=1}^{n} a_{ik}a_{jl} \underbrace{\operatorname{Cov} \left(Y_{k}, Y_{l}\right)}_{=\sigma^{2} \cdot \delta_{kl}} = \sigma^{2} \sum_{k=1}^{n} a_{ik}a_{jk} = \left(\sigma^{2}AA^{\top}\right)_{i,j} \\ &= \left(\sigma^{2}(X^{\top}X)^{-}X^{\top}X \left((X^{\top}X)^{-}\right)^{\top}\right)_{i,j}. \end{aligned}$$

#### 4.3.3 Functions that can be estimated without bias

**Definition 4.3.8.** A linear combination  $a^{\top}\beta$  of  $\beta_1, \ldots, \beta_m, a \in \mathbb{R}^m$  is called *estimable without bias*, if

$$\exists c \in \mathbb{R}^n : \quad \mathbb{E}\left(c^\top Y\right) = a^\top \beta,$$

i.e. if a linear unbiased estimator  $c^{\top}Y$  for  $a^{\top}\beta$  exists.

**Theorem 4.3.9.** The function  $a^{\top}\beta$ ,  $a \in \mathbb{R}^m$  is estimable without bias if and only if one of the following conditions is satisfied:

- 1.  $\exists c \in \mathbb{R}^n : a^{\top} = c^{\top} X$ .
- $2. \ a$  fulfills the equation

$$a^{\top} \left( X^{\top} X \right)^{-} X^{\top} X = a^{\top}. \tag{4.9}$$

Proof

1. " $\Longrightarrow$  ": If  $a^{\top}\beta$  is estimable, then there exists a vector  $d \in \mathbb{R}^n$  with  $\mathbb{E}(d^{\top}Y) = a^{\top}\beta \quad \forall \beta \in \mathbb{R}^m$ . So

$$a^{\top}\beta = d^{\top}\mathbb{E}Y = d^{\top}X\beta \Rightarrow (a^{\top} - d^{\top}X)\beta = 0, \quad \forall \beta \in \mathbb{R}^m$$
  
 $\Longrightarrow a^{\top} = d^{\top}X,$ 

Finally, setting c = d proves this implication.

2. " $\Longrightarrow$  ": If  $a^{\top}\beta$  is estimable without bias, then

$$a^{\top} (X^{\top} X)^{-} X^{\top} X \stackrel{\text{1.}}{=} c^{\top} \underbrace{X \cdot (X^{\top} X)^{-} X^{\top} X}_{=X \text{ (Lemma 4.3.4)}} = c^{\top} X \stackrel{\text{(1.)}}{=} a^{\top}.$$

Thus, relation (4.9) is satisfied.

**Remark 4.3.10.** In case of a regression with rank (X) = m the equation (4.9) is always satisfied since  $(X^{\top}X)^{-} = (X^{\top}X)^{-1}$  and thus  $a^{\top}\beta$  is estimable for all  $a \in \mathbb{R}^m$ .

**Theorem 4.3.11** (Examples of estimable functions). If rank (X) = r < m, then the following linear combinations of  $\beta$  are estimable:

- 1. The coordinates  $\sum_{j=1}^{m} x_{ij}\beta_j$ ,  $i=1,\ldots,n$  of the vector of expectations  $\mathbb{E} Y = X\beta$ .
- 2. Arbitrary linear combinations of estimable functions.

#### Proof

1. Set  $\tilde{x}_i = (x_{i1}, \dots, x_{im}), i = 1, \dots, n$ . Then

$$\sum_{j=1}^{m} x_{ij\beta_j} = \tilde{x}_i^{\top} \beta \quad \forall i = 1, \dots, n,$$
$$X\beta = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)^{\top} \beta.$$

 $\tilde{x}_i\beta$  is estimable, if  $\tilde{x}_i$  satisfies (4.9), which can be expressed in matrices for all  $i=1,\ldots,n$  as follows:

$$X \left( X^{\top} X \right)^{-} X^{\top} X = X.$$

By Lemma 4.3.4 this is valid.

2. For  $a_1, \ldots, a_k \in \mathbb{R}^m$  let  $a_1^\top \beta, \ldots, a_k^\top \beta$  be estimable functions. For all  $\lambda = (\lambda_1, \ldots, \lambda_k)^\top \in \mathbb{R}^k$  show that  $\sum_{i=1}^k \lambda_i \cdot a_i^\top \beta = \lambda^\top A \beta$  is estimable, where  $A = (a_1, \ldots, a_k)^\top$ . It needs to be shown that  $b = (\lambda^\top A)^\top$  satisfies (4.9), i.e.

$$\lambda^{\top} A \left( X^{\top} X \right)^{-} X^{\top} X = \lambda^{\top} A.$$

This equation is satisfied, since  $a_i^{\top}(X^{\top}X)^{-}X^{\top}X = a_i^{\top}, i = 1, ..., k$ . By Theorem 4.3.9, 2.)  $\lambda^{\top}A\beta$  is estimable.

**Theorem 4.3.12** ( $Gau\beta$ -Markov). Let  $a^{\top}\beta$  be an estimable function,  $a \in \mathbb{R}^m$  in the linear regression model  $Y = X\beta + \varepsilon$  with rank  $(X) \leq m$ .

1. The best linear unbiased estimator of  $a^{\top}\beta$  is given by  $a^{\top}\overline{\beta}$ , where

$$\overline{\beta} = \left( X^{\top} X \right)^{-} X^{\top} Y$$

is an OLS estimator for  $\beta$ .

2. Var 
$$(a^{\top}\overline{\beta}) = \sigma^2 a^{\top} (X^{\top}X)^{-}a$$
.

**Proof** The linearity of  $a^{\top}\overline{\beta} = a^{\top}(X^{\top}X)^{-}X^{\top}Y$  as a function of Y is clear. For the unbiasedness it holds that

$$\mathbb{E}\left(a^{\top}\overline{\beta}\right) = a^{\top}\mathbb{E}\overline{\beta} = a^{\top}(X^{\top}X)^{-}X^{\top}X\beta$$
$$= c^{\top}\underbrace{X(X^{\top}X)^{-}X^{\top}X}_{=X \text{ (Lemma 4.3.4)}}\beta = \underbrace{c^{\top}X}_{=a^{\top}}\beta = a^{\top}\beta \quad \forall \beta \in \mathbb{R}^{m}.$$

First, calculate  $\operatorname{Var}(a^{\top}\overline{\beta})$  (i.e. prove the second assertion) and show that it is minimal:

$$\begin{aligned} \operatorname{Var}\left(a^{\top}\overline{\beta}\right) &= \operatorname{Var}\left(\sum_{i=1}^{m} a_{i}\overline{\beta}_{i}\right) = \sum_{i,j=1}^{m} a_{i}a_{j} \cdot \operatorname{Cov}\left(\overline{\beta}_{i}, \overline{\beta}_{j}\right) \\ &= a^{\top}\operatorname{Cov}\left(\overline{\beta}\right) a \overset{\left(\operatorname{Satz} 4.3.7\right)}{=} a^{\top}\sigma^{2}\left((X^{\top}X)^{-}X^{\top}X(X^{\top}X)^{-}\right)^{\top} a \\ &= \sigma^{2} \cdot a^{\top}\underbrace{\left((X^{\top}X)^{-}\right)^{\top}}_{=(X^{\top}X)^{-}} X^{\top}X\underbrace{\left((X^{\top}X)^{-}\right)^{\top}}_{(X^{\top}X)^{-}} a \\ &\stackrel{\operatorname{Lemma}}{=} \overset{4.3.4, \ 1.)}{=} \sigma^{2}a^{\top}(X^{\top}X)^{-}X^{\top}X(X^{\top}X)^{-}a \\ &\stackrel{\operatorname{Theorem}}{=} \overset{4.3.9, \ 1.)}{=} \sigma^{2} \cdot c^{\top}\underbrace{X \cdot (X^{\top}X)X^{\top}X(X^{\top}X)^{-}X^{\top}}_{=X} c \\ &= \sigma^{2}\underbrace{c^{\top}X}_{=a^{\top}}(X^{\top}X)^{-}\underbrace{X^{\top}c}_{=a} = \sigma^{2}a^{\top}(X^{\top}X)^{-}a. \end{aligned}$$

Now it is shown that for an arbitrary linear unbiased estimator  $b^{\top}Y$  of  $a^{\top}\beta$  it holds that  $\operatorname{Var}(b^{\top}Y) \geq \operatorname{Var}(a^{\top}\overline{\beta})$ . Since  $b^{\top}Y$  is unbiased, it holds that  $\mathbb{E}(b^{\top}Y) = a^{\top}\beta$ . Using Theorem 4.3.9 it holds that  $a^{\top} = b^{\top}X$ . Now, consider

$$0 \le \operatorname{Var} \left( b^{\top} Y - a^{\top} \overline{\beta} \right)$$

$$= \operatorname{Var} \left( b^{\top} Y \right) - 2 \operatorname{Cov} \left( b^{\top} Y, a^{\top} \overline{\beta} \right) + \operatorname{Var} \left( a^{\top} \overline{\beta} \right)$$

$$= \operatorname{Var} \left( b^{\top} Y \right) - 2 \sigma^{2} a^{\top} (X^{\top} X)^{-} a + \sigma^{2} a^{\top} (X^{\top} X)^{-} a$$

$$= \operatorname{Var} \left( b^{\top} Y \right) - \operatorname{Var} \left( a^{\top} \overline{\beta} \right)$$

with

$$\begin{aligned} \operatorname{Cov} \ \left( b^{\top} Y, \, a^{\top} \overline{\beta} \right) &= \operatorname{Cov} \ \left( b^{\top} Y, \, a^{\top} (X^{\top} X)^{-} X^{\top} Y \right) = \sigma^{2} a^{\top} (X^{\top} X)^{-} \underbrace{X^{\top} b}_{=a} \\ &= \sigma^{2} a^{\top} (X^{\top} X)^{-} a. \end{aligned}$$

Thus,  $\operatorname{Var}\left(b^{\top}Y\right) \geq \operatorname{Var}\left(a^{\top}\overline{\beta}\right)$  and  $a^{\top}\overline{\beta}$  is a best linear unbiased estimator for  $a^{\top}\beta$ .

# Remark 4.3.13.

- 1. If rank (X) = m, then  $a^{\top}\hat{\beta}$  is the best linear unbiased estimator for  $a^{\top}\beta$ ,  $a \in \mathbb{R}^m$ .
- 2. The estimator  $a^{\top}\overline{\beta} = a^{\top}(X^{\top}X)^{-}X^{\top}Y$  does not depend on the choice of the generalized inverse as is shown in the following theorem.

**Theorem 4.3.14.** The best linear unbiased estimator  $a^{\top}\overline{\beta}$  for  $a^{\top}\beta$  is uniquely determined.

#### Proof

$$\boldsymbol{a}^{\top}\overline{\boldsymbol{\beta}} = \boldsymbol{a}^{\top}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-}\boldsymbol{X}^{\top}\boldsymbol{Y} \overset{\text{Theorem 4.3.9, 1.)}}{=} \boldsymbol{c}^{\top}\boldsymbol{X}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-}\boldsymbol{X}^{\top}\boldsymbol{Y}.$$

In order to show that  $X(X^{\top}X)^{-}X^{\top}$  does not depend on  $(X^{\top}X)^{-}$ , we prove that for arbitrary generalized inverses  $A_1$  and  $A_2$  of  $(X^{\top}X)$  it holds that  $XA_1X^{\top} = XA_2X^{\top}$ . By Lemma 4.3.4, 2.) it holds that

$$XA_1X^{\top}X = X = XA_2X^{\top}X.$$

Multiplying all parts of the equation with  $A_1X^{\top}$  on the right yields

$$XA_1 \underbrace{X^\top X A_1 X^\top}_{=X^\top} = XA_1 X^\top = XA_2 \underbrace{X^\top X A_1 X^\top}_{=X^\top}$$

Thus,  $XA_1X^{\top} = XA_2X^{\top}$ .

# 4.3.4 Normally distributed error terms

Let  $Y = X\beta + \varepsilon$  be a linear regression model with rank (X) = r < m and  $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathcal{I})$ . As in Section 4.2.3 the maximum likelihood estimator  $\tilde{\beta}$  and  $\tilde{\sigma}^2$  can be derived for  $\beta$  and  $\sigma^2$ . Exactly as in Theorem 4.2.7 it can be shown that

$$\tilde{\beta} = \overline{\beta} = (X^{\top}X)^{-}X^{\top}Y$$
 and  $\tilde{\sigma}^2 = \frac{1}{n} |Y - X\overline{\beta}|^2$ .

Now the distributional properties of  $\overline{\beta}$  and  $\tilde{\sigma}^2$  are discussed. First the unbiasedness of  $\tilde{\sigma}^2$  is discussed.  $\tilde{\sigma}^2$  is not unbiased but, the corrected estimator

$$\overline{\sigma}^2 = \frac{1}{n-r} |Y - X\beta|^2 = \frac{n}{n-r} \tilde{\sigma}^2$$

is unbiased.

**Theorem 4.3.15.** The estimator  $\overline{\sigma}^2$  is unbiased for  $\sigma^2$ .

The proof of Theorem 4.3.15 is similar to the proof of Theorem 4.2.6 in which  $\hat{\beta} = (X^\top X)^{-1} X^\top Y$  and  $\hat{\sigma}^2 = \frac{1}{n-m} |Y - X\beta|^2$  are considered for the case rank (X) = m. Thus, Theorem 4.2.6 is a special case of Theorem 4.3.15. Define  $D \coloneqq \mathcal{I} - X(X^\top X)^- X^\top$ .

**Lemma 4.3.16.** For D it holds that

- 1.  $D^{\top} = D$  (symmetry),
- 2.  $D^2 = D$  (idempotence),
- 3. DX = 0,
- 4.  $\operatorname{trace}(D) = n r$ .

# Proof

1. It holds that

$$D^{\top} = \left(\mathcal{I} - X(X^{\top}X)^{-}X^{\top}\right)^{\top} = \mathcal{I} - X\left((X^{\top}X)^{-}\right)^{\top}X^{\top}$$
$$= \mathcal{I} - X(X^{\top}X)^{-}X^{\top} = D,$$

since  $((X^{\top}X)^{-})^{\top}$  is also a generalized inverse of  $X^{\top}X$  (cf. Lemma 4.3.4, 1.)).

2. It holds that

$$D^{2} = \left(\mathcal{I} - X(X^{\top}X)^{-}X^{\top}\right)^{2}$$

$$= \mathcal{I} - 2X(X^{\top}X)^{-}X^{\top} + \underbrace{X(X^{\top}X)^{-}X^{\top}X}_{=X \text{ (Lemma 4.3.4, 2.))}}(X^{\top}X)^{-}X^{\top}$$

$$= \mathcal{I} - X(X^{\top}X)^{-}X^{\top} = D.$$

3. 
$$DX = X - \underbrace{X(X^{\top}X)^{-}X^{\top}X}_{=X \text{ (Lemma 4.3.4, 2.))}} = X - X = 0.$$

4. It holds that

$$\operatorname{trace}(D) = \operatorname{trace}(I) - \operatorname{trace}\left(X(X^{\top}X)^{-}X^{\top}\right)$$
$$= n - \operatorname{trace}\left(X(X^{\top}X)^{-}X^{\top}\right).$$

The symmetry and idempotence of the matrix A imply  $\operatorname{trace}(A) = \operatorname{rank}(A)$  as is known from linear algebra. Since  $X(X^{\top}X)^{-}X^{\top}$  is symmetric and idempotent, it is sufficient to show  $\operatorname{rank}(X(X^{\top}X)^{-}X^{\top}) = r$ . By Lemma 4.3.4 2.) it holds that

$$\operatorname{rank} (X) = r = \operatorname{rank} (X(X^{\top}X)^{-}X^{\top}X)$$

$$\leq \min \left\{ \operatorname{rank} (X(X^{\top}X)^{-}X^{\top}), \underbrace{\operatorname{rank} (X)}_{=r} \right\}$$

$$\leq \operatorname{rank} \left( X(X^{\top}X)^{-}X^{\top} \right) \leq \operatorname{rank} (X) = r$$

$$\Longrightarrow \operatorname{rank} \left( X(X^{\top}X)^{-}X^{\top} \right) = r$$

$$\Longrightarrow \operatorname{trace} \left( X(X^{\top}X)^{-}X^{\top} \right) = r.$$

**Proof of Theorem 4.3.15** By using Lemma 4.3.16 it can be shown that

$$\begin{split} \overline{\sigma}^2 &= \frac{1}{n-r} \left| Y - X \overline{\beta} \right|^2 = \frac{1}{n-r} \left| Y - X (X^\top X)^- X^\top Y \right|^2 = \frac{1}{n-r} |DY|^2 \\ &= \frac{1}{n-r} \Big| \underbrace{DX}_{=0} \beta + D\varepsilon \Big|^2 = \frac{1}{n-r} |D\varepsilon|^2 = \frac{1}{n-r} \varepsilon^\top \underbrace{D^\top D}_{=D^2 = D} \varepsilon = \frac{1}{n-r} \varepsilon^\top D\varepsilon. \end{split}$$

Thus

$$\mathbb{E}\,\overline{\sigma}^2 = \frac{1}{n-r} \mathbb{E}\left(\varepsilon^\top D \varepsilon\right) = \frac{1}{n-r} \mathbb{E}\operatorname{trace}\left(\varepsilon^\top D \varepsilon\right) = \frac{1}{n-r} \operatorname{trace}\left(D \cdot \mathbb{E}\left(\underbrace{\varepsilon \varepsilon^\top}_{\sigma^2 \mathcal{I}}\right)\right)$$
$$= \frac{\sigma^2}{n-r} \cdot \operatorname{trace}(D) = \sigma^2$$

by Lemma 4.3.16, 4.), because  $\mathbb{E} \, \varepsilon \varepsilon^{\top} = \sigma^2 \mathcal{I}$  and  $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathcal{I})$ 

**Theorem 4.3.17.** The following distributional properties hold

1. 
$$\overline{\beta} \sim N\left( (X^{\top}X)^{-}X^{\top}X\beta, \, \sigma^2(X^{\top}X)^{-}(X^{\top}X) \left( (X^{\top}X)^{-} \right)^{\top} \right),$$

2. 
$$\frac{(n-r)\overline{\sigma}^2}{\sigma^2} \sim \chi_{n-r}^2$$

3.  $\overline{\beta}$  and  $\overline{\sigma}^2$  are independent.

#### Proof

1. It holds that

$$\overline{\beta} = (X^{\top}X)^{-}X^{\top}Y = (X^{\top}X)^{-}X^{\top}(X\beta + \varepsilon)$$
$$= \underbrace{(X^{\top}X)^{-}X^{\top}X\beta}_{=\mu} + \underbrace{(X^{\top}X)^{-}X^{\top}}_{=A}\varepsilon.$$

Consequently,

$$\overline{\beta} \sim N\left(\mu, \, \sigma^2 A A^\top\right)$$

$$= N\left((X^\top X)^- X^\top X \beta, \, \sigma^2 (X^\top X)^- X^\top X ((X^\top X)^-)^\top\right)$$

with 
$$AA^{\top} = (X^{\top}X)^{-}X^{\top}X((X^{\top}X)^{-})^{\top}$$
.

2. It holds that  $\overline{\sigma}^2 = \frac{1}{n-r} \varepsilon^{\top} D \varepsilon$  by the proof of Theorem 4.3.15. Thus,

$$\frac{(n-r)\overline{\sigma}^2}{\sigma^2} = \underbrace{\left(\frac{\varepsilon}{\sigma}\right)^{\top}}_{\sim \mathcal{N}(0,\mathcal{I})} D\left(\frac{\varepsilon}{\sigma}\right) \overset{\text{(Satz 4.1.25)}}{\sim} \chi_{n-r}^2.$$

3. Consider  $A\varepsilon$  and  $\varepsilon^{\top}D\varepsilon$ . It is sufficient to show that they are independent in order to show the independence of  $\overline{\beta}$  and  $\overline{\sigma}^2$ , since  $\overline{\beta} = \mu + A\varepsilon$ ,  $\overline{\sigma}^2 = \frac{1}{n-r}\varepsilon^{\top}D\varepsilon$ . It holds that  $A \cdot \sigma^2 \mathcal{I} \cdot D = 0$ . By Theorem 4.1.26  $A\varepsilon$  and  $\varepsilon^{\top}D\varepsilon$  are independent.

#### 

#### 4.3.5 Hypothesis testing

Consider the hypothesis test  $H_0: H\beta = d$  vs.  $H_1: H\beta \neq d$ , where H is an  $(s \times m)$  matrix  $(s \leq m)$  with rank (H) = s and  $d \in \mathbb{R}^s$ . In Theorem 4.2.14 for the case rank (X) = r = m the following test statistic was considered

$$T = \frac{(H\hat{\beta} - d)^{\top} (H(X^{\top}X)^{-1}H^{\top})^{-1} (H\hat{\beta} - d)}{s\hat{\sigma}^{2}} \stackrel{(H_{0})}{\sim} F_{s,n-m}.$$

In general, we may consider

$$T = \frac{(H\overline{\beta} - d)^{\top} (H(X^{\top}X)^{-}H^{\top})^{-1} (H\overline{\beta} - d)}{s\overline{\sigma}^{2}}.$$
 (4.10)

We will show that  $T \stackrel{(H_0)}{\sim} F_{s,n-r}$ . Then, a test with confidence level  $\alpha \in (0,1)$  can be constructed by rejecting  $H_0$ , if  $T > F_{s,n-r,1-\alpha}$ .

**Definition 4.3.18.** The hypothesis  $H_0: H\beta = d$  is called *testable*, if all coordinates of the vector  $H\beta$  are estimable functions.

Theorem 4.3.9 provides conditions for H, under which  $H_0: H\beta = d$  is testable. They are formulated in the following Lemma

**Lemma 4.3.19.** The hypothesis  $H_0: H\beta = d$  is testable if and only if

- 1. There exists an  $(s \times n)$  matrix C such that H = CX, or
- 2.  $H(X^{\top}X)^{-}X^{\top}X = H$ .

First, show that the test statistic T in (4.10) is well defined, i.e. the  $(s \times s)$  matrix  $H(X^{\top}X)^{-}H^{\top}$  is positive definite and thus invertible. Corollary 4.3.3 implies

$$X^{\top}X = P^{-1} \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{pmatrix} P^{-1} \tag{4.11}$$

for an  $(m \times m)$  matrix P, which is symmetric and invertible. Thus,

$$(X^{\top}X)^{-} = P \cdot \begin{pmatrix} \mathcal{I}_{r} & 0 \\ 0 & \mathcal{I}_{m-r} \end{pmatrix} P = P \cdot P,$$

holds, i.e. there exists a unique generalized inverse  $X^{\top}X$  with this representation. This implies that the  $(s \times s)$  matrix  $HPPH^{\top} = (PH^{\top})^{\top} \cdot PH^{\top}$  is positive definite because rank  $(PH^{\top}) = s$ . Let now  $(X^{\top}X)^{-}$  be an arbitrary generalized inverse of  $X^{\top}X$ . Then, Lemma 4.3.19 implies

$$H(X^{\top}X)^{-}H^{\top} = CX(X^{\top}X)^{-}X^{\top}C^{\top} = CXPPX^{\top}C^{\top} = HPPH^{\top},$$

because  $X(X^{\top}X)^{-}X^{\top}$  is invariant with respect to the choice  $(X^{\top}X)^{-}$  by the proof of Theorem 4.3.14. Thus,  $H\left(X^{\top}X\right)^{-}H^{\top}$  is positive definite for an arbitrary generalized inverse  $\left(X^{\top}X\right)^{-}$  and the test statistic T is well defined.

**Theorem 4.3.20.** If  $H_0: H\beta = d$  is testable, then  $T \stackrel{(H_0)}{\sim} F_{s,n-r}$ .

**Proof** This proof is similar to the proof of Theorem 4.2.14. First, we compute

$$H\overline{\beta} - d = H(X^{\top}X)^{-}X^{\top}(X\beta + \varepsilon) - d$$
$$= \underbrace{H(X^{\top}X)^{-}X^{\top}X\beta - d}_{=\mu} + \underbrace{H(X^{\top}X)^{-}X^{\top}}_{=B}\varepsilon.$$

Show that  $\mu \stackrel{(H_0)}{=} 0$ .

$$\mu \overset{\text{(Lemma 4.3.19)}}{=} C \cdot \underbrace{X(X^\top X)^- X^\top X}_{=X \text{ (Lemma 4.3.4, 2.))}} \cdot \beta - d = CX\beta - d = H\beta - d \overset{(H_0)}{=} 0.$$

Using Theorem 4.3.17 yields  $(H\overline{\beta} - d)^{\top} \left( H(X^{\top}X)^{-}H^{\top} \right)^{-1} \left( H\overline{\beta} - d \right)$  and  $s \cdot \overline{\sigma}^2$  are independent and  $\frac{(n-r)\overline{\sigma}^2}{\sigma^2} \sim \chi_{n-r}^2$ . It remains to show

$$\left(\underbrace{H\overline{\beta}-d}_{=\varepsilon^{\top}B^{\top}}\right)^{\top}\left(H(X^{\top}X)^{-}H^{\top}\right)^{-1}\left(\underbrace{H\overline{\beta}-d}_{=B\varepsilon}\right)\overset{(H_{0})}{\sim}\chi_{s}^{2}.$$

It holds that

$$\varepsilon^{\top} B^{\top} \left( H(X^{\top} X)^{-} H^{\top} \right)^{-1} B \varepsilon$$

$$= \varepsilon^{\top} \underbrace{X \left( (X^{\top} X)^{-} \right)^{\top} H^{\top} \left( H(X^{\top} X)^{-} H^{\top} \right)^{-1} H(X^{\top} X)^{-} X^{\top}}_{A} \varepsilon$$

It can be shown that A is symmetric, idempotent and rank (A) = s. The idempotence can be shown as follows

$$\begin{split} A^2 &= X \Big( (X^\top X)^- \Big)^\top H^\top \Big( H(X^\top X)^- H^\top \Big)^{-1} \underbrace{H(X^\top X)^- X^\top X}_{H \text{ (Lemma 4.3.19, 2.))}} \Big( (X^\top X)^- \Big)^T H^\top \cdot \\ & \cdot \Big( H(X^\top X)^- H^\top \Big)^{-1} H(X^\top X)^- X^\top \\ &= X \left( (X^\top X)^- \Big)^\top H^\top \Big( H(X^\top X)^- H^\top \Big)^{-1} H(X^\top X)^- X^\top = A, \end{split}$$

since  $((X^{\top}X)^{-})^{\top}$  is also a generalized inverse of  $X^{\top}X$  (by Lemma 4.3.4). Thus,  $H(X^{\top}X)^{-}H^{\top} = CX(X^{\top}X)^{-}X^{\top}C^{\top}$  does not depend on the choice of  $(X^{\top}X)^{-}$ , cf. proof of Theorem 4.3.14. Using Theorem 4.1.25 yields  $\frac{\varepsilon^{\top}}{\sigma}A\frac{\varepsilon}{\sigma}\sim\chi_{s}^{2}$ , because  $\varepsilon\sim\mathcal{N}(0,\sigma^{2}\mathcal{I})$  and thus  $T\stackrel{H_{0}}{\sim}F_{s,n-r}$ .

# 4.3.6 Confidence regions

Similar to Section 4.2.5, confidence regions for different functions of the parameter vector  $\beta$  can be found. Theorem 4.3.20 directly yields the following confidence region with confidence level  $1 - \alpha \in (0, 1)$ :

**Corollary 4.3.21.** Let  $Y = X\beta + \varepsilon$  be a multivariate regression model with rank (X) = r < m, H an  $(s \times m)$  matrix with rank (H) = s,  $s \in \{1, \ldots, m\}$  and  $H_0 : H\beta = d$  testable  $\forall d \in \mathbb{R}^s$ . Then,

$$\left\{ d \in \mathbb{R}^s : \frac{\left( H\overline{\beta} - d \right)^\top \left( H(X^\top X)^- H^\top \right)^{-1} \left( H\overline{\beta} - d \right)}{s \cdot \overline{\sigma}^2} \le F_{s, n-r, 1-\alpha} \right\}$$

is a confidence region for  $H\beta$  with confidence level  $1-\alpha$ .

Corollary 4.3.22. Let  $h^{\top}\beta$  be an estimable linear function of  $\beta$ ,  $h \in \mathbb{R}^m$ . Then,

$$\left(h^{\top}\overline{\beta} - t_{n-r,\,1-\alpha/2} \cdot \overline{\sigma} \sqrt{h^{\top}(X^{\top}X)^{-}h},\, h^{\top}\overline{\beta} + t_{n-r,\,1-\alpha/2} \cdot \overline{\sigma} \sqrt{h^{\top}(X^{\top}X)^{-}h}\right)$$

is a confidence interval for  $h^{\top}\beta$  with confidence level  $1-\alpha$ .

**Proof** Set s = 1 and  $H = h^{\top}$ . Theorem 4.3.20 implies

$$T = \frac{\left(h^{\top}\overline{\beta} - d\right)^{\top} \left(h^{\top}(X^{\top}X)^{-}h\right)^{-1} \left(h^{\top}\overline{\beta} - d\right)}{\overline{\sigma}^{2}} = \frac{\left(h^{\top}\overline{\beta} - d\right) \left(h^{\top}\overline{\beta} - d\right)}{\overline{\sigma}^{2} \left(h^{\top}(X^{\top}X)^{-}h\right)}$$
$$= \frac{\left(h^{\top}\overline{\beta} - d\right)^{2}}{\overline{\sigma}^{2} \left(h^{\top}(X^{\top}X)^{-}h\right)} \sim F_{1, n-r}$$

under the condition  $h^{\top}\beta = d$ , since  $h^{\top} \left( X^{\top} X \right)^{-} h$  is one-dimensional. Thus, it holds that

$$\sqrt{T} = \frac{h^{\top}\beta - h^{\top}\overline{\beta}}{\overline{\sigma}\sqrt{h^{\top}(X^{\top}X)^{-}h}} \sim t_{n-r}.$$

Therefore,

$$P\left(-t_{n-r,1-\alpha/2} \le \sqrt{T} \le t_{n-r,1-\alpha/2}\right) = 1 - \alpha.$$

This implies the confidence interval above.

An even stronger version of Corollary 4.3.22 can be proven which holds for all h of a linear subspace:

**Theorem 4.3.23** (Confidence band of Scheffé). Let  $H = (h_1, \ldots, h_s)^{\top}$  where  $h_1, \ldots, h_s \in \mathbb{R}^m$ ,  $1 \leq s \leq m$  and  $H_0 : H\beta = d$  testable  $\forall d \in \mathbb{R}^s$ . Let rank (H) = s and  $\mathcal{L} = \langle h_1, \ldots, h_s \rangle$  the linear subspace with span  $\{h_1, \ldots, h_s\}$ . Then

$$P\left(\max_{h\in\mathcal{L}}\left\{\frac{\left(h^{\top}\beta - h^{\top}\overline{\beta}\right)^{2}}{\overline{\sigma}^{2}h^{\top}(X^{\top}X)^{-}h}\right\} \leq sF_{s,n-r,1-\alpha}\right) = 1 - \alpha$$

Thus,

$$\left[h^{\top}\overline{\beta} \pm \sqrt{sF_{s,n-r,1-\alpha}} \cdot \overline{\sigma} \sqrt{h^{\top}(X^{\top}X)^{-}h}\right]$$

is a (uniform with respect to  $h \in \mathcal{L}$ ) confidence band for  $h^{\top}\beta$ .

**Proof** Set

$$T_1 := \left(H\overline{\beta} - H\beta\right)^{\top} \left(H(X^{\top}X)^{-}H^{\top}\right)^{-1} \left(H\overline{\beta} - H\beta\right).$$

Then, Theorem 4.3.20 implies

$$P(T_1 \le s \cdot \overline{\sigma}^2 F_{s, n-r, 1-\alpha}) = 1 - \alpha$$

for all  $\alpha \in (0,1)$ . If it can be shown that

$$T_1 = \max_{x \in \mathbb{R}^s, \, x \neq 0} \left\{ \frac{\left( x^\top \left( H\overline{\beta} - H\beta \right) \right)^2}{x^\top \left( H(X^\top X)^- H^\top \right) x} \right\},\tag{4.12}$$

then the proof is concluded, since

$$1 - \alpha = P\left(T_1 \leq \underbrace{s\overline{\sigma}^2 F_{s, n-r, 1-\alpha}}_{t}\right)$$

$$= P\left(\max_{x \in \mathbb{R}^s, x \neq 0} \left\{ \frac{\left(x^\top \left(H\overline{\beta} - H\beta\right)\right)^2}{x^\top \left(H(X^\top X)^- H^\top\right) x} \right\} \leq t\right)$$

$$= P\left(\max_{x \in \mathbb{R}^s, x \neq 0} \left\{ \frac{\left((H^\top x)^\top \overline{\beta} - (H^\top x)^\top \beta\right)^2}{(H^\top x)^\top (X^\top X)^- (H^\top x)} \right\} \leq t\right)$$

$$H^\top x = h \in \mathcal{L} P\left(\max_{h \in \mathcal{L}} \left\{ \frac{\left(h^\top \overline{\beta} - h^\top \beta\right)^2}{h^\top (X^\top X)^- h} \right\} \leq s\overline{\sigma}^2 F_{s, n-r, 1-\alpha}\right).$$

In order to prove the validity of (4.12) it is sufficient to show that  $T_1$  is an upper bound of

$$\frac{\left(x^{\top}(H\overline{\beta} - H\beta)\right)^{2}}{x^{\top}\left(H(X^{\top}X) - H^{\top}\right)x},$$

which is also a maximum. Since  $H(X^{\top}X)^{-}H^{\top}$  is positive definite and invertible, there exists an invertible  $(s \times s)$  matrix B with the property  $BB^{\top} = H(X^{\top}X)^{-}H^{\top}$ . Then,

$$(x^{\top}(H\overline{\beta} - H\beta))^{2} = (\underbrace{x^{\top}B}_{(B^{\top}x)^{\top}} \cdot B^{-1}(H\overline{\beta} - H\beta))^{2}$$

$$\leq |B^{\top}x|^{2} \cdot |B^{-1}(H\overline{\beta} - H\beta)|^{2}$$

$$= x^{\top}BB^{\top}x \left(H\overline{\beta} - H\beta\right)^{\top} \cdot \underbrace{(B^{-1})^{\top}B^{-1}(H\overline{\beta} - H\beta)}_{= (B^{\top})^{-1}B^{-1} = (BB^{\top})^{-1}}$$

$$= x^{\top}H(X^{\top}X)^{-}H^{\top}x \cdot \left(H\overline{\beta} - H\beta\right)^{\top} \left(H(X^{\top}X)^{-}H^{\top}\right)^{-1} (H\overline{\beta} - H\beta).$$

Thus, it holds that

$$\frac{\left(x^{\top}(H\overline{\beta} - H\beta)\right)^{2}}{x^{\top}(H(X^{\top}X)^{-}H^{\top})x} \leq \left(H\overline{\beta} - H\beta\right)^{\top}\left(H(X^{\top}X)^{-}H^{\top}\right)^{-1}\left(H\overline{\beta} - H\beta\right)$$
$$= T_{1}.$$

For  $x = (H(X^{\top}X)^{-}H^{\top})^{-1}(H\overline{\beta} - H\beta)$  it can be shown that it actually is a maximum.

# 4.3.7 Introduction to variance analysis

In this section we discuss an example for the application of linear models with a design matrix that doesn't have full rank. It is the assertion of the variability of the expected values in the random sample  $Y = (Y_1, \ldots, Y_n)^{\top}$  in short ANOVA (analysis of variance).

First, consider the single factor variance analysis, in which it is assumed, that the random sample  $(Y_1, \ldots, Y_n)$  can be partitioned in k homogeneous subclasses  $(Y_{ij}, j = 1, \ldots, n_i), i = 1, \ldots, k$  with

1. 
$$\mathbb{E}(Y_{ij}) = \mu_i = \mu + \alpha_i, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k.$$

2. 
$$n_i > 1$$
,  $i = 1, ..., k$ ,  $\sum_{i=1}^k n_i = n$ ,  $\sum_{i=1}^k n_i \alpha_i = 0$ .

Here  $\mu$  is a factor, which is equal in all classes and  $\alpha_i \in \mathbb{R}$  are the class specific differences between the expected values  $\mu_1, \ldots, \mu_k$ . The number  $i = 1, \ldots, k$  of the classes are called levels of the influencing factor (e.g. the doses of a drug in a clinical trial) and  $\alpha_i$ ,  $i = 1, \ldots, k$  can be interpreted as the effect of the *i*-th level. The constraint  $\sum_{i=1}^{k} n_i \alpha_i = 0$  causes that the conver-

sion  $(\mu_1, \ldots, \mu_k) \longleftrightarrow (\mu, \alpha_1, \ldots, \alpha_k)$  is unique and that  $\mu = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} \mathbb{E} Y_{ij}$ . Furter, it is assumed that  $\mu_i$  can be measured with uncorrelated measurement errors  $\varepsilon_{ij}$ , i.e.

$$Y_{ij} = \mu_i + \varepsilon_{ij} = \mu + \alpha_i + \varepsilon_{ij}, \quad i = 1, \dots, k, j = 1, \dots, n_i$$
 (4.13)  
 $\mathbb{E} \, \varepsilon_{ij} = 0, \quad \text{Var} \, \varepsilon_{ij} = \sigma^2, \quad \varepsilon_{ij} \text{ uncorrelated, } i = 1, \dots, k, j = 1, \dots, n_i.$  (4.14)

Here, we want to test the classical ANOVA hypothesis that no variability in the expected values  $\mu_i$  can be found, i.e.

$$H_0: \quad \mu_1 = \mu_2 = \ldots = \mu_k,$$

which means, that

$$H_0: \quad \alpha_1 = \alpha_2 = \ldots = \alpha_k.$$

The constraint

$$\sum_{i=1}^{k} n_i \alpha_i = 0.$$

implies  $\alpha_i = 0$ .

The problem (4.13) can be rewritten in terms of multivariate linear regression as follows:

$$Y = X\beta + \varepsilon, \text{ where } Y = (Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2}, \dots, Y_{k1}, \dots, Y_{kn_k})^{\top},$$

$$\beta = (\mu, \alpha_1, \dots, \alpha_k)^{\top},$$

$$\varepsilon = (\varepsilon_{11}, \dots, \varepsilon_{1n_1}, \dots, \varepsilon_{k1}, \dots, \varepsilon_{kn_k})^{\top},$$

$$\begin{pmatrix} 1 & 1 & 0 & \dots & \dots & 0 \\ 1 & 1 & 0 & \dots & \dots & 0 \\ \vdots & & & & & \\ 1 & 1 & 0 & \dots & \dots & 0 \\ \vdots & & & & & \\ 1 & 0 & 1 & 0 & \dots & 0 \\ \vdots & & & & & \\ 1 & 0 & \dots & \dots & 0 & 1 \\ \vdots & & & & & \\ 1 & 0 & \dots & \dots & 0 & 1 \\ \vdots & & & & & \\ 1 & 0 & \dots & \dots & 0 & 1 \\ \end{pmatrix} \begin{array}{c} n_1 \\ \vdots \\ n_k \\ \vdots \\ n$$

The  $(n \times (k+1))$  matrix X has rank k < m = k+1; thus the theory of section 4.3 can be applied to this model.

Exercise 4.3.24. Show that the ANOVA hypothesis

$$H_0: \quad \alpha_i = 0, \quad \forall i = 1, \dots, k$$

is not testable!

In order to consider an equivalent but testable hypothesis

$$H_0: \quad \alpha_1 - \alpha_2 = 0, \dots, \alpha_1 - \alpha_k = 0 \quad \text{resp.} \quad H_0: \quad H\beta = 0,$$

we construct a  $(k-1) \times (k+1)$  matrix

$$H = \begin{pmatrix} 0 & 1 & -1 & 0 & \dots & 0 \\ 0 & 1 & 0 & -1 & \dots & 0 \\ \vdots & & & & & \\ 0 & 1 & 0 & \dots & -1 & 0 \\ 0 & 1 & 0 & \dots & 0 & -1 \end{pmatrix}$$

which we can use as part of our hypothesis test. (Show that!)

In the two factor variance analysis the random sample  $(Y_1, \ldots, Y_n)$  is divided in  $k_1 \cdot k_2$  homogeneous groups depending on two factors

$$Y_{i_1 i_2 j}, \quad j = 1, \dots, n_{i_1 i_2}$$

for  $i_1 = 1, ..., k_1, i_2 = 1, ..., k_2$ , such that

$$\sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} n_{i_1 i_2} = n.$$

Here it is assumed that

$$\mathbb{E} Y_{i_1 i_2 j} = \mu_{i_1 i_2} = \mu + \alpha_{i_1} + \beta_{i_2} + \gamma_{i_1 i_2}, \quad i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2,$$

thus the following linear model is constructed:

$$Y_{i_1 i_2 j} = \mu_{i_1 i_2} + \varepsilon_{i_1 i_2 j} = \mu + \alpha_{i_1} + \beta_{i_2} + \gamma_{i_1 i_2} + \varepsilon_{i_1 i_2 j},$$
  
$$j = 1, \dots, n_{i_1 i_2}, i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2.$$

**Exercise 4.3.25.** Write down the design matrix X for this case explicitly and show that it also doesn't have full rank.

# Chapter 5

# Generalized linear models

Another class of regression models usually allows for an arbitrary functional connection g between the mean of the goal variable  $EY_i$  and the linear part  $X\beta$ , which is a linear combinations of the entries of the design matrix  $X = (x_{ij})$  and the parameter vector  $\beta = (\beta_1, \ldots, \beta_m)^{\top}$ . On the other hand it allows distributions of  $Y_i$ , which are not necessarily based on the normal distribution (and functions of those). Thus it is possible to consider data  $Y_i$  that has a finite number of characteristics (e.g. "yes" and "no" in economic surveys). The class of all possible distributions is bounded by the Exponential family.

Let  $Y_1, \ldots, Y_n$  be a random sample of the goal variable of the model and let

$$X = (x_{ij})_{\substack{i=1,\dots,n\\j=1,\dots,m}}$$

the design matrix of the output variables, which are not random.

**Definition 5.0.1.** The generalized linear model is given by

$$(g(\mathbb{E}Y_1), \dots, g(\mathbb{E}Y_n))^{\top} = X\beta \quad \text{with } \beta = (\beta_1, \dots, \beta_m)^{\top}$$
 (5.1)

where  $g:G\subset\mathbb{R}\to\mathbb{R}$  is the so called *link function* with domain G. The rank is given by  $\operatorname{rank}(X)=m$ .

Under the assumption that g is known explicitly, the parameter vector  $\beta$  is desired to be estimated using  $(Y_1, \ldots, Y_n)$ . Here it is assumed that  $Y_i$ ,  $i = 1, \ldots, n$  are independent but not necessarily identically distributed. But their distribution is a member of the following family of distributions.

## 5.1 Exponential family of distributions

**Definition 5.1.1.** The distribution of a random variable Y is a member of the *exponential family*, if the functions  $a : \mathbb{R} \times \mathbb{R}_+ \to \mathbb{R}$  and  $b : \Theta \to \mathbb{R}$  exist, such that

• in the absolutely continuous case the probability density function of Y is given by

$$f_{\theta}(y) = \exp\left\{\frac{1}{\tau^2}(y\theta + a(y,\tau) - b(\theta))\right\}, \quad y \in \mathbb{R}$$
 (5.2)

• in the discrete case the probability mass function of Y is given by

$$P_{\theta}(Y=y) = \exp\left\{\frac{1}{\tau^2}(y\theta + a(y,\tau) - b(\theta))\right\}, y \in C$$
 (5.3)

where C is the (at most) countable domain of Y,  $\tau^2$  the so called error term,  $\theta \in \Theta \subset \mathbb{R}$  a parameter and

$$\Theta = \left\{ \theta \in \mathbb{R} : \int_{\mathbb{R}} \exp\left\{ \frac{y\theta + a(y,\tau)}{\tau^2} \right\} dy < \infty \right\}$$

respectively in the discrete case:

$$\Theta = \left\{ \theta \in \mathbb{R} : \sum_{y \in C} \exp\left\{ \frac{y\theta + a(y, \tau)}{\tau^2} \right\} < \infty \right\}$$

which is the natural parameter space with at least two different elements.

#### **Lemma 5.1.2.** $\Theta$ is an interval.

**Proof** Show that  $\Theta \subset \mathbb{R}$  is convex. It is then (possibly an infinite) interval. For arbitrary  $\theta_1, \theta_2 \in \Theta$  (at least one pair exists by Definition 5.1.1) it holds that  $\alpha\theta_1 + (1-\alpha)\theta_2 \in \Theta$  for all  $\alpha \in (0,1)$ . In order show that the statement above holds, suppose that the distribution of Y is absolutely continuous. Since  $\theta_i \in \Theta$ , it holds that

$$\int_{\mathbb{D}} \exp\left\{\frac{1}{\tau^2} \left(y\theta_i + a(y,\tau)\right)\right\} dy < \infty, \quad i = 1, 2.$$

The inequality

$$\alpha x_1 + (1 - \alpha)x_2 \le \max\{x_1, x_2\}, \quad x_1, x_2 \in \mathbb{R}, \quad \alpha \in (0, 1),$$

implies

$$\exp\left\{\frac{1}{\tau^2}\left(y(\alpha\theta_1 + (1-\alpha)\theta_2) + a(y,\tau)\right)\right\} \\
= \exp\left\{\alpha\frac{1}{\tau^2}\left(y\theta_1 + a(y,\tau)\right) + (1-\alpha)\frac{1}{\tau^2}\left(y\theta_2 + a(y,\tau)\right)\right\} \\
\leq \max_{i=1,2} \exp\left\{\frac{1}{\tau^2}\left(y\theta_i + a(y,\tau)\right)\right\} \\
\leq \exp\left\{\frac{1}{\tau^2}\left(y\theta_1 + a(y,\tau)\right)\right\} + \exp\left\{\frac{1}{\tau^2}\left(y\theta_2 + a(y,\tau)\right)\right\},$$

and thus

$$\int_{\mathbb{R}} \exp\left\{\frac{1}{\tau^2} \left(y(\alpha \theta_1 + (1 - \alpha)\theta_2) + a(y, \tau)\right)\right\} dy$$

$$\leq \sum_{i=1}^2 \int_{\mathbb{R}} \exp\left\{\frac{1}{\tau^2} \left(y\theta_i + a(y, \tau)\right)\right\} dy < \infty$$

by the assumptions of the lemma. In summary

$$\alpha\theta_1 + (1-\alpha)\theta_2 \in \Theta$$

and  $\Theta$  is a interval.

**Example 5.1.3.** Which distributions belong to the exponential family?

1. Normal distribution: If  $Y \sim \mathcal{N}(\mu, \sigma^2)$ , then  $\mu$  is the parameter of interest and  $\sigma^2$  the error term. It holds that

$$f_{\mu}(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

$$= \exp\left\{\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2}\left(\frac{y^2}{\sigma^2} - \frac{2y\mu}{\sigma^2} + \frac{\mu^2}{\sigma^2}\right)\right\}$$

$$= \exp\left\{\frac{1}{\sigma^2}\left(y\mu - \frac{y^2}{2} - \left(\frac{\mu^2}{2} + \frac{\sigma^2}{2}\log(2\pi\sigma^2)\right)\right)\right\}$$

and setting  $\theta = \mu$ ,  $\tau = \sigma$ ,

$$a(y,\tau) = -\frac{y^2}{2} - \frac{\sigma^2}{2}\log(2\pi\sigma^2)$$
 and  $b(\mu) = b(\theta) = \frac{\mu^2}{2}$ 

satisfies Equation (5.2)

2. Bernoulli distribution:  $Y \sim \text{Bernoulli}(p), p \in [0; 1]$ The Bernoulli distribution is usually used in surveys in market research

The Bernoulli distribution is usually used in surveys in market research where

$$Y = \begin{cases} 1, & \text{if the answer is "yes"} \\ 0, & \text{if the answer is "no"} \end{cases}$$

for a given question in the respective survey.

Here the probabilities are given by  $P(Y=1)=p,\ P(Y=0)=1-p.$ Then for  $y\in\{0,1\}$  it holds that:

$$P_{\theta}(Y = y) = p^{y} (1 - p)^{1 - y} = e^{y \log p + (1 - y) \log(1 - p)}$$
$$= e^{y \log \frac{p}{1 - p} - (-\log(1 - p))}$$

Thus the Bernoulli distribution is a member of the exponential family with  $\theta = \log \frac{p}{1-p}$ ,  $\tau = 1$ ,

$$a(y,\tau) = 0$$
,  $b(\theta) = -\log(1-p) = \log(1+e^{\theta})$ .

3. **Poisson distribution**: If  $Y \sim \text{Poisson } (\lambda), \lambda > 0$ , then for  $y \in \mathbb{N}_0$ 

$$P_{\theta}(Y = y) = e^{-\lambda} \cdot \frac{\lambda^y}{y!} = e^{y \log \lambda - \log(y!) - \lambda}.$$

Thus the Poisson distribution is a member of the exponential family with  $\theta = \log \lambda$ ,  $\tau = 1$ ,

$$a(y,\tau) = -\log(y!), \quad b(\theta) = \lambda = e^{\theta}.$$

**Lemma 5.1.4.** If the distribution of a random variable Y is a member of the exponential family,  $\mathbb{E}Y^2 < \infty$  and  $b : \Theta \to \mathbb{R}$  is two times continuously differentiable with  $b''(\theta) > 0$  for all  $\theta \in \Theta$ , then

$$\mathbb{E}Y = b'(\theta), \quad \text{Var}Y = \tau^2 b''(\theta).$$

#### Proof

1. Only the case for absolutely continuous distributions is discussed below. The discrete case can be handled simultaneously by replacing the  $\int$  with  $\sum$ . It holds that

$$\begin{aligned} & \mathrm{E}Y = \int_{\mathbb{R}} y f_{\theta}(y) dy = \int_{\mathbb{R}} y \exp\left\{\frac{1}{\tau^{2}} \left(y \theta + a(y, \tau) - b(\theta)\right)\right\} dy \\ & = e^{-\frac{b(\theta)}{\tau^{2}}} \cdot \tau^{2} \int_{\mathbb{R}} \frac{\partial}{\partial \theta} \exp\left\{\frac{1}{\tau^{2}} \left(y \theta + a(y, \tau)\right)\right\} dy \\ & = e^{-\frac{b(\theta)}{\tau^{2}}} \cdot \tau^{2} \frac{\partial}{\partial \theta} \int_{\mathbb{R}} \exp\left\{\frac{1}{\tau^{2}} \left(y \theta + a(y, \tau)\right)\right\} dy \end{aligned} \\ & = e^{-\frac{b(\theta)}{\tau^{2}}} \cdot \tau^{2} \frac{\partial}{\partial \theta} \left(e^{\frac{b(\theta)}{\tau^{2}}} \underbrace{\int_{\mathbb{R}} \exp\left\{\frac{1}{\tau^{2}} \left(y \theta + a(y, \tau) - b(\theta)\right)\right\} dy}_{\int_{\mathbb{R}} f_{\theta}(y) dy = 1} \right) \\ & = e^{-\frac{b(\theta)}{\tau^{2}}} \tau^{2} \frac{\partial}{\partial \theta} \left(e^{\frac{b(\theta)}{\tau^{2}}}\right) = e^{-\frac{b(\theta)}{\tau^{2}}} \cdot \tau^{2} \frac{b'(\theta)}{\tau^{2}} e^{\frac{b(\theta)}{\tau^{2}}} = b'(\theta). \end{aligned}$$

2. Show (analogously to 1) that

#### Exercise 5.1.5.

$$Var Y = \tau^2 b''(\theta).$$

#### 5.2 Link functions

The goal variables  $Y_i$ ,  $i=1,\ldots,n$  are i.i.d. with a distribution which is a member of the exponential family and a probability (density or) mass function as in (5.3) resp. (5.2). Assume that  $b:\Theta\to\mathbb{R}$  is two times continuously differentiable with  $b''(\theta)>0$  for all  $\theta\in\Theta$ . Additionally assume a generalized linear model as in (5.1).

**Definition 5.2.1.** (Natural link functions) The link function  $g: G \to \mathbb{R}$  is called *natural*, if  $g = (b')^{-1}$ ,  $G = \{b'(\theta) : \theta \in \Theta\}$  and g is two times continuously differentiable with  $g'(x) \neq 0$  for all  $x \in G$ .

The question why the link function is called "natural" is answered in the following Lemma.

**Lemma 5.2.2.** If the generalized linear model (5.1) has the natural link function, then  $(\theta_1, \dots, \theta_n)^{\top} = X\beta$ 

**Proof** Since  $b''(\theta) > 0$ , it holds that  $b'(\theta)$  is monotonically increasing and thus invertible. Define

$$\mu_i = \mathbb{E}Y_i, \quad \eta_i = x_i^{\top}\beta, \quad x_i = (x_{i1}, \dots, x_{im})^{\top}, \quad i = 1, \dots, n.$$

Since g is invertible it holds that

$$\mu_i = g^{-1}(x_i^{\top}\beta) = g^{-1}(\eta_i), \quad i = 1, \dots, n.$$

On the other hand, it holds that  $\mu_i = b'(\theta_i)$  by Lemma 5.1.4, so

$$b'(\theta_i) = g^{-1}(\eta_i) \stackrel{\text{Def. 5.2.1}}{=} b'(\eta_i), \quad i = 1, \dots, n.$$

Because of the monotonicity of b' the assertion  $\theta_i = \eta_i$ , i = 1, ..., n holds.

**Example 5.2.3.** In the following the natural link function for the distributions of Example 5.1.3.

1. Normal distribution: Since  $b(\mu) = \frac{\mu^2}{2}$ , it holds that

$$b'(x) = \frac{2x}{2} = x$$
 and thus  $g(x) = (b')^{-1}(x) = x$ 

The natural link function is given by g(x) = x, thus it holds that

$$(\mu_1,\ldots,\mu_n)^{\top}=(\mathbb{E}Y_1,\ldots,\mathbb{E}Y_n)^{\top}=X\beta$$

This is exactly the case of linear regression.

2. Bernoulli distribution: Since  $b(\theta) = \log(1 + e^{\theta})$ , it holds that

$$b'(x) = \frac{1}{1 + e^x} \cdot e^x = y$$

$$\Leftrightarrow \frac{1}{e^{-x} + 1} = y$$

$$\Leftrightarrow \frac{1}{y} - 1 = e^{-x}$$

$$\Leftrightarrow x = -\log \frac{1 - y}{y} = \log \frac{y}{1 - y}$$

$$\Rightarrow g(x) = (b')^{-1}(x) = \log \frac{x}{1 - x}$$

The generalized linear regression model in the case of the Bernoulli distribution is called *binary* (categorical) regression. If it is used with the natural link function, it is called *logistic regression*. In this case it holds that

$$(p_1, \dots, p_n)^{\top} = (\mathbb{E}Y_1, \dots, \mathbb{E}Y_n)^{\top},$$

$$\theta_i = \log \frac{p_i}{1 - p_i} = x_i^{\top} \beta, \quad i = 1, \dots, n$$

$$\Leftrightarrow e^{\theta_i} = \frac{p_i}{1 - p_i}$$

$$\Leftrightarrow p_i = \frac{e^{\theta_i}}{1 + e^{\theta_i}}$$

$$\Leftrightarrow p_i = \frac{e^{x_i^{\top} \beta}}{1 + e^{x_i^{\top} \beta}}, \quad i = 1, \dots, n.$$

The ratio

$$\frac{p_i}{1 - p_i} = \frac{P(Y_i = 1)}{P(Y_i = 0)}, \quad i = 1, \dots, n$$

is called *Odds*. The logarithm of the Odds is called *Logit*:

$$\log \frac{p_i}{1 - p_i}, \quad i = 1, \dots, n.$$

Logits are thus "new goal variables", which are estimated by using the linear combinations  $x_i^{\mathsf{T}}\beta$ .

An alternative link function which is often used is defined by  $g(x) = \Phi^{-1}(x)$ , which is the *Quantile function of the normal distribution*. It is however, not a natural link function. By using them, the so called *Probit model*:

$$p_i = \Phi(x_i^{\top}\beta), \quad i = 1, \dots, n$$

can be obtained.

3. **Poisson distribution**: Since  $b(\theta) = e^{\theta}$ , in this case

$$g(x) = (b')^{-1}(x) = \log x, \quad x > 0$$

is the natural link function. Thus the generalized linear model with natural link function has the representation

$$(\log \lambda_1, \dots, \log \lambda_n)^{\top} = X\beta \quad \text{or} \quad \lambda_i = e^{x_i^{\top}\beta}, i = 1, \dots, n.$$

### 5.3 Maximum likelihood estimator for $\beta$

Since the probability mass (or density) function of  $Y_i$  is given by

$$\exp\left\{\frac{1}{\tau^2}(y\theta_i + a(y,\tau) - b(\theta_i))\right\}$$

and the  $Y_i$  are independent, the log-likelihood function of the random sample  $Y = (Y_1, \ldots, Y_n)$  can be written as follows

$$\log L(Y, \theta) = \log \prod_{i=1}^{n} f_{\theta_i}(Y_i) = \frac{1}{\tau^2} \sum_{i=1}^{n} (Y_i \theta_i + a(Y_i, \tau) - b(\theta_i)).$$
 (5.4)

The proof of Lemma 5.2.2 implies that

$$\theta_i = (b')^{-1} (g^{-1}(x_i^{\mathsf{T}}\beta)), \quad i = 1, \dots, n,$$
 (5.5)

which implies that  $\log L(Y, \theta)$  is a function of the parameter  $\beta$ . From now on the notation  $\log L(Y, \beta)$  is used to emphasize this fact.

The maximum likelihood estimator  $\hat{\beta}$  for  $\beta$  is desired:

$$\hat{\beta} = \operatorname*{argmax}_{\beta} \log L(Y, \beta)$$

The necessary condition for an extrema

$$\frac{\partial \log L(Y,\beta)}{\partial \beta_i} = 0, \quad i = 1, \dots, m,$$

needs to be checked. Introduce the following notation

$$U_{i}(\beta) = \frac{\partial \log L(Y, \beta)}{\partial \beta_{i}}, \quad i = 1, \dots, m,$$
  

$$U(\beta) = (U_{1}(\beta), \dots, U_{m}(\beta))^{\top},$$
  

$$I_{ij}(\beta) = \mathbb{E}[U_{i}(\beta)U_{j}(\beta)], \quad i, j = 1, \dots, m,$$

#### Definition 5.3.1.

- 1. The matrix  $I(\beta) = (I_{ij}(\beta))_{i,j=1}^m$  is called Fisher information matrix.
- 2. Introduce the Hesse matrix  $W(\beta)$  as a random matrix

$$W(\beta) = (W_{ij}(\beta))_{i,j=1}^m \text{ with } W_{ij}(\beta) = \frac{\partial^2}{\partial \beta_i \partial \beta_j} \log L(Y, \beta).$$

This  $(m \times m)$  matrix contains the 2nd order derivative of the loglikelihood function, which will be relevant for solving the maximisation problem

$$\log L(Y,\beta) \to \max_{\beta}$$
.

**Theorem 5.3.2.** It can be shown that  $U(\beta)$  and  $I(\beta)$  have the following explicit form.

1.

$$U_j(\beta) = \sum_{i=1}^n x_{ij} \left( Y_i - \mu_i(\beta) \right) \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \frac{1}{\sigma_i^2(\beta)}, \quad j = 1, \dots, m$$

2.

$$I_{jk}(\beta) = \sum_{i=1}^{n} x_{ij} x_{ik} \left( \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \right)^2 \frac{1}{\sigma_i^2(\beta)}, \quad j, k = 1, \dots, m ,$$

where  $\eta_i = x_i^{\top} \beta$ ,  $\mu_i(\beta) = g^{-1}(x_i^{\top} \beta)$  is the expectation of  $Y_i$  and

$$\sigma_i^2(\beta) \stackrel{\text{Lemma 5.1.4}}{=} \tau^2 b''(\theta_i) \stackrel{(5.5)}{=} \tau^2 b''((b')^{-1}(g^{-1}(x_i^{\top}\beta))), \quad i = 1, \dots, n,$$

is the variance of  $Y_i$ .

#### Proof

1. Introduce the notation

$$l_i(\beta) = \frac{1}{\tau^2} \left( Y_i \theta_i + a \left( Y_i, \tau \right) - b(\theta_i) \right), \quad i = 1, \dots, n.$$

Then,

$$U_j(\beta) = \sum_{i=1}^n \frac{\partial l_i(\beta)}{\partial \beta_j}, \quad j = 1, \dots, m.$$

Applying the chain rule several times yields

$$\frac{\partial l_i(\beta)}{\partial \beta_j} = \frac{\partial l_i(\beta)}{\partial \theta_i} \cdot \frac{\partial \theta_i}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_j}, \quad i = 1, \dots, n, \ j = 1, \dots, m.$$

Since

$$\frac{\partial l_i(\beta)}{\partial \theta_i} = \frac{1}{\tau^2} \Big( Y_i - b'(\theta_i) \Big) \overset{\text{Lemma 5.1.4}}{=} \frac{1}{\tau^2} \Big( Y_i - \mu_i(\beta) \Big),$$

$$\frac{\partial \theta_i}{\partial \mu_i} = \left( \frac{\partial \mu_i}{\partial \theta_i} \right)^{-1} = \left( \left( b'(\theta_i) \right)' \right)^{-1} = \left( b''(\theta_i) \right)^{-1} \overset{\text{Lemma 5.1.4}}{=} \frac{\sigma_i^2(\beta)}{\sigma_i^2(\beta)} \frac{\partial \mu_i}{\partial \eta_i} = \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i}$$

because  $\mu_i = EY_i = g^{-1}(\eta_i)$ ,

$$\frac{\partial \eta_i}{\partial \beta_j} = \frac{\partial (x_i^\top \beta)}{\partial \beta_j} = x_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, m,$$

we have

$$U_{j}(\beta) = \frac{1}{\tau^{2}} \sum_{i=1}^{n} x_{ij} (Y_{i} - \mu_{i}(\beta)) \cdot \frac{\tau^{2}}{\sigma_{i}^{2}(\beta)} \cdot \frac{\partial g^{-1}(\eta_{i})}{\partial \eta_{i}}$$
$$= \sum_{i=1}^{n} x_{ij} (Y_{i} - \mu_{i}(\beta)) \frac{\partial g^{-1}(\eta_{i})}{\partial \eta_{i}} \cdot \frac{1}{\sigma_{i}^{2}(\beta)}, \quad j = 1, \dots, m.$$

2. For all  $i, j = 1, \ldots, m$  it holds that

$$\begin{split} I_{ij}(\beta) &= \mathrm{E}(U_i(\beta)U_j(\beta)) \\ &= \sum_{k,l=1}^n x_{ki} x_{lj} \underbrace{\mathrm{Cov} \ (Y_k, Y_l)}_{\delta_{k_l} \sigma_k^2(\beta)} \cdot \frac{\partial g^{-1}(\eta_k)}{\partial \eta_k} \frac{\partial g^{-1}(\eta_l)}{\partial \eta_l} \frac{1}{\sigma_k^2(\beta) \sigma_l^2(\beta)} \\ &= \sum_{k=1}^n x_{ki} x_{kj} \left( \frac{\partial g^{-1}(\eta_k)}{\partial \eta_k} \right)^2 \frac{1}{\sigma_k^2(\beta)}. \end{split}$$

**Remark 5.3.3.** In case of the natural link function, simplify the equation above so that the log-likelihood function is given by

$$\log L(Y, \beta) = \frac{1}{\tau^2} \sum_{i=1}^{n} \left( Y_i x_i^{\top} \beta + a(Y_i, \tau) - b(x_i^{\top} \beta) \right).$$

Since in this case  $g^{-1}(\eta_i) = b'(\eta_i), \ \eta_i = x_i^{\top} \beta = \theta_i \text{ holds},$ 

$$\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} = b''(\theta_i) \stackrel{\text{Lemma 5.1.4}}{=} \frac{1}{\tau^2} \sigma_i^2(\beta),$$

and thus

$$U_{j}(\beta) = \frac{1}{\tau^{2}} \sum_{i=1}^{n} x_{ij} (Y_{i} - \mu_{i}(\beta)), \quad j = 1, \dots, m,$$
  
$$I_{jk}(\beta) = \frac{1}{\tau^{4}} \sum_{i=1}^{n} x_{ij} x_{ik} \sigma_{i}^{2}(\beta), \quad j, k = 1, \dots, m.$$

#### Theorem 5.3.4.

$$W_{jk}(\beta) = \sum_{i=1}^{n} x_{ij} x_{ik} \left( \left( Y_i - \mu_i(\beta) \right) \nu_i - \frac{u_i^2}{\sigma_i^2(\beta)} \right), \quad j, k = 1, \dots, m$$

where

$$u_{i} = \frac{\partial g^{-1}(\eta_{i})}{\partial \eta_{i}}$$

$$\nu_{i} = \frac{1}{\tau^{2}} \cdot \frac{\partial^{2}((b')^{-1} \circ g^{-1}(\eta_{i}))}{\partial \eta_{i}^{2}}$$

$$\mu_{i}(\beta) = EY_{i}, \ \sigma_{i}^{2}(\beta) = VarY_{i},$$

$$\eta_{i} = x_{i}^{\top}\beta$$

for  $i = 1, \ldots, n$ .

**Proof** For arbitrary j, k = 1, ..., m it holds that

$$\begin{split} W_{jk}(\beta) &= \frac{\partial}{\partial \beta_k} U_j(\beta) \overset{\text{Theorem 5.3.2}}{=} \frac{\partial}{\partial \beta_k} \sum_{i=1}^n x_{ij} \left( Y_i - \mu_i(\beta) \right) \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \frac{1}{\sigma_i^2(\beta)} \\ &= \sum_{i=1}^n x_{ij} \left( (Y_i - \mu_i(\beta)) \frac{\partial}{\partial \beta_k} \left( \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \frac{1}{\sigma_i^2(\beta)} \right) - \\ &\qquad - \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \frac{1}{\sigma_i^2(\beta)} \frac{\partial \mu_i(\beta)}{\partial \beta_k} \right) \\ &= \sum_{i=1}^n \left( x_{ij} (Y_i - \mu_i(\beta)) \frac{\partial}{\partial \beta_k} \left( \frac{1}{\tau^2} \frac{\partial \theta_i}{\partial \eta_i} \right) \right. \\ &\qquad - \left( \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \right)^2 \frac{1}{\sigma_i^2(\beta)} x_{ik} \right) \\ &= \sum_{i=1}^n x_{ij} x_{ik} \left( (Y_i - \mu_i(\beta)) \nu_i - u_i^2 \frac{1}{\sigma_i^2(\beta)} \right), \end{split}$$

where

$$\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot \frac{1}{\sigma_i^2(\beta)} \xrightarrow{\text{Lemma 5.1.4}}_{\text{and Theorem 5.3.2}} \frac{\partial b'(\theta_i)}{\partial \eta_i} \cdot \frac{1}{\tau^2} \cdot \frac{1}{b''(\theta_i)}$$

$$= \frac{\partial b'(\theta_i)}{\partial \theta_i} \cdot \frac{\partial \theta_i}{\partial \eta_i} \frac{1}{\tau^2} \frac{1}{b''(\theta_i)} = \frac{1}{\tau^2} \frac{\partial \theta_i}{\partial \eta_i}$$

and

$$\frac{\partial}{\partial \beta_k} \left( \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot \frac{1}{\sigma_i^2(\beta)} \right) = \frac{1}{\tau^2} \frac{\partial^2 \theta_i}{\partial \eta_i^2} \cdot \frac{\partial \eta_i}{\partial \beta_k} \stackrel{\eta_i = x_i^\top \beta}{=} \frac{1}{\tau^2} \frac{\partial^2 \theta_i}{\partial \eta_i^2} \cdot x_{ik},$$

with

$$\frac{\partial \overbrace{g^{-1}(\eta_i)}^{\mu_i(\beta)}}{\partial \beta_k} = \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_k} = \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot x_{ik}$$

and 
$$\theta_i = (b')^{-1} \circ g^{-1}(\eta_i), i = 1, \dots, n.$$

Moreover for generalized linear models with natural link function

$$W(\beta) = -I(\beta) = \left(-\frac{1}{\tau^4} \sum_{i=1}^n x_{ij} x_{ik} \sigma_i^2(\beta)\right)_{i,k=1}^m$$
 (5.6)

holds, since in this case  $\nu_i = 0$  for all i = 1, ..., n.  $W(\beta)$  is therefore deterministic. Indeed, by Lemma 5.2.2  $\theta_i = x_i^{\top}\beta = \eta_i$  and thus  $\frac{\partial^2 \theta_i}{\partial \eta_i^2} = 0$ , i = 1, ..., n.

Remark 5.3.3 implies  $u_i^2 = \frac{1}{\tau^4} \sigma_i^4(\beta)$ .

**Example 5.3.5.** What do  $U(\beta)$ ,  $I(\beta)$  and  $W(\beta)$  look like for the models introduced in Example 5.2.3 (natural link function)?

1. Normal distribution: this case corresponds to the usual multivariate linear regression with normally distributed error terms. In this case it holds that  $\mu = X\beta$ ,  $\tau^2 = \sigma^2$ .

Remark 5.3.3 implies

$$U(\beta) = \frac{1}{\sigma^2} X^{\top} (Y - X\beta),$$
  

$$I(\beta) = (\mathbb{E} (U_i(\beta) \cdot U_j(\beta)))_{i,j=1,\dots,m} = \frac{1}{\sigma^2} X^{\top} X,$$
  

$$W(\beta) = -I(\beta).$$

2. **Logistic regression**: It holds that  $\tau^2 = 1$ ,  $\mu_i = p_i$ ,  $\sigma_i^2 = p_i(1 - p_i)$ ,  $i = 1, ..., n, p_i \in (0, 1)$  and thus

$$U(\beta) = X^{\top}(Y - p)$$
  

$$I(\beta) = X^{\top} diag(p_i(1 - p_i))X$$
  

$$W(\beta) = -I(\beta)$$

where  $p = (p_1, \ldots, p_n)^{\top}$ .

3. **Poisson regression**: It holds that  $\tau^2 = 1$ ,  $\mu_i = \lambda_i = \sigma_i^2$ ,  $i = 1, \ldots, n$  and thus

$$U(\beta) = X^{\top}(Y - \lambda),$$
  

$$I(\beta) = X^{\top} diag(\lambda_i)X,$$
  

$$W(\beta) = -I(\beta),$$

where 
$$\lambda = (\lambda_1, \dots, \lambda_n)^{\top}$$

When is the solution to the equation  $U(\beta) = 0$  maximizing the function  $\log L(Y, \beta)$ ?

In other words: When does an unique MLE  $\hat{\beta}$  of  $\beta$  exist?

$$\hat{\beta} = \operatorname*{argmax}_{\beta} \log L(Y, \beta).$$

The sufficient condition of a maximum implies that the Hesse matrix  $W(\beta)$  is negative definite.

Consider the case of the natural link function. Then Remark 5.3.3 implies that

- The system of equations  $U(\beta)=0$  can be rewritten as  $U(\beta)=\frac{1}{\tau^2}X^\top(Y-\mu(\beta))=0$
- The matrix  $W(\beta) = -\frac{1}{\tau^4} X^{\top} diag(\sigma_i^2(\beta)) X$  is negative definite, has  $\operatorname{rank}(X) = m$  and  $0 < \sigma_i^2(\beta) < \infty$  for all  $i = 1, \ldots, n$ . Under those conditions there exists an unique MLE  $\hat{\beta}$  for  $\beta$ .

In the following, two numerical algorithms for solving the system of (in general nonlinear) equations  $U(\beta) = 0$  are introduced. These approaches are iterative, i.e. they approximate the MLE  $\hat{\beta}$  incrementally.

#### 1. Newton's method

Choose a suitable starting value  $\hat{\beta}_0 \in \mathbb{R}^m$ .

In step k+1, calculate  $\hat{\beta}_{k+1}$  from  $\hat{\beta}_k$ ,  $k=0,1,\ldots$  as follows:

• Take the first order Taylor expansion of  $U(\beta)$  at  $\hat{\beta}_k : U(\beta) \approx U(\hat{\beta}_k) + W(\hat{\beta}_k)(\beta - \hat{\beta}_k)$ .

- Solving for 0:  $U(\hat{\beta}_k) + W(\hat{\beta}_k)(\beta \hat{\beta}_k) = 0$
- The solution for this system of equations is  $\hat{\beta}_{k+1}$ :

$$\hat{\beta}_{k+1} = \hat{\beta}_k - W^{-1}(\hat{\beta}_k) \cdot U(\hat{\beta}_k), \quad k = 0, 1, 2, \dots$$

assuming that  $W(\hat{\beta}_k)$  is invertible.

Stop the iteration process once  $|\hat{\beta}_{k+1} - \hat{\beta}_k| < \delta$  for a predetermined boundary  $\delta > 0$ .

The convergence of this method heavily depends on the choice of  $\hat{\beta}_0$ , since  $\hat{\beta}_0$  has to be close enough to  $\hat{\beta}$ . Another disadvantage of this method is that the random matrix  $W(\beta)$  might not be invertible. That is why a modification of the Newton method is presented in which  $W(\beta)$  is replaced by the expectation

$$\mathbb{E}W(\beta) = -I(\beta). \tag{5.7}$$

It can be shown that the identity (5.7) holds by using Theorem 5.3.4 and the fact that  $\mathbb{E}Y_i = \mu_i$ , i = 1, ..., n. If it is assumed that  $\operatorname{rank}(X) = m$  and  $u_i \neq 0$ , i = 1, ..., n, then by Theorem 5.3.2,  $I(\beta)$  is invertible. This method is called Fisher's scoring method.

The only difference of Newton's method compared to Fisher's Scoring is that in the second step the iterative equation

$$\hat{\beta}_{k+1} = \hat{\beta}_k + I^{-1}(\hat{\beta}_k)U(\hat{\beta}_k), k = 0, 1, \dots$$

is used.

In the case of the natural link function (cf. Remark 5.3.3)

$$\hat{\beta}_{k+1} = \hat{\beta}_k + \tau^4 \left( X^\top diag(\sigma_i^2(\hat{\beta}_k)) X \right)^{-1} \frac{1}{\tau^2} \left( X^\top (Y - \mu(\hat{\beta}_k)) \right)$$
$$= \hat{\beta}_k + \tau^2 \left( X^\top diag(\sigma_i^2(\hat{\beta}_k)) X \right)^{-1} \left( X^\top (Y - \mu(\hat{\beta}_k)) \right).$$

## 5.4 Asymptotic tests for $\beta$

The goal of this section is to construct a test for the hypotheses

$$H_0: \beta = \beta_0 \text{ vs.}$$
  
 $H_1: \beta \neq \beta_0$ 

with  $\beta = (\beta_1, \dots, \beta_m)^{\top}$  and  $\beta_0 = (\beta_{01}, \dots, \beta_{0m})^{\top}$ . In particular, the hypotheses  $H_0: \beta = 0$  resp.  $H_0: \beta_i = 0$  are of interest, because they imply

that the test variables  $Y = (Y_1, \dots, Y_n)^{\top}$  does not depend on several output variables (e.g.  $(x_{1j}, \dots, x_{nj})^{\top}$  in case of the hypothesis  $\beta_j = 0$ ).

In order to test these kinds of hypotheses, the test statistics  $T_n$  are used, which have an asymptotic (for  $n \to \infty$ ) reference distribution (e.g. multivariate normal or  $\chi^2$  distribution). Some groundwork has to be done beforehand though. Let

$$g(\mathbb{E}Y_i) = X_i\beta, \quad i = 1, \dots, n$$

be a generalized linear model with natural link function g. Let  $L(Y, \beta)$  be the likelihood function,  $U(\beta)$  the partial derivatives of  $\log L(Y, \beta)$  and  $I(\beta)$  the Fisher information matrix in this model.

 $\hat{\beta}_n = \hat{\beta}(Y_1, \dots, Y_n, X)$  denotes a sequence of maximum likelihood estimators for  $\beta$ .

Assume that

- 1.  $\exists$  compact subspace  $K \subset \mathbb{R}^m$ , such that all rows  $X_i$ , i = 1, ..., n,  $n \in \mathbb{N}$ , of X are in K. Here  $\theta = x^{\top}\beta \in \Theta$  for all  $\beta \in \mathbb{R}^m$  and  $x \in K$ .
- 2. There exists a sequence  $\{\Gamma_n\}_{n\in\mathbb{N}}$  of diagonal  $(m\times m)$  matrices  $\Gamma_n = \Gamma_n(\beta)$  with the properties
  - (a)  $\gamma_{i,i}^n > 0, i \in \{1, \dots, m\}$
  - (b)  $\lim_{n\to\infty} \Gamma_n = 0$ ,
  - (c)  $\lim_{n\to\infty} \Gamma_n^{\top} I_n(\beta) \Gamma_n = K^{-1}(\beta)$ , where  $K(\beta)$  is a symmetric positive definite  $(m \times m)$  matrix for all  $\beta \in \mathbb{R}^m$ .

**Theorem 5.4.1.** Under the conditions above, there exists a  $\Gamma_n$  consistent sequence of MLE  $\{\hat{\beta}_n\}$  for  $\beta$ ,

(i.e. 
$$P\left(\Gamma_n^{-1}|\hat{\beta}_n - \beta| \le \varepsilon, U(\hat{\beta}_n) = 0\right) \to 1 \text{ for } n \to \infty$$
), such that

1. 
$$T_n^* = \Gamma_n^{-1}(\hat{\beta}_n - \beta) \xrightarrow[n \to \infty]{d} \mathcal{N}(0, K(\beta)),$$

2. 
$$T_n = 2(\log L(Y, \hat{\beta}_n) - \log L(Y, \beta)) \xrightarrow[n \to \infty]{d} \chi_m^2$$

where  $m = \dim \beta$ .

**Remark 5.4.2.** (cf. [27], p.288-292)

- 1. Usually  $\Gamma_n = diag\left(\frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}}\right)$  is used.
- 2. Until now we always assumed that the dispersion term  $\tau^2$  is known. If that is not the case, then  $\tau^2$  can be estimated by

$$\hat{\tau}^2 = \frac{1}{n-m} \sum_{i=1}^n \frac{\left(Y_i - \mu_i(\hat{\beta}_n)\right)^2}{b''(\hat{\theta}_{ni})}$$

where  $\hat{\theta}_{ni} = (b')^{-1}(\mu_i(\hat{\beta}_n)), i = 1, \dots, n$ . This estimator is an empirical analogue to the equation  $\tau^2 = \frac{\text{Var}Y_i}{b''(\theta_i)}$  of Lemma 5.1.4.

3. The second assertion of Theorem 5.4.1 also holds, if the unknown parameter  $\tau^2$  can be replaced by a consistent estimator  $\tau_n^2$ .

How can Theorem 5.4.1 be used to test the hypothesis

$$H_0: \beta = \beta_0 \text{ vs.}$$
  
 $H_1: \beta \neq \beta_0$ 

or component wise

$$H_0: \beta_j = \beta_{j0}, \ j = 1, \dots, m \text{ vs.}$$
  
 $H_1: \exists j_1: \beta_{j_1} \neq \beta_{j_10}$ ?

Let

$$g(\mathbb{E}Y_i) = \sum_{j=1}^{m} x_{ij}\beta_j, \quad i = 1, \dots, n$$

be a generalized linear model with natural link function g. Using Remark 5.3.3, it holds that

$$\log L(Y, \beta) = \frac{1}{\tau^2} \sum_{i=1}^n \left( Y_i x_i^{\top} \beta + a(Y_i, \tau) - b(x_i^{\top} \beta) \right)$$

where  $Y = (Y_1, \dots, Y_n)^{\top}$  and  $x_i = (x_{i1}, \dots, x_{im})^{\top}$ . Thus is holds that

$$T_n = \frac{2}{\tau^2} \sum_{i=1}^n \left( Y_i x_i^{\top} (\hat{\beta}_n - \beta_0) - b(x_i^{\top} \hat{\beta}_n) + b(x_i^{\top} \beta_0) \right)$$

By specifying an exponential model  $(\tau, b \text{ are known})$ , with respect to the random sample of the goal variable Y and the design matrix X,  $H_0$  is rejected if  $T_n > \chi^2_{m,1-\alpha}$ , where m is the number of parameters in the model,  $\chi^2_{m,1-\alpha}$  the  $(1-\alpha)$  quantile of the  $\chi^2_m$ - distribution and  $\alpha \in (0,1)$  is the significance level of the asymptotic test. This test can only be applied for relatively large n. Type I errors have the (for  $n \to \infty$ ) asymptotic probability  $\alpha$ . If a simple hypothesis

$$H_0: \beta_j = 0 \text{ vs.}$$
  
 $H_1: \beta_j \neq 0$ 

is tested, the test statistic  $T_n^1$  is used which can be derived from  $T_n^*$ :  $H_0$  is rejected, if

$$|T_n^1| = \frac{|\hat{\beta}_{nj}|}{(\Gamma_n(\hat{\beta}_n))_{jj}} > z_{1-\frac{\alpha}{2}}$$

where  $z_{1-\frac{\alpha}{2}}$  is the  $(1-\frac{\alpha}{2})$  quantile of the N(0,1) distribution. Here  $\{\Gamma_n\}$  is chosen in a way that  $K(\beta) = Id$ , for all  $\beta \in \mathbb{R}^m$ . This is an asymptotic test with confidence level  $\alpha$ , since

with confidence level 
$$\alpha$$
, since 
$$P_{H_0}(|T_n^1| > z_{1-\frac{\alpha}{2}}) = 1 - P_{H_0}(|T_n^1| \le z_{1-\frac{\alpha}{2}}) \xrightarrow[n \to \infty]{} 1 - \Phi(z_{1-\frac{\alpha}{2}}) + \underbrace{\Phi(-z_{1-\frac{\alpha}{2}})}_{1-\Phi(z_{1-\frac{\alpha}{2}})}$$

$$=1-\left(1-\frac{\alpha}{2}\right)+1-\left(1-\frac{\alpha}{2}\right)=\alpha,$$

where

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-\frac{t^2}{2}} dt$$

is the cumulative distribution function of the N(0,1) distribution.

#### Example 5.4.3. (Credit risk assessment)<sup>1</sup>

The following data is provided by a southern German bank from the 90's: Results from credit risk assessment for n = 1000 credit applications (ca. 700 good credits and 300 bad credits) analysed:

Goal variable  $Y_i = \begin{cases} 0, & \text{if the credit of customer } i \text{ has been paid} \\ 1, & \text{if the credit of customer } i \text{ has not been paid} \end{cases}$ 

The design matrix X contains the following additional information about the customer:

$$x_{i1}$$
 - Account management with the bank = 
$$\begin{cases} 1, & \text{no account} \\ 0, & \text{else} \end{cases}$$

$$x_{i2}$$
 - Assessment of account management = 
$$\begin{cases} 1, & \text{good account} \\ 0, & \text{no- or bad account} \end{cases}$$

 $x_{i3}$  - Term of credit in months

 $x_{i4}$  - Value of Credit in DM

$$x_{i5}$$
 - Payment history of customer = 
$$\begin{cases} 1, & \text{good} \\ 0, & \text{else} \end{cases}$$

$$x_{i6}$$
 - Reference = 
$$\begin{cases} 1, & \text{private} \\ 0, & \text{business} \end{cases}$$

**Question**: How should  $\hat{\beta}$  be estimated?

As a model, the logit model is used with  $p_i = P(Y_i = 1), i = 1, ..., n$ :

$$\log \frac{p_i}{1 - p_i} = \beta_0 + x_{i1}\beta_1 + x_{i2}\beta_2 + x_{i3}\beta_3 + x_{i4}\beta_4 + x_{i5}\beta_5 + x_{i6}\beta_6$$

for 
$$i = 1, ..., n$$
, where  $\beta = (\beta_0, ..., \beta_6)^{\top}, m = 7$ .

<sup>&</sup>lt;sup>1</sup>cf. Fahrmeir, L., Kneib, T., Lang, S. - Regression, p.208

		Y = 1	Y = 0
$x_1$	no account	45.0	20.0
$x_2$	$\operatorname{good}$	15.3	49.8
	bad	39.7	30.2
$x_4$	Credit value	Y = 1	Y = 0
	$0 < \ldots \le 500$	1.00	2.14
	$500 < \ldots \le 1000$	11.33	9.14
	$1000 < \ldots \le 1500$	17.00	19.86
	$1500 < \ldots \le 2500$	19.67	24.57
	$2500 < \ldots \le 5000$	25.00	28.57
	$5000 < \ldots \le 7500$	11.33	9.71
	$7500 < \ldots \le 10000$	6.67	3.71
	$10000 < \ldots \le 15000$	7.00	2.00
	$15000 < \ldots \le 20000$	1.00	0.29
$x_5$	Credit history	Y = 1	Y = 0
	good	82.33	94.95
	bad	17.66	5.15
$x_6$	Reference	Y = 1	Y = 0
	private	57.53	69.29
	business	42.47	30.71
		-	-

Table 5.1: Abstract of the data

**Goal**: Estimate  $\beta_0, \ldots, \beta_6$  and check, which factors are important for future credit risk assessment.

 $H_0: \beta_i = 0$  (feature  $x_i$  does not affect the credit risk assessment) is rejected, if the p-value  $\leq \alpha$ . It can also be noticed that  $\beta_4$  is not relevant for credit risk assessment, contrary to belief. A refinement of the model is necessary:

#### New model:

$$g(\mathbb{E}Y_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3^1 x_{i3} + \beta_3^2 x_{i3}^2 + \beta_4^1 x_{i4} + \beta_4^2 x_{i4}^2 + \beta_5 x_{i5} + \beta_6 x_{i6}$$

Question: Which model is better?

$\overline{x}_1$	$\overline{x}_2$	$\overline{x}_3$	$\overline{x}_4$	$\overline{x}_5$	$\overline{x}_6$
0.274	0.393	20.903	3271	0.911	0.657

Table 5.2: Means  $\overline{x}_j$  of  $x_{ij}$  in the data

	Value	$\sqrt{(I_n^{-1}(\hat{\beta}))_{ii}}$	$T_n^1$	<i>p</i> -value
$\beta_0$	0.281	0.303	-0.94	0.347
$\beta_1$	0.618	0.175	3.53	< 0.001
$\beta_2$	-1.338	0.201	-6.65	< 0.001
$\beta_3$	0.033	0.008	4.29	< 0.001
$\beta_4$	0.023	0.033	0.72	0.474
$\beta_5$	-0.986	0.251	-3.93	< 0.001
$\beta_6$	-0.426	0.266	-2.69	0.007

Table 5.3: Results of the ML estimation by using the Fisher Scoring method, where  $\sqrt{(I_n^{-1}(\hat{\beta}))_{ii}}$  is used as an asymptotic standard deviation of  $\hat{\beta}_i$ . Significance level :  $\alpha = 0.001$ .

In other words, the following hypotheses are tested:

 $H_0: \beta_3^2 = 0$  (linear model) vs.  $H_1: \beta_3^2 \neq 0$  (quadratic model) resp.

 $H_0: \beta_4^2 = 0$  (linear model) vs.  $H_1: \beta_4^2 \neq 0$  (quadratic model)

Here the type of statistical hypothesis is generalized as follows:

$$H_0: C\beta = d \text{ vs. } H_1: C\beta \neq d$$

is tested, where C is a  $(r \times m)$  - matrix with rank  $C = r \leq m$  and  $d \in \mathbb{R}^r$ . For comparison, the hypothesis

$$H_0: \beta = \beta_0 \text{ vs. } H_1: \beta \neq \beta_0, \quad \beta, \beta_0 \in \mathbb{R}^m$$

was tested before. Obviously  $\beta = \beta_0$  is a special case of  $C\beta = d$  with C =Id,  $d = \beta_0$ . The new hypotheses include assertions about the linear combinations of the parameters. How should  $H_0$  vs.  $H_1$  be tested?

Let  $\widetilde{\beta}_n$  be the MLE of  $\beta$  under  $H_0$ , i.e.  $\widetilde{\beta}_n = \underset{\beta \in \mathbb{R}^m: \ C\beta = d}{\operatorname{argmax}} \log L(Y, \beta)$ Let  $\widehat{\beta}_n$  be the MLE of  $\beta$  unrestricted, i.e.  $\widehat{\beta}_n = \underset{\beta \in \mathbb{R}^m}{\operatorname{argmax}} \log L(Y, \beta)$ .

The idea behind the following tests is to compare  $\tilde{\beta}_n$  with  $\hat{\beta}_n$ .

	Value	$\sqrt{(I_n^{-1}(\hat{\beta}))_{ii}}$	$T_n^1$	p-Wert
$\beta_0$	-0.488	0.390	-1.25	0.211
$\beta_1$	0.618	0.176	3.51	< 0.001
$\beta_2$	-1.337	0.202	-6.61	< 0.001
$\beta_3^1$	0.092	0.025	3.64	< 0.001
$\beta_3^2$	-0.001	< 0.001	-2.20	0.028
$eta_4^1$	-0.264	0.099	-2.68	0.007
$\beta_4^1$	0.023	0.007	3.07	0.002
$\beta_5$	-0.995	0.255	-3.90	< 0.001
$\beta_6$	-0.404	0.160	-2.52	0.012

Table 5.4: p-values for the regression coefficients of the new model

If the deviation  $\hat{\beta}_n - \tilde{\beta}_n$  is big,  $H_0$  should be rejected.

**Theorem 5.4.4.** Let  $\log L(Y,\beta)$  be the log-likelihood function of the random sample of the goal variable  $Y = (Y_1, \ldots, Y_n)^{\top}$ ,  $I_n(\beta)$  be the fisher information matrix,  $U(\beta)$  be the score function of the generalized linear model with natural link function g:

$$g(\mathbb{E}Y_i) = X_i\beta, \quad i = 1, \dots, n.$$

Consider the following test statistics

1. likelihood-ratio test statistic:

$$\widetilde{T}_n = 2(\log L(Y, \hat{\beta}_n) - \log L(Y, \widetilde{\beta}_n))$$

2. Wald statistic:

$$\widetilde{T}_n^* = (C\widehat{\beta}_n - d)^\top (CI_n^{-1}(\widehat{\beta}_n)C^\top)^{-1}(C\widehat{\beta}_n - d)$$

3. Score statistic:

$$\overline{T}_n^* = U(\widetilde{\beta}_n)^{\top} I_n^{-1}(\widetilde{\beta}_n) U(\widetilde{\beta}_n)$$

Under certain conditions for the estimators  $\hat{\beta}$  and  $\widetilde{\beta}$  (cf. Theorem 5.4.1) the test statistics 1 - 3 are asymptotically  $\chi_m^2$  distributed, e.g. for the likelihood-ratio-test statistic it holds that

$$\widetilde{T}_n \xrightarrow[n \to \infty]{d} \chi_m^2$$
.

Corollary 5.4.5. Theorem 5.4.4 provides the following decision rule:  $H_0$  is rejected, if

 $\widetilde{T}_n(\widetilde{T}_n^*, \overline{T}_n) > \chi_{m,1-\alpha}^2.$ 

This is an asymptotic test with confidence level  $\alpha$ .

**Example 5.4.6** (Continued). The following values for the test statistics are obtained:

- $\tilde{T}_n = 12.44 \ p$ -value: 0.0020
- $\widetilde{T}_n^* = 11.47 \ p$ -value: 0.0032

for  $\alpha = 0.005$  it holds that the *p*-value  $\leq \alpha$ , thus  $H_0: \beta_4^2 = 0$  is rejected  $\Rightarrow$  the quadratic generalized linear model is preferred.

## 5.5 Criteria for model selection or model adjustment

It is known that the goodness of fit of a parametric model to the data rises, if the number of parameters increases. A goal of a statistician is to find a well fitted model with as little variables as possible in order to avoid overfitting. The Akaike information criterion can help achieving this goal by comparing models with (possibly) different parameter estimators.

Akaike information coefficient:

$$AIC = -2\log L(Y, \hat{\beta}) + 2m$$

where  $Y = (Y_1, ..., Y_n)$  is the random sample of the goal variable in the generalized linear model and  $\hat{\beta}$  the corresponding MLE. The value of the AIC takes the required maximality of the log-likelihood function  $\log L(Y, \hat{\beta})$  into account and punishes models with a large number of parameters m. The models with the smallest AIC is considered the better model. Sometimes instead of the AIC, the standardized AIC given by  $\frac{\text{AIC}}{n}$  is used.

**Example 5.5.1** (Continued). Calculate the AIC for the linear and quadratic Logit model in the example of credit risk assessment:

 $\begin{aligned} & \text{Linear model}: \text{AIC} = 1043.815 \\ & \text{Quadratic model}: \text{AIC} = 1035.371 \end{aligned}$ 

By considering the AIC, it can be noticed that the quadratic model seems to be better.

A disadvantage of making a decision based on the AIC alone is, that the final decision is up to the statistician. Thus it is desirable to construct a

statistical test which can asses the goodness of fit for the model. The  $\chi^2$ -test aims to solve this issue.

Let

$$g(\mathbb{E}Y_i) = X_i\beta, \quad i = 1, \dots, n$$

be a generalized linear model with link function g and parameter vector  $\beta = (\beta_1, \dots, \beta_m)^{\top}$ . Split the goal variables  $Y_1, \dots, Y_n$  in k groups, such that they are as homogeneous as possible with respect to the parameters that need to be estimated. A said partition can be achieved by splitting the domain of the goal variable  $Y_i$  "skillfully" in  $k > m^2$  intervals  $(a_l, b_l]$ :

$$-\infty \le a_1 < b_1 = a_2 < b_2 = a_3 < \dots < b_{k-1} = a_k < b_k \le +\infty$$

Group l contains all observations  $Y_i$ , which are in  $(a_l, b_l]$ . Here  $(a_l, b_l]$  need to be chosen in a way, such that  $\hat{\mu}_j = g^{-1}(X_j\hat{\beta})$  are constant within the respective groups:  $\hat{\mu}_j \equiv \hat{\mu}_l \ \forall \ j$  of group l.<sup>3</sup> Let

- $n_l = \# \{Y_j : Y_j \in (a_l, b_l]\}$  be class size of class l.
- $\overline{Y}_l = \frac{1}{n_l} \sum Y_j$  be the arithmetic mean of class l.
- $\hat{\beta}$  be the MLE of  $\beta$ , which was obtained by using Y.
- $l_l(\beta) = \sum \log f_{\theta}(Y_j)$  be the log-likelihood function of the goal variables  $Y_i$  within the group l.
- $\hat{\mu}_l = g^{-1}(X_l\hat{\beta})$  and  $v(\hat{\mu}_l)$  be the expectation estimator resp. variance estimator  $\mu_l = \mathbb{E}Y_l$ , which are obtained by using the MLE  $\hat{\beta}$ .

Here  $v(\hat{\mu}_l) = \tau^2 b''(b'^{-1}(\hat{\mu}_l))$ , where  $b(\cdot)$  is the corresponding coefficient in the probability density function  $f_{\theta}$  in the exponential family. The following test statistic is obtained:

$$\chi^2 = \sum_{l=1}^k \frac{(\overline{Y}_l - \hat{\mu}_l)^2}{v(\hat{\mu}_l)/n_l},$$

$$D = -2\tau^2 \sum_{l=1}^k \left( l_l(\hat{\mu}_l) - l_l(\overline{Y}_l) \right).$$

**Theorem 5.5.2.** If  $n \to \infty$  and  $n_l \to \infty \ \forall \ l$ , then under certain conditions it holds that

$$\chi^2 \xrightarrow[n \to \infty]{d} \chi^2_{k-m-1}$$

$$D \xrightarrow[n \to \infty]{d} \chi^2_{k-m-1}$$

$$^2k \le m \Rightarrow D \xrightarrow[n \to \infty]{d} \chi^2_{k-m-1}$$

$$^3\text{This is an informal description of the methodology}$$

<sup>3</sup>This is an informal description of the methodology, in which for each  $Y_i$ ,  $n_i$  independent copies of  $Y_i$  are created, which compose the *i*-th class.

Corollary 5.5.3. The hypotheses

$$H_0: Y = (Y_1, \dots, Y_n)$$
 is from the model  $g(\mathbb{E}Y_i) = X_i$ 

vs

$$H_1: Y = (Y_1, \ldots, Y_n)$$
 is not from the model  $g(\mathbb{E}Y_i) = X_i\beta$ 

for i = 1, ..., n can be tested as follows:

 $H_0$  is rejected (for large n) at significance level  $\alpha$ , if

$$\chi^2 > \chi^2_{k-m-1,1-\alpha}$$
 resp.  $D > \chi^2_{k-m-1,1-\alpha}$ .

Those tests should not be used if the class sizes  $n_l$  are small.

**Example 5.5.4.** What do the tests described above look like in the Logit-resp. Poisson regression?

1. logit model:  $y_i \sim \text{Bernoulli}(p_i), i = 1, ..., n$ 

$$\Rightarrow$$
 generalized linear model  $\log \frac{p_i}{1-p_i} = x_i \beta$ 

for  $i=1,\ldots,n$ . Divide  $y_1,\ldots,y_n$  in k classes, such that the probability of occurring 1 in each class is estimated as good as possible by  $\overline{y}_l = \frac{1}{n_l} \sum y_i$ . Thus it holds that

- $\hat{\mu}_l = \hat{p}_l = g^{-1}(X_l \hat{\beta}) = \frac{e^{X_l^{\top} \hat{\beta}}}{1 + e^{X_l^{\top} \hat{\beta}}},$
- $v(\hat{p}_l) = \hat{p}_l(1 \hat{p}_l),$
- $\chi^2 = \sum_{l=1}^k \frac{(\overline{Y}_l \hat{p}_l)^2}{\hat{p}_l(1 \hat{p}_l)/n_l}$ .
- 2. Poisson model:  $Y_i \sim \text{Poisson}(\lambda)$ ,

$$\Rightarrow$$
 generalized linear model  $\log \lambda_i = X_i \beta$ 

for i = 1, ..., n. Thus it holds that  $\hat{\mu}_l = \hat{\lambda}_l = e^{X_l \hat{\beta}}, \ v(\hat{\lambda}_l) = \hat{\lambda}_l$  and

$$\chi^2 = \sum_{l=1}^k \frac{(\overline{Y}_l - \hat{\lambda}_l)^2}{\hat{\lambda}_l / n_l}.$$

# Chapter 6

# Principal Component Analysis

In this chapter, methods for reducing the complexity of big statistical data is presented in form of the principal component analysis (PCA). PCA aims to reduce a high dimensional datasets  $X = (X_1, \ldots, X_n)^{\top} \in \mathbb{R}^n$  to very few but important components  $\varphi = AX \in \mathbb{R}^d$  with  $d \ll n$ . Those components should then explain most of the variability of the original dataset X. Here, A is a  $(d \times n)$  matrix which can be found if some restrictions (given in (6.1)) are fulfilled. Another applications of PCA is the visualization of complex datasets, outlier detection, cluster analysis and so on. For an overview see [18].

#### 6.1 Introduction

In order to motivate the following problem, consider the example of text mining in automotive:

**Example 6.1.1.** A car manufacturer is interested in minimizing its losses which are due to fraud and incompetence in warranty repairs in one of the subsidiaries. That's why he wants to introduce a conspicuousness analysis of repair visits in said warranty subsidiary, which is supposed to find suspicious reports with the help of a computer that can be manually checked afterwards. Another motivation for the automatized early detection system is the comprehensive examination of few subsidiaries in irregular time intervals (due to high costs) which could be marginalized. A typical repair log consists of a maximum of 300.000 technical terms. That's why the logs should be saved as vectors  $x = (x_1, \ldots, x_n)^{\top}$  of length n = 300.000, where

$$x_i = \begin{cases} 1, & \text{the word } i \text{ is in text } x \\ 0, & \text{else} \end{cases}$$

These vectors x are normed such that they are on the sphere  $S^{n-1}$ . Within one year a huge dataset of those vectors is created with several million entries. The task of a statistician is the drastic reduction of the dimension n-1 of the data such that a visualization is possible. A possible solution is the usage of PCA. The groundwork for PCA has been done by Beltran (1873) and Jordan (1874) who introduced the singular value decomposition (SVD). In a more or less modern form (cf. (6.1)) it is presented in the work of K. Pearson (1901) and H. Hotelling (1933). Also, the name PCA was introduced by Hotelling. A more developed version of the method was introduced by Girshick (1939), Anderson (1963), Rao (1964) and some others. Without a computer the calculation of principal components for n > 4 turns out to be rather difficult, thus this methodology has found its practical applications after their invention.

Since the 1980's there is a rapid increase in applications of PCA in the whole knowledge domain (especially in in engineering), where multivariate datasets are analyzed.

#### 6.2 PCA on model level

This section aims to introduce the main concept of PCA for random samples  $X = (X_1, \ldots, X_n)^{\top}$  with known covariance structure. Let  $X = (X_1, \ldots, X_n)^{\top}$  be a random sample of random variables  $X_i$  with known covariance matrix  $\Sigma$  and  $\operatorname{Var} X_i \in (0, \infty)$ ,  $i = 1, \ldots, n$ . Let  $\lambda_1 > \lambda_2 > \ldots > \lambda_n > 0$  be the eigenvalues of  $\Sigma$ , which are sorted in descending order and all different from each other. The goal is to find linear combinations  $\alpha^{\top} X$  of  $X_i$  which have the biggest variance, whereas the vectors  $\alpha$  are normed respectively, e.g. such that  $\alpha \in S^{n-1}$  with the euclidean norm.

**Definition 6.2.1.** The linear combination  $\alpha_i^\top X$ , i = 1, ..., n, is called *i*-th principal component of X, if it has the biggest variance under the condition that  $\alpha_i \in S^{n-1}$  and  $\alpha_1^\top X, \alpha_2^\top X, ..., \alpha_{i-1}^\top X$  and  $\alpha_i^\top X$  are uncorrelated:

$$\begin{cases} \operatorname{Var} \alpha^{\top} X \to \max_{\alpha}, \\ |\alpha| = 1, \\ \operatorname{Cov} (\alpha^{\top} X, \alpha_{j}^{\top} X) = 0, \quad j = 1, \dots, i - 1. \end{cases}$$

$$(6.1)$$

Here  $\alpha_i$  is the coefficient vector of the *i*-th principal component  $\alpha_i^{\top} X$ .

**Theorem 6.2.2.** The i-th principal component of X is given by

$$Y_i = \alpha_i^\top X,$$

where  $\alpha_i$  is the eigenvector of  $\Sigma$  with eigenvalue  $\lambda_i$ . Here

$$Var(Y_i) = \lambda_i, \quad i = 1, \dots, n.$$

**Proof** Show that the assertion holds for i = 1, 2. For i > 2 the proof works in the same spirit.

For i=1 there is a constraint  $|\alpha|=1$  in (6.1), which is taken into the Lagrange function

$$f(\alpha) = \operatorname{Var}(\alpha^{\top} X) + \lambda(|\alpha|^2 - 1).$$

Furthermore

$$\operatorname{Var}(\alpha^{\top}X) = \mathbb{E}(\alpha^{\top}X - \mathbb{E}\alpha^{\top}X)^{2} = \mathbb{E}(\alpha^{\top}(X - \mathbb{E}X))^{2}$$
$$= \mathbb{E}\alpha^{\top}(X - \mathbb{E}X)(X - \mathbb{E}X)^{\top}\alpha = \alpha^{\top}\mathbb{E}(X - \mathbb{E}X)(X - \mathbb{E}X)^{\top}\alpha$$
$$= \alpha^{\top}\Sigma\alpha,$$

 $|\alpha|^2 = \alpha^\top \cdot \alpha$ , and  $f(\alpha) = \alpha^\top \Sigma \alpha + \lambda (\alpha^\top \alpha - 1)$ .

The necessary conditions for a maximum is given by

$$\frac{\partial f}{\partial \alpha} = 0, \quad \frac{\partial f}{\partial \lambda} = 0,$$

where the second equation represents the constraint  $|\alpha| = 1$ .  $\frac{\partial f}{\partial \alpha} = \left(\frac{\partial f}{\partial \alpha^1}, \dots, \frac{\partial f}{\partial \alpha^n}\right)$ , where  $\alpha = (\alpha^1, \dots, \alpha^n)^{\top}$  and  $\frac{\partial f}{\partial \alpha} = 0$  can be rewritten as  $\Sigma \alpha - \lambda \alpha = 0$  vectorial or  $\Sigma \alpha = \lambda \alpha$ , which means, that  $\alpha$  is an eigenvector of  $\Sigma$  with eigenvalue  $\lambda$ . Since  $\operatorname{Var}(\alpha^{\top}X) = \alpha^{\top}\Sigma\alpha$  is supposed to be maximized, it holds that

$$\operatorname{Var}(\alpha^{\top} X) = \alpha^{\top} \lambda \alpha = \lambda \underbrace{\alpha^{\top} \alpha}_{1} = \lambda$$

and  $\lambda = \lambda_1 > \lambda_2 > \ldots > \lambda_n \Rightarrow \lambda = \lambda_1$  and  $\alpha = \alpha_1$ . For i = 2, the maximization task

$$\begin{cases} \alpha^{\top} \Sigma \alpha \to \max_{\alpha}, \\ \alpha^{\top} \cdot \alpha = 1, \\ \text{Cov } (\alpha_1^{\top} X, \alpha^{\top} X) = 0 \end{cases}$$

needs to be solved with respect to  $\alpha$ , where

Cov 
$$(\alpha_1 X, \alpha^\top X) = \alpha_1^\top \Sigma \alpha = \alpha^\top \Sigma \alpha_1 = \alpha^\top \lambda_1 \alpha_1 = \lambda_1 \alpha^\top \alpha_1.$$

That means, the following function needs to be maximized:

$$f(\alpha) = \alpha^{\top} \Sigma \alpha + \lambda (\alpha^{\top} \alpha - 1) + \delta \alpha^{\top} \alpha_1.$$

Similarly as above it holds that

$$\frac{\partial f}{\partial \alpha} = \Sigma \alpha + \lambda \alpha + \delta \alpha_1 = 0.$$

The constraint  $\alpha_1^{\top} \Sigma \alpha = 0$  and  $\alpha_1^{\top} \alpha = 0$  (see above) imply

$$\alpha_1^{\top} \frac{\partial f}{\partial \alpha} = \delta \underbrace{\alpha_1^{\top} \alpha_1}_{1} = \delta = 0,$$

which means, that  $\Sigma \alpha = \lambda \alpha$  and  $\alpha$  is again, an eigenvector of  $\Sigma$  with eigenvalue  $\lambda$ . Since  $\alpha$  is supposed to be orthogonal to  $\alpha_1$  and  $\operatorname{Var}(\alpha^{\top} X) = \lambda$  is supposed to be maximized, it holds that

$$\alpha = \alpha_2 \text{ and } \lambda = \lambda_2 \Rightarrow Y_2 = \alpha_2^\top X.$$

**Exercise 6.2.3.** Work out the proof for i > 2!

Let now  $A = (\alpha_1, \dots, \alpha_n)$ . This is a orthogonal  $(n \times n)$  matrix, for which it holds that (by Theorem 6.2.2)

$$\Sigma A = A\Lambda, \quad \Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_n),$$

or equivalently

$$A^{\top} \Sigma A = \Lambda, \quad \Sigma = A \Lambda A^{\top}. \tag{6.2}$$

**Theorem 6.2.4.** For a  $(n \times m)$  matrix B, with orthogonal columns  $b_i$ , i = 1, ..., m,  $m \le n$ , let  $Y = B^{\top}X$  and  $\Sigma_Y = \text{Cov }(Y) = B^{\top}\Sigma B$  be the covariance matrix of Y. Then

$$A_m = \operatorname*{argmax}_{B} \operatorname{trace}(\Sigma_Y),$$

where  $A_m = (\alpha_1, \ldots, \alpha_m)$ .

**Proof** Since  $\alpha_1, \ldots, \alpha_n$  is a basis of  $\mathbb{R}^n$ , it holds that

$$b_k = \sum_{i=1}^n c_{ik} \alpha_i, \quad k = 1, \dots, m,$$

where  $B = (b_1, \ldots, b_m)$ , or matrix wise, B = AC, with  $C = (c_{ij})$ ,  $i = 1, \ldots, n, j = 1, \ldots, m$ . Thus it holds that

$$\Sigma_Y = B^{\top} \Sigma B = C^{\top} \underbrace{A^{\top} \Sigma A}_{\Lambda} C = C^{\top} \Lambda C = \sum_{j=1}^n \lambda_j c_j c_j^{\top},$$

where  $c_j^{\top}$  is the *j*-th row of C. Thus it holds that

$$\operatorname{trace}(\Sigma_Y) = \sum_{j=1}^n \lambda_j \operatorname{trace}(c_j c_j^\top) = \sum_{j=1}^n \lambda_j \operatorname{trace}(c_j^\top c_j) = \sum_{j=1}^n \lambda_j |c_j|^2.$$

Since  $C = A^{-1}B = A^{\top}B$ , it holds that

$$C^{\top}C = B^{\top} \underbrace{AA^{\top}}_{I_n} B = \underbrace{B^{\top}B}_{I_m} = I_m,$$

where

$$I_k = \operatorname{diag}\left(\underbrace{1,\ldots,1}_{k}\right).$$

Thus

$$\sum_{i=1}^{n} \sum_{j=1}^{m} c_{ij}^2 = m,$$

and the columns of C are orthonormal. Thus C can be seen as a part (the first m columns) of an orthonormal  $(n \times n)$  matrix D. Since the rows of D are also orthonormal vectors and  $c_i^{\top}$  are the first m entries of the rows of D, it holds that

$$c_i^{\top} c_i = \sum_{i=1}^m c_{ij}^2 \le 1, \quad i = 1, \dots, n.$$

Since

$$\operatorname{trace}(\Sigma_Y) = \sum_{i=1}^n \lambda_i \sum_{j=1}^m c_{ij}^2 = \sum_{i=1}^n \beta_i \lambda_i,$$

where  $\beta_i \leq 1$ , i = 1, ..., n,  $\sum_{i=1}^n \beta_i = m$  and

$$\lambda_1 > \lambda_2 > \dots > \lambda_n, \quad \sum_{i=1}^n \beta_i \lambda_i \to \max$$

for  $\beta_1 = \ldots = \beta_m = 1$ ,  $\beta_{m+1} = \ldots = \beta_n = 0$ . If  $B = A_m$ , then

$$c_{ij} = \begin{cases} 1 & , 1 \le i = j \le m \\ 0 & , \text{ else} \end{cases},$$

which implies  $\beta_1 = \ldots = \beta_m = 1$ ,  $\beta_{m+1} = \ldots = \beta_n = 0$ . Thus  $A_m$  is the solution of  $\operatorname{trace}(\Sigma_Y) \to \max_B$ .

The assertion of Theorem 6.2.4 implies that

$$\operatorname{Var}\left(\sum_{i=1}^{m} Y_i\right) = \operatorname{Var}\left(\sum_{i=1}^{m} \alpha_i^{\top} X\right)$$

is maximized for all m = 1, ..., n, if  $Y_i$  are principal components of X.

Corollary 6.2.5. Spectral representation of  $\Sigma$ .

It holds that

$$\Sigma = \sum_{i=1}^{n} \lambda_i \cdot \alpha_i \cdot \alpha_i^{\top}. \tag{6.3}$$

**Proof** The representation is obtained by using (6.2), since

$$\Sigma = (\alpha_1, \dots, \alpha_n) \cdot \operatorname{diag}(\lambda_1, \dots, \lambda_n) \cdot (\alpha_1, \dots, \alpha_n)^{\top}.$$

#### Remark 6.2.6.

- 1. Since  $\lambda_1 > \lambda_2 > \ldots > \lambda_n$  with  $|\alpha_i| = 1$ ,  $\forall i$ , the representation (6.3) implies, that the first principal components do not only explain the biggest ratio of the variance of  $X_i$ , but also to the covariance. This value decreases with an increasing  $i = 1, \ldots, n$ .
- 2. If rank  $(\Sigma) = r < n$ , then (6.3) implies, that  $\Sigma$  can be completely determined by considering the first r principal components and coefficient vectors.

**Lemma 6.2.7.** Let  $\Sigma$  be a positive definite and symmetric  $(n \times n)$  matrix with eigenvalues  $\lambda_1 > \lambda_2 > \ldots > \lambda_n > 0$  and corresponding eigenvectors  $\alpha_1, \ldots, \alpha_n, |\alpha_i| = 1, i = 1, \ldots, n$ . Then

$$\lambda_k = \sup_{\alpha \in S_k, \alpha \neq 0} \frac{\alpha^\top \Sigma \alpha}{|\alpha|^2},$$

where  $S_k = \langle \alpha_1, \dots, \alpha_{k-1} \rangle^{\perp}$  for arbitrary  $k = 1, \dots, n$ .

**Proof** Let

$$c = \sup_{\alpha \in S_k} \frac{\alpha^{\top} \Sigma \alpha}{|\alpha|^2}.$$

Show that  $\lambda_k \leq c \leq \lambda_k$ .

1.  $c \geq \lambda_k$ : For  $\alpha = \alpha_k$  prove that

$$c \geq \frac{\alpha_k^\top \Sigma \alpha_k}{\alpha_k^\top \alpha_k} = \frac{\lambda_k \alpha_k^\top \alpha_k}{\alpha_k^\top \alpha_k} = \lambda_k.$$

2.  $c \leq \lambda_k$ : It needs to be shown that

$$\alpha^{\top} \Sigma \alpha \leq \lambda_k |\alpha|^2, \quad \forall \alpha \in S_k, \quad \alpha \neq 0, \quad \forall \alpha \in \mathbb{R}^n \quad \alpha = \sum_{i=1}^n c_i \alpha_i,$$

 $<sup>^1\</sup>langle\dots\rangle$  denotes the span

since  $\{\alpha_i\}_{i=1}^n$  form an orthonormal basis.

$$\alpha \in S_k \quad \Rightarrow \quad c_1 = \ldots = c_{k-1} = 0.$$

That means

$$\alpha = \sum_{i=k}^{n} c_i \alpha_i, \quad \Sigma \alpha = \sum_{i=1}^{n} c_i \Sigma \alpha_i = \sum_{i=1}^{n} c_i \lambda_i \alpha_i,$$

$$\alpha^{\top} \Sigma \alpha = \left(\sum_{i=1}^{n} c_i \alpha_i\right)^{\top} \left(\sum_{i=1}^{n} \lambda_i c_i \alpha_i\right)$$

$$= \sum_{i,j=1}^{n} c_i c_j \lambda_i \underbrace{\alpha_j^{\top} \alpha_i}_{\delta_{i,i}} = \sum_{i=1}^{n} c_i^2 \lambda_i, \quad |\alpha|^2 = \sum_{i=1}^{n} c_i^2.$$

Thus it holds that  $\alpha \in S_k$ 

$$\alpha^{\top} \Sigma \alpha = \sum_{i=k}^{n} c_i^2 \lambda_i \le \sum_{i=k}^{n} \lambda_k c_i^2 = \lambda_k \sum_{i=k}^{n} c_i^2 = \lambda_k |\alpha|^2,$$

and  $c \leq \lambda_k$  since  $\lambda_k > \lambda_j$ , j > k.

**Theorem 6.2.8.** Let B, Y and  $\Sigma_Y$  such as in Theorem 6.2.4. Then

$$A_m = \operatorname*{argmax}_{B} \det(\Sigma_Y),$$

where  $A_m = (\alpha_1, \ldots, \alpha_m)$ 

**Proof** Let  $k \in \{1, ..., m\}$  be fix. Introduce  $S_k = \langle \alpha_1, ..., \alpha_{k-1} \rangle^{\perp} \subset \mathbb{R}^k$  (as in Lemma 6.2.7). Let  $\mu_1 > \mu_2 > ... > \mu_m$  be the eigenvalues of  $\Sigma_Y = B^{\top} \Sigma B$  with corresponding eigenvectors  $\gamma_1, ..., \gamma_m$ , which are orthonormal. Let  $T_k = \langle \gamma_{k+1}, ..., \gamma_m \rangle \subset \mathbb{R}^m$ . It obviously holds, that

$$Dim (S_k) = n - k + 1, \quad Dim T_k = k.$$

As in Lemma 6.2.7, it can be shown, that  $\forall \gamma \neq 0, \gamma \in T_k$  it holds that

$$\frac{\gamma^{\top} \Sigma \gamma}{|\gamma|^2} \ge \mu_k.$$

Consider  $\widetilde{S}_k = B(T_k) \subset \mathbb{R}^n$ . Since B is an orthonormal transformation, it is thus unique Dim  $(\widetilde{S}_k) = \text{Dim } (T_k) = k$ . The formula

$$\operatorname{Dim} (S_k \cup \widetilde{S}_k) + \operatorname{Dim} (S_k \cap \widetilde{S}_k) = \operatorname{Dim} S_k + \operatorname{Dim} \widetilde{S}_k$$

implies

$$\operatorname{Dim} (S_k \cap \widetilde{S}_k) = \underbrace{\operatorname{Dim} S_k}_{n-k+1} + \underbrace{\operatorname{Dim} \widetilde{S}_k}_{k} - \underbrace{\operatorname{Dim} (S_k \cup \widetilde{S}_k)}_{\leq n} \geq n-k+1+k-n = 1$$

that means,  $\exists \alpha \in S_k \cap \widetilde{S}_k$ ,  $\alpha \neq 0$ . For this  $\alpha$  it holds that  $\alpha = B\gamma$ ,  $\gamma \in T_k$  and thus

$$\mu_k \le \frac{\gamma^\top \Sigma \gamma^2}{|\gamma|^2} = \frac{\gamma^\top B^\top \Sigma B \gamma}{\underbrace{\gamma^\top \gamma}_{\gamma^\top B^\top B \gamma}} = \frac{\alpha^\top \Sigma \alpha}{\alpha^\top \alpha} \le \lambda_k$$

since  $|\gamma| = |B\gamma|$ , because B is preserving distances. That's why  $\mu_k \leq \lambda_k$  for all  $k = 1, \ldots, m$  and

$$\det(\Sigma_Y) = \prod_{i=1}^m \mu_k \le \prod_{k=1}^m \lambda_k \quad \Rightarrow \quad \max_B \det(\Sigma_Y) \le \prod_{k=1}^m \lambda_k.$$

However, since  $B = A_m$ ,  $\mu_k = \lambda_k$ , k = 1, ..., m, it holds that

$$A_m = \operatorname*{argmax}_{B} \det(\Sigma_Y).$$

Now geometric properties of principal components are considered.

**Proposition 6.2.9.** The principal component coefficients  $\alpha_1, \ldots, \alpha_n$  are the principle axis of the ellipsoids  $x^{\top} \Sigma^{-1} x = c$ , with semi-axis length  $\sqrt{c\lambda_i}$ ,  $i = 1, \ldots, n$ .

**Proof** The principal components of X are given by  $Z = A^{\top}X$ , where  $A = (\alpha_1, \ldots, \alpha_n)$  is an orthonormal transformation and thus  $A^{\top} = A^{-1}$ , X = AZ. Therefore for the ellipsoid it holds that

$$x^{\top} \Sigma^{-1} x = \sum_{Subst. x = Az} z^{\top} A^{\top} \Sigma^{-1} A z = z^{\top} \Lambda^{-1} z = c,$$

where

$$A^{\top} \Sigma^{-1} A = \Lambda^{-1} = \operatorname{diag} \left( \frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_n} \right), \quad \Lambda = \operatorname{diag} \left( \lambda_1, \dots, \lambda_n \right),$$

since  $\Sigma^{-1}$  has the same eigenvectors with eigenvalues  $\frac{1}{\lambda_i}$ . That's why the ellipsoid  $z^{\top}\Lambda^{-1}z=c$  can be represented in its normed form as

$$\sum_{k=1}^{n} \frac{z_k^2}{c\lambda_k} = 1.$$

That implies that the  $\alpha_i$  point towards the principal axis and that the half-axis are given by  $\sqrt{c\lambda_i}$ .

Remark 6.2.10. (Multivariate normal distribution). If  $X \sim N(0, \Sigma)$  then  $x^{\top}\Sigma^{-1}x = c$  is an ellipsoid of constant probability for X, since the probability density function of X

$$f_X(x) = \frac{1}{\sqrt{\det \Sigma}} \exp\left\{-\frac{1}{2}x^{\top} \Sigma^{-1} x\right\} \cdot \frac{1}{(2\pi)^{\frac{n}{2}}}, \quad x \in \mathbb{R}^n,$$

is constant on this ellipsoid. Else  $x^{\top}\Sigma^{-1}x = c$  defines contours of the constant probability for X. Here the vector  $\alpha_1$  points towards the largest variance of  $\alpha^{\top}X$  (it is the biggest principal axis with length  $\sqrt{c\lambda_1}$  of the ellipsoid);  $\alpha_2$  points towards the second largest variance (half-axis with length  $\sqrt{c\lambda_2}$ ), and so on (cf. Condition (6.1)).

**Remark 6.2.11.** Another form of PCA is possible, if instead of  $X = (X_1, \ldots, X_n)^{\top}$  the normed random sample  $X_{\omega} = (X_1/\omega_1, \ldots, X_n/\omega_n)^{\top}$  is used, where the weights  $\omega = (\omega_1, \ldots, \omega_n)^{\top}$  contain advance information which represent a certain preference in the analysis. A usual choice is

$$\omega_i = \sqrt{\sigma_{ii}} = \sqrt{\operatorname{Var} X_i},$$

which leads to a PCA of  $X^* = (X_1^*, \dots, X_n^*)$ ,  $X_i^* = \frac{X_i}{\sqrt{\text{Var}X_i}}$ ,  $i = 1, \dots, n$  by using the correlation matrix  $\Sigma^* = (\text{Corr }(X_j, X_i))_{i,j=1}$  with

$$\operatorname{Corr}(X_i, X_j) = \frac{\operatorname{Cov}(X_i, X_j)}{\sqrt{\operatorname{Var} X_i \operatorname{Var} X_j}} = \operatorname{Cov}(X_i^*, X_j^*), \quad i, j = 1, \dots, n.$$

By doing so, other principal components  $\alpha_i^{*T}X^*$  can be obtained for which  $\alpha_i^* \neq \alpha_i$  holds for i = 1, ..., n.

What are the advantages and disadvantages of PCA based on  $(X, \Sigma)$  and  $(X^*, \Sigma^*)$ ?

#### Disadvantages of $(X, \Sigma)$ -PCA:

- 1. PCA based on  $(X^*, \Sigma^*)$  does not depend on the unit measurements of X. Thus comparisons between results of PCA for several samples of different origin are possible.
- 2. If the variances  $X_i$  are varying a lot, the variables  $X_i$  with the largest variance are determined to be the first principal component. The PCA based on  $(X^*, \Sigma^*)$  does not have this disadvantage. The  $(X, \Sigma)$ -PCA is not significant in such a case, because it sorts the variables  $X_i$  (in a slightly different form) in a variance wise descending order.

**Example 6.2.12.** Let  $X = (X_1, X_2)$ , where  $X_1$  is the length and  $X_2$  the weight.  $X_1$  can be measured in cm or m,  $X_2$  only in kg. In those two cases, the covariance matrices X are given by

$$\Sigma_1 = \begin{pmatrix} 80 & 44 \\ 44 & 80 \end{pmatrix}$$
 resp.  $\Sigma_2 = \begin{pmatrix} 8000 & 4400 \\ 4400 & 8800 \end{pmatrix}$ .

Calculating the first principal components in both cases yields

$$\alpha_1^{\top} X = 0,707X_1 + 0,707X_2 \text{ for } \Sigma_1 \text{ resp.}$$
  
 $\alpha_1^{\top} X = 0,998X_1 + 0,055X_2 \text{ for } \Sigma_2.$ 

Note that, in the first case, both  $X_1$  and  $X_2$  have the same absolute value with respect to the 1. principal component, whereas in the 2. case  $X_1$  is a dominating factor. Moreover it holds that  $\frac{\lambda_1}{\lambda_1+\lambda_2}\cdot 100\% = 77,5\%$  in the first case and  $\frac{\lambda_1}{\lambda_1+\lambda_2}\cdot 100\% = 99,3\%$  in the 2. case (it is the ratio of the variation of the first principal component to the whole variation).

3. If random variables  $X_i$  in X have a differing origin (as in the example above), then the interpretation of the ratio of the variation is rather problematic, since the sum  $\lambda_1 + \ldots + \lambda_n$  contains  $m^2$ ,  $kg^2$  and so on. The PCA based on  $(X^*, \Sigma^*)$  only considers values without unit, such that the sum  $\lambda_1 + \ldots + \lambda_n$  can be interpreted.

#### Advantages of $(X, \Sigma)$ -PCA:

- 1. If instead of  $\Sigma$  resp.  $\Sigma^*$  the empirical analogues  $\hat{\Sigma}$  resp.  $\hat{\Sigma}^*$  are used (if  $\Sigma(\Sigma^*)$  are not known, they have to be estimated by using the available data), then  $(X, \hat{\Sigma})$ -PCA has some advantages, since the statistical methods are easier to use in this case compared to using them in  $(X^*, \hat{\Sigma}^*)$ -PCA.
- 2. If all  $X_i$  in X have the same unit, then the  $(X, \Sigma)$ -PCA is sometimes preferable, since during the standardisation of  $(X, \Sigma)$  to  $(X^*, \Sigma^*)$  the relation to the units of X are lost.

**Remark 6.2.13.** Sometimes instead of  $|\alpha| = 1$  the standardization  $|\alpha_k| = \sqrt{\lambda_k}$ , k = 1, ..., n in Definition 6.2.1 is used (cf. optimisation problem (6.1)). This is in particular the case in the correlation based PCA.

Remark 6.2.14. (Equal eigenvalues  $\lambda_i$ ). If some eigenvalues of  $\Sigma$  are equal, e.g.  $\lambda_1 = \lambda_2 = \ldots = \lambda_k > \lambda_{k+1} > \ldots > \lambda_m$  implies that there exists a linear subspace of dimension k, in which an arbitrary basis represents the first k eigenvectors. This means, with respect to PCA that the first k eigenvectors can not be defined uniquely. Geometrical interpretation: The first k half-axis of  $x^{\top}\Sigma^{-1}x = c$  are equal, i.e., the ellipsoid  $x^{\top}\Sigma^{-1}x = c$  has a spherical k-dimensional cross section through the origin, where the directions of the half-axis can be chosen (orthogonal to each other) arbitrarily.

**Remark 6.2.15** ( $\lambda_i = 0$ ). If  $\lambda_1 > \ldots > \lambda_{n-k} > \lambda_{n-k+1} = \ldots = \lambda_n = 0$ , then there are only n-k linear independent random vectors  $X_i$  in the random sample X. That's why only those n-k variables should be good for the analysis.

#### 6.3 PCA on data level

In this section it is not assumed that the covariance matrix  $\Sigma$  is known. That's why it is replaced by empirical covariance matrix  $\hat{\Sigma}$ .

Let  $X^1, X^2, \ldots, X^m$  be independent realizations of a n-dimensional random vector  $X = (X_1, \ldots, X_n)^{\top}, X^i = (X_1^i, \ldots, X_n^i)^{\top}, i = 1, \ldots, m.$   $X^i$  is interpreted as a sample of X.

**Definition 6.3.1.** Define the *n*-dimensional random vector  $a_k$  by

$$a_k = \operatorname*{argmax}_{a \in \mathbb{R}^n} \frac{1}{m-1} \sum_{i=1}^m (Y_i - \overline{Y})^2$$

with constraint |a| = 1, a uncorrelated to  $a_1, \ldots, a_{k-1}$  for all  $k = 1, \ldots, n$ , where

$$Y_i = a^{\mathsf{T}} X^i, \quad i = 1, \dots, m, \quad \overline{Y} = \frac{1}{m} \sum_{i=1}^m Y_i.$$

Thus  $a_k^{\top}X$  defines the k-th principal component of X with coefficient vector  $a_k$ .  $Y_{ik} = a_k^{\top}X^i$  is the evaluation of the k-ten principal component of the i-th observation  $X^i$  of  $X_i$ ,  $i = 1, \ldots, m$ ,  $k = 1, \ldots, n$ .

Lemma 6.3.2. It holds that

$$\frac{1}{m-1} \sum_{i=1}^{m} (Y_{ik} - \overline{Y}_k)^2 = l_k, \quad k = 1, \dots, n,$$

where

$$\overline{Y}_k = \frac{1}{m} \sum_{i=1}^m Y_{ik}, \ \overline{X}_k = \frac{1}{m} \sum_{i=1}^m X_k^i, \quad k = 1, \dots, n$$

and  $l_k$  is the eigenvalue of the empirical covariance matrix  $\hat{\Sigma} = (\hat{\sigma}_{ij})_{i,j=1}^n$ ,

$$\hat{\sigma}_{ij} = \frac{1}{m-1} \sum_{t=1}^{m} (X_i^t - \overline{X}_i)(X_j^t - \overline{X}_j), \quad i, j = 1, \dots, n, \quad l_1 > l_2 > \dots > l_n.$$

 $a_k$  is the eigenvector of  $\hat{\Sigma}$  with eigenvalue  $l_k$ ,  $k=1,\ldots,n$ .

#### Proof

Exercise 6.3.3. cf. proof of Theorem 6.2.2.

In the following, replace  $X^i$  with  $X^i - \overline{X}$  but keep the notation  $X^i$ ,  $i = 1, \ldots, n$ .

**Remark 6.3.4.** The properties of PCA as formulated in Theorem 6.2.4, Corollary 6.2.5 and Proposition 6.2.9 are preserved in the statistical version (cf. Definition 6.3.1) by using the following modification:  $\Sigma$  is replaces by  $\hat{\Sigma}$ ,  $A = (\alpha_1, \ldots, \alpha_n)$  by  $A = (a_1, \ldots, a_n)$ ,  $A_m = (\alpha_1, \ldots, \alpha_m)$  by  $A_m = (a_1, \ldots, a_m)$  and  $\Sigma_Y$  by the empirical covariance matrix  $\hat{\Sigma}_Y$  of Y. Thus use the spectral representation of  $\hat{\Sigma}$ :

$$\hat{\Sigma} = \sum_{i=1}^{n} l_i a_i a_i^{\top}. \tag{6.4}$$

#### Exercise 6.3.5. Prove that!

In the following another property of the empirical PCA, which can also be seen as an equivalent definition is presented:

**Theorem 6.3.6.** Let B be a  $(n \times p)$  matrix,  $p \leq n$ , with orthogonal columns. Let  $Z_i = B^{\top} X^i$ , i = 1, ..., m be a projection of  $X^i$ , i = 1, ..., m, on a p-dimensional subspace  $L_B$ . Define

$$G(B) = \sum_{i=1}^{m} \left| X^i - Z_i \right|^2.$$

Then

$$A_p = (a_1, \dots, a_p) = \operatorname*{argmin}_B G(B).$$

**Proof** By the Pythagorean Theorem it holds that  $|X^i|^2 = |Z_i|^2 + |X^i - Z_i|^2$ , that's why

$$G(B) = \sum_{i=1}^{m} |X^i|^2 - \sum_{i=1}^{m} |Z_i|^2 \to \min$$

if

$$\widetilde{G}(B) = \sum_{i=1}^{m} |Z_i|^2 = \sum_{i=1}^{m} Z_i^{\top} Z_i = \sum_{i=1}^{m} X^{iT} B B^{\top} X^i \to \max_{B}.$$

It holds that

$$\widetilde{G}(B) = \operatorname{trace}\left(\sum_{i=1}^{m} \left(X^{iT}BB^{\top}X^{i}\right)\right) = \sum_{i=1}^{m} \operatorname{trace}\left(X^{iT}BB^{\top}X^{i}\right)$$

$$= \sum_{i=1}^{m} \operatorname{trace}\left(B^{\top}X^{i}X^{iT}B\right) = \operatorname{trace}\left(B^{\top}\underbrace{\left(\sum_{i=1}^{m} X^{i}X^{iT}\right)}_{2(m-1)\hat{\Sigma}}B\right)$$

$$= (m-1)\operatorname{trace}(B^{\top}\hat{\Sigma}B).$$

In summary:

$$\widetilde{G}(B) = (m-1)\operatorname{trace}\left(B^{\top}\widehat{\Sigma}B\right),$$

since it is maximized by Remark 6.3.4 and Theorem 6.2.4, if  $B = A_p$ .

**Remark 6.3.7.** How can Theorem 6.3.6 be used as an equivalent definition of the empirical PCA?  $a_i$  are defined as orthogonal vectors, which is the span of a linear subspace  $L_p = \langle a_1, \ldots, a_p \rangle$ ,  $p = 1, \ldots, n-1$ , with the property, that the sum of the quadratic orthogonal distances of  $X^i$  to  $L_p$  are minimized. Thus for p = 1,  $L_1$  would be the best line for approximating the data set  $X^1, \ldots, X^m$ . For p = n-1,  $L_{n-1}$  would be the best hyperplane with the same property (cf. linear regression).

The following theorem provides a new interpretation for PCA as well as more efficient method for the calculation.

**Theorem 6.3.8.** (Singular value decomposition.)

Let  $\widetilde{X} = \left(X^1 - \overline{X}, X^2 - \overline{X}, \dots, X^m - \overline{X}\right)^{\top}$  be a  $(m \times n)$  matrix, with centered observations  $X^i$  of X. Let rank  $(\widetilde{X}) = r \leq n, m$ . The following decomposition holds:

$$\widetilde{X} = ULA_r^{\top}, \tag{6.5}$$

where U is a  $(m \times r)$  matrix with orthonormal columns.

$$L = \operatorname{diag}(\widetilde{l}_1, \dots, \widetilde{l}_r)$$
 where  $\widetilde{l}_i = \sqrt{(m-1)l_i}$ 

is the square root of the *i*-th (non trivial) eigenvalue of  $\widetilde{X}^{\top}\widetilde{X} = (m-1)\hat{\Sigma}$ , i = 1, ..., r.  $A_r = (a_1, ..., a_r)$  is the  $(n \times r)$  matrix with columns  $a_i$ .

**Proof** Define  $U = (u_1, \ldots, u_r)$  with columns  $u_i = Xa_i/l_i$ ,  $i = 1, \ldots, r$ . In the following it is shown, that the representation (6.5) holds. Using the spectral representation 6.4 it holds that

$$(m-1)\hat{\Sigma} = \tilde{X}^{\top}\tilde{X} = \sum_{i=1}^{r} \tilde{l}_{i}^{2} a_{i} a_{i}^{\top}, \text{ since } l_{i} = 0, i = r+1, \dots, n.$$

Thus

$$ULA_r^{\top} = U \begin{pmatrix} \widetilde{l}_1 a_1^{\top} \\ \vdots \\ \widetilde{l}_r a_r^{\top} \end{pmatrix} = \sum_{i=1}^r \widetilde{X} \frac{a_i}{\widetilde{l}_i} \widetilde{l}_i a_i^{\top} = \sum_{i=1}^r \widetilde{X} a_i a_i^{\top} \stackrel{l_i = 0, i > r}{=} \sum_{i=1}^n \widetilde{X} a_i a_i^{\top}.$$

<sup>&</sup>lt;sup>2</sup>Since  $X^i$  was replaced by  $X^i - \overline{X}$ .

It holds that  $\widetilde{X}a_i = 0$ ,  $i = r + 1, \ldots, n$ , since  $\operatorname{rank}(\widetilde{X}) = r$  and because of the centered columns of  $\widetilde{X}$  by  $\overline{X}$ . By the orthogonality of the  $a_i$  it holds that

$$ULA_r^{\top} = \widetilde{X} \sum_{i=1}^n a_i a_i^{\top} = \widetilde{X}I = \widetilde{X}.$$

**Remark 6.3.9.** The matrix U provides the following versions of evaluations

$$Y_{ik} = a_k^{\top} X^i = X^{iT} a_k, \quad Y_{ik} = u_{ik} \tilde{l_k}, \quad i = 1, \dots, m, \ k = 1, \dots, n.$$

It holds that

$$Var(u_{ik}) = \frac{Var(Y_{ik})}{\tilde{l}_k^2} = \frac{l_k}{(m-1)l_k} = \frac{1}{m-1}, \quad \forall i, k.$$

# 6.4 Asymptotic distributions of principal components for normal distributed random samples

Let now  $X \sim N(\mu, \Sigma)$ ,  $\Sigma$  have the eigenvalues  $\lambda_1 > \lambda_2 > \ldots > \lambda_n > 0$  and corresponding eigenvectors  $\alpha_k$ ,  $k = 1, \ldots, n$ . Calculate

$$\lambda = (\lambda_1, \dots, \lambda_n)^\top, \quad l = (l_1, \dots, l_n)^\top,$$
  

$$\alpha_k = (\alpha_{k1}, \dots, \alpha_{kn})^\top, \quad a_k = (a_{k1}, \dots, a_{kn})^\top,$$
  

$$k = 1, \dots, n$$

#### Theorem 6.4.1.

- 1. l is asymptotically (for  $m \to \infty$ ) independent of  $a_k$ ,  $k = 1, \ldots, n$ .
- 2. l and  $a_k$ ,  $k=1,\ldots,n$  are asymptotic  $m\to\infty$  multivariate normal distributed, with asymptotic expectation

$$\lim_{m \to \infty} \mathbb{E}(l) = \lambda \quad \text{and} \quad \lim_{m \to \infty} \mathbb{E}(a_k) = \alpha_k, \quad k = 1, \dots, n.$$

3. It holds that

Cov 
$$(l_k, l_{k'}) \sim \begin{cases} \frac{2\lambda_k^2}{m-1}, & k = k' \\ 0, & k \neq k' \end{cases}$$
 for  $m \to \infty$ ,

$$\operatorname{Cov}\left(a_{kj}, a_{k'j'}\right) \sim \begin{cases} \frac{\lambda_k}{m-1} \sum_{l=1, l \neq k}^{n} \frac{\lambda_l \alpha_{lj} \alpha_{lj'}}{(\lambda_l - l_k)^2}, & k = k' \\ -\frac{\lambda_k \lambda_{k'} \alpha_{kj} \alpha_{k'j'}}{(m-1)(\lambda_k - \lambda_{k'})^2}, & k \neq k' \end{cases}$$
 for  $m \to \infty$ .

### Without proof!

The assertion of Theorem 6.4.1 can be used for constructing MLE and confidence intervals for  $\lambda$  and  $\alpha_k$ .

#### Exercise 6.4.2.

- 1. Show that an MLE of  $\Sigma$  is given by  $\frac{m-1}{m}\hat{\Sigma}$ .
- 2. Show, that the MLE

$$\begin{cases} \text{ for } \lambda \text{ is given by} & \hat{\lambda} = \frac{m-1}{m}l. \\ \text{ for } \alpha_k \text{ is given by} & \hat{\alpha}_k = a_k, k = 1, \dots, n. \end{cases}$$

3. Show that the MLE in 2. coincide with the moment estimators  $\lambda$  and  $\alpha_k$ , which can be obtained from Theorem 6.4.1.

**Corollary 6.4.3** (Confidence intervals for  $\lambda_k$ ). An asymptotic confidence interval for  $\lambda_k$   $(m \to \infty)$  with confidence level  $1 - \alpha$  is given by

$$\left[l_k\left(1-\sqrt{\frac{2}{m-1}}z_{\frac{\alpha}{2}}\right)^{-1},l_k\left(1+\sqrt{\frac{2}{m-1}}z_{\frac{\alpha}{2}}\right)^{-1}\right],$$

where m is large enough such that  $-\sqrt{\frac{2}{m-1}}z_{\frac{\alpha}{2}} < 1$ .

**Proof** Since  $l_k \sim N\left(\lambda_k, \frac{2\lambda_k^2}{m-1}\right)$  for  $m \to \infty$  by Theorem 6.4.1, 2. and 3., it holds that

$$\frac{l_k - \lambda_k}{\sqrt{\frac{2}{m-1}} \lambda_k} \sim N(0, 1) \quad \text{for} \quad m \to \infty.$$

This implies, that

$$\lim_{m \to \infty} P\left(z_{\frac{\alpha}{2}} \le \frac{l_k - \lambda_k}{\lambda_k} \sqrt{\frac{m-1}{2}} \le z_{1-\frac{\alpha}{2}}\right) = 1 - \alpha,$$

or for  $m \to \infty$ 

$$\sqrt{\frac{2}{m-1}}z_{\frac{\alpha}{2}} \leq \frac{l_k}{\lambda_k} - 1 \leq \sqrt{\frac{2}{m-1}}\underbrace{z_{1-\frac{\alpha}{2}}}_{=-z_{\frac{\alpha}{2}}},$$

$$\frac{l_k}{1 - \sqrt{\frac{2}{m-1}} z_{\frac{\alpha}{2}}} \le \lambda_k \le \frac{l_k}{1 + \sqrt{\frac{2}{m-1}} z_{\frac{\alpha}{2}}}$$

with probability  $1 - \alpha$ .

Since all  $l_k$ , k = 1, ..., n are asymptotically  $(m \to \infty)$  independent, a simultaneous confidence interval for l can be denoted by a cartesian product of the confidence intervals for  $l_k$  as in Corollary 6.4.3.

Lemma 6.4.4. It holds that

$$(m-1)\alpha_k^{\top} \left(l_k \hat{\Sigma}^{-1} + l_k^{-1} \hat{\Sigma} - 2I_n\right) \alpha_k \xrightarrow[m \to \infty]{d} \chi_{n-1}^2.$$

#### Without proof!

As a consequence of the lemma above, the (asymptotic) confidence ellipsoid for  $\alpha_k$  with confidence level  $1-\beta$ 

$$\left\{ y \in \mathbb{R}^n : (m-1)y^{\top} \left( l_k \hat{\Sigma}^{-1} + l_k^{-1} \hat{\Sigma} - 2I_n \right) y \le \chi_{n-1,\beta}^2 \right\}$$

is obtained.

**Remark 6.4.5.** Corollary 6.4.3 resp. Lemma 6.4.4 can be used to construct statistical tests for  $\lambda_k$  resp.  $\alpha_k$  as follows:

1. Test  $H_0: \lambda_k = \lambda_{k_o}$  v.s.  $H_1: \lambda_k \neq \lambda_{k_0}$   $H_0$  can be rejected, if

$$\left| \frac{l_k - \lambda_{k_0}}{\sqrt{\frac{2}{m-1}} \lambda_{k_0}} > z_{\frac{\alpha}{2}} \right|.$$

This is an asymptotic test  $(m \to \infty)$  with confidence level  $\alpha$ .

2. Test  $H_0: \alpha_k = \alpha_{k_0}$  v.s.  $H_1: \alpha_k \neq \alpha_{k_0}$   $H_0$  can be rejected, if

$$(m-1)\alpha_{k_0}^{\top} \left( l_k \hat{\Sigma}^{-1} + l_k^{-1} \hat{\Sigma} - 2I_n \right) \alpha_{k_0} \ge \chi_{n-1,\alpha}^2.$$

This is an asymptotic test  $(m \to \infty)$  with confidence level  $\alpha$ .

#### 6.5 Outlier detection

In this section it is assumed that the random sample  $X^1, X^2, \ldots, X^m$  can contain some outliers. How can an outlier be defined? In statistical literature there is no coherent definition. Generally speaking, an observation  $X^i$  is an outlier if it attains an unusual value (with respect to the distribution of X). For example, an unusual value of some coordinates  $X_i$  could be significantly bigger or smaller than the others. An outlier could also occur in form of an unusual combination of the coordinate values of some coordinates  $X^i$ . A reason for those anomalies could lie in the data, or simply occur because of measurement errors.

**Example 6.5.1.** Let  $X = (X_1, X_2)$ , where  $X_1 =$  "height" (in cm) and  $X_2 =$  "weight" (in kg) of children between the age of 5 and 15. The feature  $X_i$  is obtained  $X_i$  in a medical survey. Here the features  $X_i$  = (250, 80) and  $X_i$  = (175, 25) are considered as outliers, because  $X_i$  = 250cm is an abnormal height and for  $X_i$ ,  $X_i$  = 175 and  $X_i$  = 25 as a combination are highly unlikely.

How can outliers be detected? One way to identify outliers of  $X^i$  is to plot the dataset  $X^1, \ldots, X^m$  and spot values which are outside of a larger agglomeration of values. If the dimension n of X is high, it is rather difficult to visualize the data. It can thus be helpful to generate a data point of the first 2-3 principal components of  $(X^1, \ldots, X^m)$ . By looking at them, outliers of  $X^i_k$  can also be identified quickly. In order to detect unusual relationships between coordinate values  $X^i_k$ , the last few principal components should be considered. Let  $a_1, \ldots, a_n$  be the coefficient vectors of the principal components of  $(X^1, \ldots, X^m)$ ,  $Y_{ik} = a_k^\top X^i$ ,  $i = 1, \ldots, m$ ,  $k = 1, \ldots, n$  be realizations of the principal components of the observation  $X^i$  and  $l_k$ ,  $k = 1, \ldots, n$  be the eigenvalues of the empirical covariance matrix  $\hat{\Sigma}$  of  $(X^1, \ldots, X^m)$ . For  $1 \leq n_0 \leq n$ , define the statistic

$$d_i^{(1)}(n_0) = \sum_{k=n-n_0+1}^n Y_{ik}^2, \quad d_i^{(2)}(n_0) = \sum_{k=n-n_0+1}^n \frac{Y_{ik}^2}{l_k},$$

$$d_i^{(3)}(n_0) = \sum_{k=n-n_0+1}^n l_k Y_{ik}^2, \quad d_i^{(4)}(n_0) = \max_{n-n_0+1 \le k \le n} \frac{|Y_{ik}|}{\sqrt{l_k}},$$

for  $i = 1, \ldots, m$ .

Lemma 6.5.2. It holds that

$$d_j^{(2)}(n) = \left(X^i - \overline{X}\right)^{\top} \hat{\Sigma}^{-1} \left(X^i - \overline{X}\right), \quad i = 1, \dots, m,$$

where  $Y_{ik}$  are centered, i.e.  $Y_{ik}$  is replaced by  $Y_{ik} - \overline{Y_k}$ , k = 1, ..., n, i = 1, ..., m.

**Proof** It holds that

$$\hat{\Sigma} = ALA^{\top}$$
, where  $L = \text{diag } (l_1, \dots, l_n)$  and  $A = (a_1, \dots, a_n)$ .

Thus

$$\hat{\Sigma}^{-1} = AL^{-1}A^{\top}$$
 with  $L^{-1} = \text{diag } (l_1^{-1}, \dots, l_n^{-1}).$ 

Since additionally  $Y_i = A^{\top} X^i$  for  $Y_i = (Y_{i1}, \dots, Y_{in})^{\top}$ ,  $i = 1, \dots, n$ , it holds that

$$X^{i} = A^{\top^{-1}}Y_{i} = AY_{i}, \quad X^{i}^{\top} = Y_{i}^{\top}A^{\top}, \quad i = 1, \dots, n$$

and thus

$$\overline{X} = \frac{1}{m} \sum_{i=1}^{m} X^i = A \overline{Y}, \quad \overline{Y} = \frac{1}{m} \sum_{i=1}^{m} Y^i, \quad \overline{X}^\top = \overline{Y}^\top A^\top.$$

This implies that

$$\begin{split} \left(X^{i} - \overline{X}\right)^{\top} \hat{\Sigma}^{-1} \left(X^{i} - \overline{X}\right) &= \left(Y_{i} - \overline{Y}\right)^{\top} \underbrace{A^{\top} A}_{I} L^{-1} \underbrace{A^{\top} A}_{I} \left(Y_{i} - \overline{Y}\right) \\ &= \left(Y_{i} - \overline{Y}\right)^{\top} L^{-1} \left(Y_{i} - \overline{Y}\right) = \sum_{k=1}^{n} \frac{Y_{ik}^{2}}{l_{k}} = d_{i}^{(2)}(n). \end{split}$$

In order to identify outliers in  $(X^1, \ldots, X^m)$ , the values  $d_i^{(j)}(n)$ ,  $i = 1, \ldots, m$ ,  $j = 1, \ldots, n$  for n = 1, 2, 3 are calculated. Observations  $X^i$  with the largest value  $d_i^{(j)}(n)$  are classified as possible outliers. Additionally the plot of the point cloud, defined by

$$D = \left\{ \left( d_i^{(2)}(n) - d_i^{(2)}(n_0), d_i^{(2)}(n_0) \right), i = 1, \dots, m \right\}$$

can be helpful.  $X^i$  is considered an outlier, if

$$\left(d_i^{(2)}(n) - d_i^{(2)}(n_0), d_i^{(2)}(n_o)\right)$$

is isolated from the remaining point cloud D.

Remark 6.5.3. If  $X \sim N(\mu, \Sigma)$  with known  $\mu$  and  $\Sigma$  and PCA is conducted on model level, the distributions of  $d_i^{(j)}(n_0)$  can be explicitly stated. They are (except for  $d_i^{(4)}$ ) gamma distributed with known parameters e.g.  $d_i^{(2)}(n_0) \sim \chi_{n_0}^2$ ,  $i=1,\ldots,m$ . The distribution function of  $d_j^{(4)}(n_0)$  is given by  $\Phi^{n_0}(x)$ , where  $\Phi(x)$  is the distribution function of a N(0,1) distribution. Confidence intervals for  $d_i^{(j)}(n_0)$  can provide a decision rule, whether  $X^i$  is an outlier. Even though this approach is based on a strict mathematical basis, it is rather uncommon in practice, since normally distributed data (with known  $\mu$  and  $\Sigma$ !) are relatively rare.

**Remark 6.5.4.** The statistics  $d_i^{(2)}, d_i^{(4)}$  emphasize the last statistics more than  $d_i^{(1)}$  (because of the corresponding standardization). That's why they are sufficient for the detection of unusual correlations in the data (cf. Example 6.5.1, observation  $X^j = (175, 25)$ . The statistic  $d_j^{(3)}$  emphasizes the first principal component. Thus it can be used to detect unusual large (small) values of the coordinates  $X_k^i$  (cf. Example 6.5.1  $X_1^i = 250$ ).

П

### 6.6 PCA and regression

Consider the multivariate regression model:  $Y = X\beta + \varepsilon$ , where  $Y = (Y_1, \ldots, Y_n)^{\top}$  is the vector of goal variables,

$$X = (X_{ij})_{\substack{i=1,\dots,n\\j=1,\dots,m}}$$

the  $(n \times m)$  matrix of output variables, rank (X) = m,  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^{\top}$  the vector of error terms, where  $\varepsilon_i$  are independent of  $\mathbb{E}\varepsilon_i = 0$ ,  $\operatorname{Var}\varepsilon_i = \sigma^2$ ,  $i = 1, \dots, n$ . W.l.o.g. assume, that X (as in Theorem 6.3.8) is centered, i.e., the empirical mean of X is zero, or in detail,  $X_{ij}$  is replaced by  $X_{ij} - \overline{X_j}$ , where

$$\overline{X_j} = \frac{1}{n} \sum_{i=1}^{n} X_{ij}, \quad j = 1, \dots, m.$$

Assuming that some of the variables  $X_{ij}$  in X are almost linearly dependent, i.e.  $\det(X^{\top}X) \approx 0$ , causes the estimators  $\hat{\beta}$  of  $\beta$  to be affected in form of an instability for the calculation, since  $\operatorname{Cov}(\hat{\beta}) = \sigma^2(X^{\top}X)^{-1}$  (cf. Theorem 6.3.8) only contains little variance of  $\hat{\beta}_j$ . A solution to this problem is the usage of generalizations as in Section 6.3. Another application of PCA is the detection of linear dependencies in X by looking at the last principal components and eliminating variables  $\beta_j$  based on those. This application will be discussed in more detail below.

Let  $a_1, \ldots, a_m$  be the coefficient vectors of the principal components (i.e. the eigenvectors) of  $X^{\top}X$ . Let  $Z_{ik} = a_k^{\top}X^i$  be the realization of the k-th principal component of the i-th row  $X^i$  of X,  $i = 1, \ldots, n$ ,  $k = 1, \ldots, m$ . With  $Z = (Z_{ik})$  it holds that Z = XA, where  $A = (a_1, \ldots, a_m)$  is an orthogonal  $(m \times m)$  matrix. The regression equation  $Y = X\beta + \mathcal{E}$  is given by:

$$Y = X \underbrace{AA^{\top}}_{I} \beta + \mathcal{E} = \underbrace{XA}_{Z} \underbrace{A^{\top}\beta}_{\gamma} + \mathcal{E} = Z\gamma + \mathcal{E}, \text{ where } \gamma = A^{\top}\beta.$$
 (6.6)

By doing so, the old output variables  $\beta$  are replaced by the transformation  $\gamma = A^{\top}\beta$ . The estimation of  $\gamma$  is obtained with Theorem 4.2.1:

$$\hat{\gamma} = \left(Z^{\top}Z\right)^{-1}Z^{\top}Y = L^{-1}Z^{\top}Y,\tag{6.7}$$

where  $L = \text{diag } (l_1, \ldots, l_m)$  contains the eigenvalues  $l_i$  of  $X^{\top}X$ . This holds, since Z has orthogonal columns. Thus

$$\hat{\beta} = A\hat{\gamma} = AL^{-1}Z^{\top}Y = \underbrace{AL^{-1}A^{\top}}_{(X^{\top}X)^{-1}}X^{\top}Y = \sum_{k=1}^{m} l_{k}^{-1}a_{k}a_{k}^{\top}X^{\top}Y,$$

where in the last part of the equation the terms (6.6), (6.7) and the spectral representation (Corollary 6.2.5) of  $(X^{\top}X)^{-1}$  have been used. Theorem 6.2.4 implies furthermore, that

$$\operatorname{Var}(\hat{\beta}) = \sigma^2 \sum_{k=1}^m l_k^{-1} a_k a_k^{\top}.$$

Thus the following assertion is proved:

**Lemma 6.6.1.** The solution of the OLS equation  $Y = X\beta + \mathcal{E}$  is given by

$$\hat{\beta} = \sum_{k=1}^{m} l_k^{-1} a_k a_k^{\top} X^{\top} Y.$$

Here it holds that

$$\operatorname{Cov}(\hat{\beta}) = \sigma^2 \sum_{k=1}^{m} l_k^{-1} a_k a_k^{\top}.$$

**Remark 6.6.2.** What are the advantages of the in (6.6)-(6.7) introduced methodology?

- 1. After calculating the principal components of  $X^{\top}X$ , the calculation of  $\hat{\gamma} = L^{-1}Z^{\top}Y$  is fast and easy, since (6.7) does not include any inverted matrices  $(L^{-1} = \text{diag }(l_1^{-1}, \ldots, l_m^{-1}))$  is explicitly known).
- 2. If some of the  $l_k$  are close to zero or rank (X) < m, some of the last few principal components (with small or even zero variance) of  $X^{\top}X$  can simply be excluded from the regression. This can be realized with the new estimator given by

$$\widetilde{\beta} = \sum_{k=1}^{p} l_k^{-1} a_k a_k^{\top} X^{\top} Y$$

p < m.

**Lemma 6.6.3.** Let rank (X) = m:

1. The estimator  $\tilde{\beta}$  is biased:

$$\mathbb{E}\widetilde{\beta} = \left(I - \sum_{k=p+1}^{m} a_k a_k^{\top}\right) \beta.$$

2. It holds that:

$$\operatorname{Cov}(\widetilde{\beta}) = \sigma^2 \sum_{k=1}^p l_k^{-1} a_k a_k^{\top}$$

Proof

1. Since

$$\widetilde{\beta} = \widehat{\beta} - \sum_{k=p+1}^{m} l_k^{-1} a_k a_k^{\top} X^{\top} Y$$

and  $\hat{\beta}$  is biased, it holds that

$$\mathbb{E}\widetilde{\beta} = \mathbb{E}\widehat{\beta} - \sum_{k=p+1}^{m} l_k^{-1} a_k a_k^{\top} X^{\top} \mathbb{E} Y$$

$$= \beta - \sum_{k=p+1}^{m} l_k^{-1} a_k \underbrace{a_k^{\top} X^{\top} X}_{l_k a_k^{\top}} \beta = \beta - \sum_{k=p+1}^{m} a_k a_k^{\top} \beta$$

$$= \left( I - \sum_{k=p+1}^{m} a_k a_k^{\top} \right) \beta$$

2.

#### Exercise 6.6.4.

Another equivalent formulation for regression with PCA is given in the following. Instead of using  $\gamma = A^{\top}\beta$ , use singular value decomposition (cf. Theorem 6.3.8) for X:

$$X = UL^{\frac{1}{2}}A^{\top},$$

where U is a  $(n \times m)$  matrix with orthonormal columns and L a diagonal matrix with  $L^{\frac{1}{2}} = \text{diag }(\sqrt{l_1}, \dots, \sqrt{l_m})$ . Define

$$\delta = L^{\frac{1}{2}} A^{\top} \beta, \tag{6.8}$$

then

$$Y = X\beta + \mathcal{E} = U \underbrace{L^{\frac{1}{2}}A^{\top}\beta}_{\delta} + \mathcal{E} = U\delta + \mathcal{E}.$$

The MLE for  $\delta$  is given by

$$\hat{\delta} = \underbrace{(U^{\top}U)^{-1}}_{I}U^{\top}Y = U^{\top}Y,$$

since U has orthonormal columns. (6.8) implies  $\beta = AL^{-\frac{1}{2}}\delta$  and thus

$$\hat{\beta} = AL^{-\frac{1}{2}}\hat{\delta} = AL^{-\frac{1}{2}}U^{\top}Y.$$

Here the relationship between  $\gamma$  and  $\delta$  is given by:

$$\gamma = A^{\top} \beta = A^{\top} \left( A L^{-\frac{1}{2}} \delta \right) = \underbrace{A^{\top} A}_{I} L^{-\frac{1}{2}} \delta = L^{-\frac{1}{2}} \delta.$$

Thus the following Lemma has been proven.

**Lemma 6.6.5.** The principal components  $Y = U\delta + \mathcal{E}$  of the regression  $Y = X\beta + \mathcal{E}$  has the MLE solution  $\hat{\delta} = U^{\top}Y$  resp.

$$\hat{\beta} = AL^{-\frac{1}{2}}U^{\top}Y. \tag{6.9}$$

Here the parameter vector  $\delta$  is simply a standardized version of  $\gamma$ :  $\delta = L^{\frac{1}{2}}\gamma$ .

#### Remark 6.6.6.

- 1. Since there are efficient algorithms for calculating a singular value decomposition, the term (6.9) can be calculated more efficiently compared to  $\hat{\beta} = (X^{\top}X)^{-1}X^{\top}Y$ , since  $X^{\top}X$  has to be inverted in the latter
- 2. Instead of removing the last m-p principal components of  $X^{\top}X$  from the regression (cf. Remark 6.6.2, 2.), it is generally possible to calculate  $\widetilde{\beta}$  on a subset M of  $\{1, \ldots, m\}$ :

$$\widetilde{\beta}_M = \sum_{k \in M} l_k^{-1} a_k a_k^{\top} X^{\top} Y.$$

Here, only principal components  $l_k$ ,  $k \in M$ , are used for the regression. Then it also holds that

$$\operatorname{Cov}(\widetilde{\beta}_M) = \sigma^2 \sum_{k \in M} l_k^{-1} a_k a_k^{\top},$$

cf. Exercise 6.6.4. This approach uses the elimination of components  $\gamma_k$ ,  $k \notin M$  of  $\gamma = (\gamma_1, \dots, \gamma_m)^{\top}$  of the ML estimation. Equivalently it can be thought of the exclusion of the components  $\delta_k$ ,  $k \notin M$  of  $\delta = (\delta_1, \dots, \delta_m)^{\top}$ , since  $\delta = L^{\frac{1}{2}}$ , with  $\delta_k = \sqrt{l_k} \gamma_k$  for all k.

What are possible strategies for determining M?

1.  $M = \{k : l_k > l^*\}$  for a predetermined threshold  $l^* > 0$ . If

$$\bar{l} = \frac{1}{m} \sum_{i=1}^{m} l_i$$

are close to 1,  $l^* \in (0.01, 0.1)$ . The disadvantage of this methodology is that some of the (possibly important for the forecast of Y) principal components, might have a small variance and are thus eliminated from the model.

2. Let  $\sigma_{ii}^2$  be the *i*-th diagonal element of  $(X^\top X)^{-1}$ . It holds that  $\sigma_{ii}^2 = \frac{\operatorname{Var}\hat{\beta}_i}{\sigma^2}$  (cf. Theorem 6.2.4),  $i = 1, \ldots, m$ . Then  $M = \{k : \sigma_{kk}^2 > \sigma^*\}$  can be chosen for a sufficient threshold  $\sigma^*$ . For the choice of  $\sigma^*$  see [18], p. 174. This methodology has the same disadvantages as 1.

- 3. Define  $M = \{1, ..., p\}$ , where p is the biggest number  $\leq m$ , for which one of the following criteria is met:
  - a) It holds that:

$$\sum_{i=1}^{m} \mathbb{E}(\widetilde{\beta}_{M_i} - \beta_i)^2 \le \sum_{i=1}^{m} \mathbb{E}(\widehat{\beta}_i - \beta_i)^2, \tag{6.10}$$

for all 
$$\beta = (\beta_1, \dots, \beta_m)^{\top} \in \mathbb{R}^m$$
.

b) It holds that:

$$\mathbb{E}(c^{\top} \widetilde{\beta}_{M} - c^{\top} \beta)^{2} \leq \mathbb{E}(c^{\top} \widehat{\beta} - c^{\top} \beta)^{2} \quad \forall \beta \in \mathbb{R}^{m}, c \in \mathbb{R}^{m}$$

c) It holds that:

$$\mathbb{E}\left|X\widetilde{\beta}_{M}-X\beta\right|^{2}\leq\mathbb{E}\left|X\widehat{\beta}-X\beta\right|^{2}$$

Here the criteria a) is similar to the task of estimating  $\beta$  as precise as possible. Criteria b) and c) on the other hand deliver the best possible estimation of of  $\mathbb{E}Y = X\beta$  with  $X\hat{\beta}_M$  resp.  $X\hat{\beta}$ . All terms in a)-c) are mean squared errors, which contain both the bias and the variance of  $\tilde{\beta}_M$ .

Many more strategies are described in statistical literature, which provide a better estimator  $\tilde{\beta}_M$  compared to  $\hat{\beta}$  depending on the given situation. The question on how to choose M is still unanswered.

An alternative approach of eliminating principal components in the regression is given by the following estimator  $\widetilde{\beta}_R$ :

$$\widetilde{\beta}_R = \sum_{k=1}^m (l_k + K_k)^{-1} a_k a_k^\top X^\top Y,$$

where  $K_1, \ldots, K_m > 0$  are weights, which represent additional influencing factors with respect to the regression. By using those weights it can be achieved, that  $l_k \approx 0$  does not have a destabilizing influence on the estimation

Exercise 6.6.7. Show that

Cov 
$$(\widetilde{\beta}_R) = \sigma^2 \sum_{k=1}^m \frac{l_k}{(l_k + K_k)^2} a_k a_k^{\top}$$

•  $\widetilde{\beta}_R$  is a biased estimator of  $\beta$ . Find the bias of  $\widetilde{\beta}_R$ !

 $\widetilde{\beta}_R$  is called *Ridge Regression*. Here the question arises on how to choose  $K_k$ , k = 1, ..., m. In practice,  $K_k = K$ , k = 1, ..., m is usually used, where K is to be chosen small.

Another application of PCA in regression is the so called *latent root regression*. This form of regression aims to only remove principal components, if they have a small variance  $l_k$  and do not add any additional value to the estimation of  $\mathbb{E}Y$  with  $X\beta$ . Here the PCA is applied to the  $(m+1)\times (m+1)$  matrix  $\widetilde{X}^{\top}\widetilde{X}$  with  $\widetilde{X}=(Y,X)$ . Let  $\widetilde{a}_k,\ k=0,\ldots,m$  be the coefficients of the PCA of  $\widetilde{X}^{\top}\widetilde{X}$ , with corresponding eigenvalues  $\widetilde{l}_k,\ k=0,\ldots,m$ . Let  $\widetilde{a}_k=(a_{k0},\ldots,a_{km})^{\top},\ k=0,\ldots,m$ .

Define the index set of the principal components that are to be eliminated as  $M_L = \{k = 0, \dots, m : \tilde{l}_k \leq l^*, |a_{k0}| \leq a^*\}$ . This is the index set of those principal components, that have small variance and do not influence the estimation of Y a lot. Let  $M = \{0, \dots, m\} \backslash M_L$ . Define  $\hat{\beta}_L = \sum_{k \in M} \tilde{c}_k \tilde{a}_k$ , where  $\{\tilde{c}_k, k \in M\} = \operatorname{argmin}_{\beta} |Y - X\beta|^2$  with  $\beta = \sum_{k \in M} c_k \tilde{a}_k$ .

Theorem 6.6.8. It holds that

$$\widetilde{c_k} = -\frac{a_{k0}\sqrt{\sum_{i=1}^n \left(Y_i - \overline{Y}\right)^2}}{\widetilde{l_k} \sum_{i \in M} \frac{a_{i0}^2}{\widetilde{l_i}}}, \quad k \in M.$$

#### Without proof!

Thresholds  $l^*$  and  $a^*$  are still to be chosen empirically.

## 6.7 Numeric calculation of principal components

In order to understand how statistical software packages calculates principal components, it is important to know the algorithms. By knowing the algorithms one can gain awareness about why some results might be bad (e.g. with eigenvalues that are almost equal) or what kind of restrictions there are with respect to size of the datasets (e.g. storage wise or runtime wise). In the following a short overview for those methods is given. Since the PCA is mainly based on calculating eigenvalues  $\lambda_i$  and eigenvectors  $\alpha_i$  of a positive semi-definite  $(m \times m)$  matrix  $\Sigma$ , the focus will mainly be on this calculation.

Let  $\Sigma$  thus be a  $(m \times m)$  matrix with eigenvectors  $\alpha_1, \ldots, \alpha_m$  and eigenvalues  $\lambda_1, \ldots, \lambda_m$ , which is positive semi-definite. In statistical literature there are at least four approaches for calculating  $\alpha_i$  and  $\lambda_i$ :

- 1. Power iteration.
- 2. QR decomposition,

- 3. Singular value decomposition,
- 4. Neural networks.

Here only the essence of power iteration is presented: it represents an iterative algorithm for determining  $\lambda_1$  and  $\alpha_1$ , if  $\lambda_1 >> \lambda_2 > \ldots > \lambda$ . Let  $u_0 \in \mathbb{R}^m$  be a starting vector. Define  $u_r = \Sigma u_{r-1} = \Sigma^r u_0$  for all  $r \in \mathbb{N}$ . If

$$u_0 = \sum_{i=1}^m c_i \alpha_i,$$

where  $\alpha_1, \ldots, \alpha_m$  are the orthonormal basis vectors and  $c_1, \ldots, c_m$  are coordinates then

$$u_r = \Sigma^r u_0 = \sum_{i=1}^m c_i \Sigma^r \alpha_i = \sum_{i=1}^m c_i \lambda_i^r \alpha_i, \quad r \in \mathbb{N}.$$

Let  $u_r = (u_{r1}, \dots, u_{rm})^\top$ ,  $\alpha_i = (\alpha_{i1}, \dots, \alpha_{im})^\top$ .

**Lemma 6.7.1.** It holds that

$$\frac{u_{ri}}{u_{r-1,i}} \xrightarrow[r \to \infty]{} \lambda_1$$

for  $i = 1, \ldots, m$  and

$$\frac{u_r}{c_i \lambda_1^r} \xrightarrow[r \to \infty]{} \alpha_1.$$

**Proof** For j = 1, ..., m it holds that

$$u_{rj} = \sum_{i=1}^{m} c_i \lambda_i^r \alpha_{ij}$$

and thus

$$\frac{u_{rj}}{u_{r-1,j}} = \frac{\sum_{i=1}^{m} c_i \lambda_i^r \frac{\alpha_{ij}}{\lambda_1^{r-1}}}{\sum_{i=1}^{m} c_i \lambda_i^{r-1} \frac{\alpha_{ij}}{\lambda_1^{r-1}}}$$

$$= \frac{c_1 \alpha_{1j} \lambda_1 + \sum_{i=2}^{m} c_i \left(\frac{\lambda_i}{\lambda_1}\right)^{r-1} \lambda_i \alpha_{ij}}{c_1 \alpha_{1j} + \sum_{i=2}^{m} c_i \left(\frac{\lambda_i}{\lambda_1}\right)^{r-1} \alpha_{ij}} \xrightarrow{r \to \infty} \frac{c_1 \alpha_{1j}}{c_1 \alpha_{1j}} \lambda_1 = \lambda_1,$$

since  $\frac{\lambda_i}{\lambda_1} < 1, i = 2, \dots$  Furthermore,

$$\frac{u_r}{u\lambda_1^r} = \alpha_1 + \sum_{i=2}^m \frac{c_i}{c_1} \left(\frac{\lambda_i}{\lambda_1}\right)^r \alpha_i \xrightarrow[r \to \infty]{} \alpha_1.$$

The fact, that  $c_1$  is unknown, is nothing to be worried about, since  $\frac{u_r}{\lambda_1^r}$  can be standardized. The proof of Lemma 6.6.5 implies, that the rate of convergence of  $\frac{u_{ri}}{u_{r-1,i}}$  to  $\lambda_1$  and  $\frac{u_r}{c_1\lambda_1^r}$  to  $\alpha_1$  is getting worse, if and only if  $\lambda_1 \approx \lambda_2$ , or in this case  $\frac{\lambda_2}{\lambda_1} \approx 1$ .

What should be done in the case that  $\lambda_1 \approx \lambda_2$ , in order to increase the rate of convergence? Instead of using  $\Sigma$ ,  $\Sigma - \rho I$  can be used for the iteration, in order to decrease the ratio  $\frac{\lambda_2 - \rho}{\lambda_1 - \rho}$ . Furthermore  $\Sigma$  can be replaced with  $(\Sigma - \rho I)^{-1}$ , which leads to solving the system of equations  $(\Sigma - \rho I) u_r = u_{r-1}$  for every  $r \in \mathbb{N}$ . Thus for a suitable choice of  $\rho$  a convergence to  $\alpha_k$ ,  $k = 1, \ldots, m$  is possible (in the second case).

#### Exercise 6.7.2. Construct those vectors and proof the convergence!

An increase in the rate of convergence can also be achieved, if one considers the sequence  $\{u_{2r}\}$  instead of  $\{u_r\}$  where  $u_{2r} = T^{2^r}u_0$ ,  $r \in \mathbb{N}$ . Further methodologies for improving the algorithm of power iteration can be found in [18], p. 410-411.

# **Bibliography**

- [1] H. Albrecher, S. A. Ladoucette, and J. L. Teugels. Asymptotics of the sample coefficient of variation and the sample dispersion. *J. Statist. Plann. Inference*, 140(2):358–368, 2010.
- [2] H. Dehling, B. Haupt. Einf"uhrung in die Wahrscheinlichkeitstheorie und Statistik. Springer, Berlin, 2003.
- [3] P. Bickel and K. Doksum. *Mathematical Statistics: Basic Ideas and Selected Topics*. Prentice Hall, London, 2001. 2nd ed., Vol. l.
- [4] A. A. Borovkov. Mathematical Statistics. Gordon & Breach, 1998.
- [5] M. Burkschat, E. Cramer, and U. Kamps. Beschreibende Statistik, Grundlegende Methoden. Springer, Berlin, 2004.
- [6] G. Casella and R. L. Berger. Statistical Inference. Pacific Grove (CA), Duxbury, 2002.
- [7] E. Cramer, K. Cramer, U. Kamps, and Zuckschwerdt. *Beschreibende Statistik, Interaktive Grafiken*. Springer, Berlin, 2004.
- [8] E. Cramer and U. Kamps. Grundlagen der Wahrscheinlichkeitsrechnung und Statistik. Springer, Berlin, 2007.
- [9] P. Dalgaard. Introductory Statistics with R. Springer, Berlin, 2002.
- [10] A.J Dobson. An Introduction to Generalizes Linear Models. Chapmen & Hall, Boca Raton, 2002.
- [11] Joseph L Doob. The limiting distributions of certain statistics. *The Annals of Mathematical Statistics*, 6(3):160–169, 1935.
- [12] L. Fahrmeir, T. Kneib, and S. Lang. Regression. Modelle, Methoden und Anwendungen. Springer, Berlin, 2007.
- [13] L. Fahrmeir, R. Künstler, I. Pigeot, and G. Tutz. *Statistik. Der Weg zur Datenanalyse*. Springer, Berlin, 2001.
- [14] H. O. Georgii. Stochastik. de Gruyter, Berlin, 2002.

BIBLIOGRAPHY 226

[15] J. Hartung, B. Elpert, and K. H. Klösener. *Statistik*. R. Oldenbourg Verlag, München, 1993. 9. Auflage.

- [16] C. C. Heyde and E. Seneta. Statisticians of the Centuries. Springer, Berlin, 2001.
- [17] A. Irle. Wahrscheinlichkeitstheorie und Statistik, Grundlagen, Resultate, Anwendungen. Teubner, 2001.
- [18] I. T. Jolliffe. *Principal component analysis*. Springer, 2nd edition edition, 2002.
- [19] L. J. Kazmir. Wirtschaftsstatistik. McGraw-Hill, 1996.
- [20] K. R. Koch. Parameter Estimation and Hypothesis Testing in Linear Models. Springer, Berlin, 1999.
- [21] A. Krause and M. Olson. *The Basics of S-PLUS*. Springer, Berlin, 2002. Third Ed.
- [22] E. L. Lehmann. *Elements of Large-Sample Theory*. Springer, New York, 1999.
- [23] J. Lehn and H. Wegmann. *Einführung in die Statistik*. Teubner, Stuttgart, 2000. 3. Auflage.
- [24] J. Maindonald and J. Braun. *Data Analysis and Graphics Using R.* Cambridge University Press, 2003.
- [25] M. Overbeck-Larisch and W. Dolejsky. *Stochastik mit Mathematica*. Vieweg, Braunschweig, 1998.
- [26] H. Pruscha. Angewandte Methoden der Mathematischen Statistik. Teubner, Stuttgart, 2000.
- [27] H. Pruscha. Vorlesungen über Mathematische Statistik. Teubner, Stuttgart, 2000.
- [28] L. Sachs. Angewandte Statistik. Springer, 2004.
- [29] L. Sachs and J. Hedderich. Angewandte Statistik, Methodensammlung mit R. Springer, Berlin, 2006.
- [30] Robert J Serfling. Approximation theorems of mathematical statistics, volume 162. John Wiley & Sons, 2009.
- [31] M. R. Spiegel and L. J. Stephens. Statistik. McGraw-Hill, 1999.
- [32] E. Spodarev. Wahrscheinlichkeitstheorie und stochastische Prozesse. Ulm, 2020.

BIBLIOGRAPHY 227

[33] E. Spodarev. Elementare Wahrscheinlichkeitsrechnung und Statistik. Ulm, 2022.

- [34] V. Spokoiny and T. Dickhaus. *Basics of modern mathematical statistics*. Springer, 2015.
- [35] W. A. Stahel. Statistische Datenanalyse. Vieweg, 1999.
- [36] W. Venables and D. Ripley. *Modern applied statistics with S-PLUS*. Springer, 1999. 3rd ed.
- [37] L. Wasserman. All of Statistics. A Concise Course in Statistical Inference. Springer, 2004.

# Index

Acceptance region, 75	Convolution stability of the multivariate
AIC criteria, 196	normal distribution, 133
Akaike infromation coefficient, 196	Cramér-Rao, inequality, 36
ANOVA, see analysis of variance	Critical region, see Rejection region
asymptotic tests, 189	
	Decision rule, 74
Bayes estimator, 29	design matrix, 128, 141
Bayes formula, 29	Distribution with monotone likelihood
Bernoulli distribution	ratio, 97
Asymptotic confidence interval, 65	
best linear unbiased estimator, 144	effect, 174
bilinear form, 134	Eindeutigkeit der besten er-
Binomial distribution, 98	wartungstreuen Schätzer, 49
Blackwell-Rao, Inequality of, 51	empirical
Bonferroni inequality, 155	Probability distribution function,
Bootstrap	11
confidence intervals, 36	Erlang distribution, 5
estimator, 34	error terms, 141
Bootstrap estimators	Estimators
Monte-Carlo methods, 35	sufficient, 42
	Exceedance probability, 81
$\chi^2$ distribution, 5	exponential family, 177
class specific differences, 174	
Class-sizes, 109	Fisher
classical ANOVA hypothesis, 174	Fisher information, 23
coefficient of determination, 153	Fisher-Snedecor distribution, F dis-
Complete, 47	tribution, 8
Confidence interval, 58	Fisher information matrix, 117, 184, 195
Asymptotic	Fisher Scoring, 194
Bernoulli distribution, 65	Fisher's scoring method, 189
Poisson distribution, 66	Fundamental theorem for two-sided
asymptotical, 59	tests, 107
Bootstrap, 36	
length, 59	Gamma distribution
minimal, 59	Moment generating and character-
Confidence Intevals	istic function, 3
Asymptotic, 64	Stability, 4
Confidence level, 58	Theorem of Gauß-Markov, 165
Confusion matrix, 76	generalized inverse matrix, 158

Hesse matrix, 184	Neyman-Fisher Factorisation Theorem,
Hoeffding inequality, 63	45
Hypotheses, 74	Neyman-Fisher, Factorisation Theorem,
Alternative, 74	45
Main hypothesis, 74	Neyman-Pearson
hypotheses	Fundamental Lemma, 93
testable, 170	Optimality theorem, 92
testable, 170	non-centered $\chi_{n,\mu}^2$ distribution, 137
identifiable, 1	Normal distribution
iteration test, 125	
	Confidence interval
Jackknife estimator for the	for two samples, 68
Bias, 33	confidence interval
Expectation, 32	One sample, 60
Variance, 33	multivariate, 129
variance, so	Significance tests, 84
Karl Popper, 75	normal equation, 142
Kullback-Leibler information, 19	
,	Odds, 182
least squares	OLS estimator, 142
ordinary, 143	one-parametric exponential class, 97
Lehmann-Scheffé, Theorem, 50	
level of influencing factor, 174	p-value, 80
Likelihood function, 16	parameter space, 1
likelihood ratio test, 195	parameter vector, 1
linear form, 133	Pearson test statistic is introduced, 110
linear regression, 128	Performance function, 76
multiple, 144	Plug-in estimator, 12
without full rank, 158	Plug-in method, 12
multivariate with full rank, 141	point estimation, 1
	point estimator, 1
Linear transformation of $\mathcal{N}(\mu, K)$ , 133 link function, 177	Poisson distribution, 70, 86, 88
Link functions	Asymptotic confidence interval, 66
	Neyman-Fisher test, 120
natural, 181	Neyman-Pearson test, 95
logit model, 198	Poisson model, 198
Logit modell, 182	Poisson regression, 188
Logit-Modell, 192	
loss function, 29	posteriori distribution, 29
Mi I:ll:llt: 15 16	Power function, 76
Maximum-Likelihood estimator, 15, 16	predictor variables, 128
weak consistency, 21	prior distribution, 29
Method of least squares, 142	Probit model, 182
Method of moments estimation, 13	Procedure of Cramér-Wold, 131
Mistakes of type I and II, 76	1 6 104
mixed moments, 134	quadratic form, 134
models	Covariance, 134
generalized linear, 177	Quantile function of the normal distri-
Moment estimator, 14	bution, 182
Multinomial distribution, 109	<b>.</b>
	Randomization region, 75
Newton's method, 188	regression

BIBLIOGRAPHY 230

binary categorical, 182 logistic, 182, 188 Rejection region, 75 related sample, 68 Resampling methods, 32 residual distribution, 153 residuals, 152 response variable, 141 Score function, 195 Score statistic, 195 Significance level, 58 Sufficient estimators, 42 $t$ distribution, 6 Test  Asymptotic, 78, 85 Binomial test, 123 $\chi^2$ goodness-of-fit test, 109 $\chi^2$ -Pearson-Fisher test, 115 for connectivity, 152 for regression parameters, 151 Goodness-of-fit test, 108 Iteration test, 125 Kolmogorov-Smirnov, 109 Monte-Carlo test, 78 most powerful, 90	one-sided, 77 right-sided, 77 two-sided, 77 Parametric significance test, 84 Power, 76 powerful, 90 randomized, 75, 89 Scope, 90 of Shapiro-Francia, 122 of Shapiro-Wilk, 123 Shapiros goodness-of-fit test, 121 unbiased, 82 Wald test, 85 of Wald-Wolfowitz, 127 Test statistic, 60 Theorem $\chi^2 \text{ Distribution, special case, 5}$ $\text{Cramér-Rao inequality, 36}$ $\text{Density of the t distribution, 6}$ Factorisation Theorem of Neyman- Fisher, 45 Lehmann-Scheffé, 50 Moment generating and characteristic function of the Gamma distribution, 3 weak consistency of ML estimators, 21
Neyman-Pearson test, 91 One sided, 96 parameter of the Poisson distribution, 95 Rejection region, 91 Scope, 91 Neyman-Pearson-Test modifizierter, 104 NP test, see Neyman-Pearson test Parameters of the normal distribution, 84 parametric, 77 left-sided, 77	unimodal, 21 Uniqueness theorem for characteristic functions, 130 for moment generating functions, 137  variability of the expected values, 174 Variance analysis, 174 single factor, 174 two factor, 176  Wald statistic, 195