



---

# Stochastics III

---

Ulm University  
Institute of Stochastics

Lecture Notes  
Prof. Dr. Volker Schmidt  
Winter Term 2012/2013

ULM, FEBRUARY 2013

## Contents

<b>1</b>	<b>Introduction and Mathematical Foundations</b>	<b>6</b>
1.1	Some Basic Notions and Results of Matrix Algebra . . . . .	6
1.1.1	Trace and Rank . . . . .	6
1.1.2	Eigenvalues and Eigenvectors . . . . .	7
1.1.3	Diagonalization Method . . . . .	8
1.1.4	Symmetric and Definite Matrices; Factorization . . . . .	9
1.2	Multivariate Normal Distribution . . . . .	10
1.2.1	Definition and Fundamental Properties . . . . .	10
1.2.2	Characteristics of the Multivariate Normal Distribution . . . . .	12
1.2.3	Marginal Distributions and Independence of Subvectors; Convolution Properties . . . . .	14
1.2.4	Linear Transformation of Normally Distributed Random Vectors . . . . .	16
1.2.5	Degenerate Multivariate Normal Distribution . . . . .	17
1.3	Linear and Quadratic Forms of Normally Distributed Random Vectors . . . . .	18
1.3.1	Definition, Expectation and Covariance . . . . .	18
1.3.2	Noncentral $\chi^2$ -Distribution . . . . .	21
1.3.3	Distributional Properties of Linear and Quadratic Forms . . . . .	23
<b>2</b>	<b>Linear Models; Design Matrix with Full Rank</b>	<b>26</b>
2.1	Method of Least Squares . . . . .	27
2.1.1	Normal Equation . . . . .	27
2.1.2	Properties of the LS-Estimator $\hat{\beta}$ . . . . .	29
2.1.3	Unbiased Estimation of the Variance $\sigma^2$ of the Error Terms . . . . .	31
2.2	Normally Distributed Error Terms . . . . .	33
2.2.1	Maximum-Likelihood Estimation . . . . .	33
2.2.2	Distributional Properties of $\hat{\beta}$ and $S^2$ . . . . .	35
2.2.3	Statistical Tests for the Regression Coefficients . . . . .	36
2.2.4	Confidence Regions; Prediction of Response Variables . . . . .	41
2.2.5	Confidence Band . . . . .	42
<b>3</b>	<b>Arbitrary Design Matrix; Generalized Inverse</b>	<b>45</b>
3.1	Analysis of Variance as a Linear Model . . . . .	45
3.1.1	One-Factor Analysis of Variance; ANOVA Null Hypothesis . . . . .	45
3.1.2	Reparametrization of the Expectations . . . . .	48
3.1.3	Two-Factor Analysis of Variance . . . . .	50
3.2	Estimation of Model Parameters . . . . .	52
3.2.1	LS-Estimator for $\beta$ . . . . .	53

3.2.2	Expectation Vector and Covariance Matrix of the LS-Estimator $\bar{\beta}$ . . . . .	57
3.2.3	Estimable Functions . . . . .	59
3.2.4	Best Linear Unbiased Estimator; Gauss–Markov Theorem . . . . .	62
3.3	Normally Distributed Error Terms . . . . .	65
3.3.1	Maximum–Likelihood Estimation . . . . .	66
3.3.2	Testing Linear Hypotheses . . . . .	69
3.3.3	Confidence Regions . . . . .	73
3.4	Examples . . . . .	75
3.4.1	F–Test for the ANOVA Null Hypothesis . . . . .	75
3.4.2	F–Tests for the Two–Factor Analysis of Variance . . . . .	77
3.4.3	Two–Factor Analysis of Variance with Hierarchical Classification . . . . .	81
<b>4</b>	<b>Generalized Linear Models</b> . . . . .	<b>84</b>
4.1	Definition and Basic Properties . . . . .	84
4.1.1	Exponential Family . . . . .	84
4.1.2	Link of the Parameters; Natural Link Function . . . . .	86
4.2	Examples . . . . .	86
4.2.1	Linear Model with Normally Error Terms . . . . .	86
4.2.2	Binary Categorical Regression . . . . .	87
4.2.3	Poisson–Distributed Sample Variables with Natural Link Function . . . . .	88
4.3	Maximum–Likelihood Estimator for $\beta$ . . . . .	88
4.3.1	Loglikelihood Function and its Partial Derivatives . . . . .	88
4.3.2	Hessian Matrix . . . . .	90
4.3.3	Maximum–Likelihood Equation and Numerical Approach . . . . .	92
4.3.4	Asymptotic Normality of ML Estimators; Asymptotic Tests . . . . .	94
4.4	Weighted LS Estimator for Categorical Regression . . . . .	95
4.4.1	Estimation of the Expectation Vector . . . . .	95
4.4.2	Asymptotic Normality of the LS–Estimator . . . . .	97
4.4.3	Evaluation of the Goodness of Fit . . . . .	99
<b>5</b>	<b>Goodness–of–Fit Tests</b> . . . . .	<b>100</b>
5.1	Kolmogorov–Smirnov Test . . . . .	100
5.1.1	Empirical Distribution Function; KS Test Statistic . . . . .	100
5.1.2	Asymptotic Distribution . . . . .	101
5.1.3	Pointwise and Uniform Consistency . . . . .	104
5.2	$\chi^2$ –Goodness–of–Fit Test . . . . .	106
5.2.1	Aggregation; Pearson–Statistic . . . . .	106
5.2.2	Asymptotic Distribution . . . . .	108

5.2.3	Goodness-of-Fit; Local Alternatives . . . . .	109
5.3	Pearson-Fisher Test . . . . .	111
5.3.1	Pearson-Fisher Test Statistic . . . . .	112
5.3.2	Multivariate Central Limit Theorem for ML Estimators . . . . .	113
5.3.3	Fisher-Information Matrix and Central Limit Theorem in the Coarsened Model . . . . .	114
5.3.4	Asymptotic Distribution of the Pearson-Fisher Statistic . . . . .	116
5.4	Examples . . . . .	119
5.4.1	Pearson-Fisher Test for Poisson-Distribution . . . . .	119
5.4.2	Pearson-Fisher Test for Normal Distribution . . . . .	120
5.4.3	Goodness-of-Fit Tests of Shapiro-Wilk-Type . . . . .	121
<b>6</b>	<b>Nonparametric Localization Tests</b> . . . . .	<b>124</b>
6.1	Two Simple Examples of One-Sample Problems . . . . .	124
6.1.1	Binomial Test . . . . .	124
6.1.2	Run Test for Randomness . . . . .	126
6.2	Wilcoxon-Rank Test . . . . .	128
6.2.1	Model Description; Median Test . . . . .	128
6.2.2	Distribution of the Test Statistic $T_n^+$ for Small Sample Sizes . . . . .	129
6.2.3	Asymptotic Distribution . . . . .	132
6.3	Two-Sample Problems . . . . .	134
6.3.1	Run Test of Wald-Wolfowitz . . . . .	134
6.3.2	Wilcoxon Rank-Sum Test for Location Alternatives . . . . .	135

## References

- [1] Büning, H., Trenkler, G. (1994)  
Nichtparametrische statistische Methoden  
de Gruyter, Berlin
- [2] Cressie, N.A. (1993)  
Statistics for Spatial Data  
J. Wiley & Sons, New York
- [3] Dobson, A.J. (2002)  
An Introduction to Generalized Linear Models  
Chapman & Hall, Boca Raton
- [4] Falk, M., Marohn, F., Tewes, B. (2002)  
Foundations of Statistical Analyses and Applications with SAS  
Birkhäuser, Basel
- [5] Hastie, T., Tibshirami, R., Friedman, J. (2001)  
The Elements of Statistical Learning  
Springer, New York
- [6] Koch, K.R. (1997)  
Parameterschätzung und Hypothesentests in linearen Modellen  
Dümmlers-Verlag, Bonn
- [7] Lehmann, E.L. (1999)  
Elements of Large-Sample Theory  
Springer, New York
- [8] Lehmann, E.L., Romano, J.P. (2005)  
Testing Statistical Hypotheses  
Springer, New York
- [9] McCullagh, P., Nelder, J.A. (1989)  
Generalized Linear Models  
Chapman & Hall, London.
- [10] Pruscha, H. (2000)  
Vorlesungen über mathematische Statistik  
Teubner-Verlag, Stuttgart.
- [11] Van der Vaart, A., Wellner, J. (1996)  
Weak Convergence and Empirical Processes  
Springer-Verlag, New York
- [12] Vapnik, V.N. (1998)  
Statistical Learning Theory  
J. Wiley & Sons, New York

# 1 Introduction and Mathematical Foundations

These lecture notes are made for students who already have a basic knowledge of mathematical statistics. Estimation and statistical test methods which have been discussed in "Stochastik I" are assumed to be known.

The present lecture notes consist of the following parts:

- multivariate normal distribution (nondegenerate and degenerate normal distribution, linear and quadratic forms)
- linear models (multiple regression, normally distributed error terms, single- and multiple-factor analysis of variance)
- generalized linear models (logistic regression, maximum-likelihood equation, weighted least squares estimator, evaluation of the goodness of fit)
- tests for distribution assumptions (Kolmogorow-Smirnow test,  $\chi^2$ -goodness-of-fit test of Pearson-Fisher)
- nonparametric location tests (binomial test, iteration tests, linear rank tests)

In particular, we will use notions and results which have been introduced in the lecture notes "Elementare Wahrscheinlichkeitsrechnung und Statistik" and "Stochastik I": we will indicate references to these lecture notes by "WR" and "I" in front of the section number of the cited lemmas, theorems, corollaries and formulas.

## 1.1 Some Basic Notions and Results of Matrix Algebra

First, we recall some basic notions and results of matrix algebra, which are needed in these lecture notes.

### 1.1.1 Trace and Rank

- The *trace*  $\text{tr}(\mathbf{A})$  of a quadratic  $n \times n$  matrix  $\mathbf{A} = (a_{ij})$  is given by

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii}. \quad (1)$$

- Let  $\mathbf{A}$  be an arbitrary  $n \times m$  matrix. The *rank*  $\text{rk}(\mathbf{A})$  is the maximum number of linearly independent rows (or columns) of  $\mathbf{A}$ .
  - The vectors  $\mathbf{a}_1, \dots, \mathbf{a}_\ell \in \mathbb{R}^m$  are called *linearly dependent* if there exist real numbers  $c_1, \dots, c_\ell \in \mathbb{R}$ , which are not all equal to zero and  $c_1 \mathbf{a}_1 + \dots + c_\ell \mathbf{a}_\ell = \mathbf{o}$ .
  - Otherwise the vectors  $\mathbf{a}_1, \dots, \mathbf{a}_\ell \in \mathbb{R}^m$  are called *linearly independent*.

From the definition of the trace of a matrix in (1) and from the definition of matrix multiplication the next lemma directly follows.

**Lemma 1.1** *Let  $\mathbf{C}$  be an arbitrary  $n \times m$  matrix and  $\mathbf{D}$  an arbitrary  $m \times n$  matrix. Then  $\text{tr}(\mathbf{CD}) = \text{tr}(\mathbf{DC})$ .*

It can be proved that a quadratic matrix  $\mathbf{A}$  is invertible if and only if  $\mathbf{A}$  has full rank or  $\det \mathbf{A} \neq 0$ , respectively. The following result is also useful in this context.

**Lemma 1.2** *Let  $\mathbf{A}$  be an  $n \times m$  matrix with  $n \geq m$  and  $\text{rk}(\mathbf{A}) = m$ . Then  $\text{rk}(\mathbf{A}^\top \mathbf{A}) = m$ .*

**Proof**

- It is obvious that the rank  $\text{rk}(\mathbf{A}^\top \mathbf{A})$  of the  $m \times m$  matrix  $\mathbf{A}^\top \mathbf{A}$  cannot exceed  $m$ .
- Now, we assume that  $\text{rk}(\mathbf{A}^\top \mathbf{A}) < m$ . Then, there exists a vector  $\mathbf{c} = (c_1, \dots, c_m)^\top \in \mathbb{R}^m$  with  $\mathbf{c} \neq \mathbf{o}$  and  $\mathbf{A}^\top \mathbf{A} \mathbf{c} = \mathbf{o}$ .
- From this follows that  $\mathbf{c}^\top \mathbf{A}^\top \mathbf{A} \mathbf{c} = \mathbf{o}$  and  $(\mathbf{A} \mathbf{c})^\top (\mathbf{A} \mathbf{c}) = \mathbf{o}$ , i.e.,  $\mathbf{A} \mathbf{c} = \mathbf{o}$ .
- However, this is contradictory to the assumption that  $\text{rk}(\mathbf{A}) = m$ . □

Furthermore, it can be proved that the following properties of trace and rank are valid.

**Lemma 1.3** *Let  $\mathbf{A}$  and  $\mathbf{B}$  be arbitrary  $n \times n$  matrices. Then  $\text{tr}(\mathbf{A} - \mathbf{B}) = \text{tr}(\mathbf{A}) - \text{tr}(\mathbf{B})$  always holds. If  $\mathbf{A}$  is idempotent and symmetric, i.e.,  $\mathbf{A} = \mathbf{A}^2$  and  $\mathbf{A} = \mathbf{A}^\top$ , it also holds that  $\text{tr}(\mathbf{A}) = \text{rk}(\mathbf{A})$ .*

**1.1.2 Eigenvalues and Eigenvectors**

**Definition** Let  $\mathbf{A}$  be an arbitrary  $n \times n$  matrix. Each (complex) number  $\lambda \in \mathbb{C}$  is called an *eigenvalue* of the matrix  $\mathbf{A}$  if and only if there exists a vector  $\mathbf{x} \in \mathbb{C}^n$  with  $\mathbf{x} \neq \mathbf{o}$  and

$$(\mathbf{A} - \lambda \mathbf{I})\mathbf{x} = \mathbf{o}. \quad (2)$$

We call  $\mathbf{x}$  an *eigenvector* corresponding to  $\lambda$ .

**Remark**

- Only if  $\lambda$  is a solution of the so-called *characteristic equation*

$$\det(\mathbf{A} - \lambda \mathbf{I}) = 0, \quad (3)$$

there is a solution  $\mathbf{x} \in \mathbb{C}^n$  with  $\mathbf{x} \neq \mathbf{o}$  for (2). The left-hand side  $P(\lambda) = \det(\mathbf{A} - \lambda \mathbf{I})$  of (3) is called the *characteristic polynomial* of matrix  $\mathbf{A}$ .

- Let  $\lambda_1, \dots, \lambda_k \in \mathbb{R}$  be the real-valued solutions of (3). Then the characteristic polynomial can be written in the form

$$P(\lambda) = (-1)^n (\lambda - \lambda_1)^{a_1} \dots (\lambda - \lambda_k)^{a_k} q(\lambda), \quad (4)$$

where  $a_1, \dots, a_k \in \mathbb{N}$  are positive natural numbers, the so-called *algebraic multiplicities* of  $\lambda_1, \dots, \lambda_k$ , and  $q(\lambda)$  is a polynomial of order  $n - \sum_{i=1}^k a_i$  which has no real solutions.

**Lemma 1.4** *Let  $\mathbf{A} = (a_{ij})$  be a symmetric  $n \times n$  matrix with real-valued entries  $a_{ij}$ . Then every eigenvalue is a real number and eigenvectors  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^n$  which correspond to different eigenvalues  $\lambda_i, \lambda_j \in \mathbb{R}$  are orthogonal to each other.*

**Proof**

- The determinant  $\det(\mathbf{A} - \lambda \mathbf{I})$  in (3) is given by

$$\det(\mathbf{A} - \lambda \mathbf{I}) = \sum_{\boldsymbol{\pi}} (-1)^{r(\boldsymbol{\pi})} \prod_{i: i \neq \pi_i} a_{i\pi_i} \prod_{i: i = \pi_i} (a_{i\pi_i} - \lambda), \quad (5)$$

where the summation extends over all  $m!$  permutations  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_m)$  of the natural numbers  $1, \dots, m$  and  $r(\boldsymbol{\pi})$  is the number of pairs in  $\boldsymbol{\pi}$ , which are not in the natural order.

- Since the elements of  $\mathbf{A}$  are real numbers, every solution  $\lambda = a + ib$  of (3) implies another solution  $\bar{\lambda} = a - ib$  of (3).
- Let  $\mathbf{x} = \mathbf{a} + i\mathbf{b}$  and  $\bar{\mathbf{x}} = \mathbf{a} - i\mathbf{b}$  be eigenvectors which correspond to  $\lambda$  or  $\bar{\lambda}$ , respectively. Then  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  and  $\mathbf{A}\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}}$  or

$$\bar{\mathbf{x}}^\top \mathbf{A}\mathbf{x} = \bar{\mathbf{x}}^\top \lambda\mathbf{x} = \lambda\bar{\mathbf{x}}^\top \mathbf{x}$$

and

$$\bar{\mathbf{x}}^\top \mathbf{A}\mathbf{x} = (\mathbf{A}^\top \bar{\mathbf{x}})^\top \mathbf{x} = (\mathbf{A}\bar{\mathbf{x}})^\top \mathbf{x} = (\bar{\lambda}\bar{\mathbf{x}})^\top \mathbf{x} = \bar{\lambda}\bar{\mathbf{x}}^\top \mathbf{x}.$$

- From this it follows that  $\lambda\bar{\mathbf{x}}^\top \mathbf{x} = \bar{\lambda}\bar{\mathbf{x}}^\top \mathbf{x}$ .
- Since  $\bar{\mathbf{x}}^\top \mathbf{x} = |\mathbf{a}|^2 + |\mathbf{b}|^2 > 0$ , it holds that  $\lambda = \bar{\lambda}$ , i.e.,  $\lambda$  is a real number.
- In a similar way it can be proved that for different eigenvalues  $\lambda_i, \lambda_j \in \mathbb{R}$  there exist eigenvectors  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^n$  with real-valued components which are orthogonal to each other.
- Since the matrix  $\mathbf{A} - \lambda_i\mathbf{I}$  only contains real-valued elements, it holds that if  $\mathbf{x}_i$  is an eigenvector which corresponds to  $\lambda_i$ , then also  $\bar{\mathbf{x}}_i$  and  $\mathbf{x}_i + \bar{\mathbf{x}}_i \in \mathbb{R}^n$  are eigenvectors that correspond to  $\lambda_i$ .
- Therefore we can (and will) assume w.l.o.g. that  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^n$ . Furthermore, if

$$(\mathbf{A} - \lambda_i\mathbf{I})\mathbf{x}_i = \mathbf{0} \quad \text{and} \quad (\mathbf{A} - \lambda_j\mathbf{I})\mathbf{x}_j = \mathbf{0},$$

it follows that  $\mathbf{A}\mathbf{x}_i = \lambda_i\mathbf{x}_i$  and  $\mathbf{A}\mathbf{x}_j = \lambda_j\mathbf{x}_j$  as well as

$$\mathbf{x}_j^\top \mathbf{A}\mathbf{x}_i = \lambda_i\mathbf{x}_j^\top \mathbf{x}_i \quad \text{and} \quad \mathbf{x}_i^\top \mathbf{A}\mathbf{x}_j = \lambda_j\mathbf{x}_i^\top \mathbf{x}_j.$$

- On the other hand it is obvious that  $\mathbf{x}_j^\top \mathbf{x}_i = \mathbf{x}_i^\top \mathbf{x}_j$  and with the symmetry of  $\mathbf{A} = (a_{ij})$  we get the identity  $\mathbf{x}_j^\top \mathbf{A}\mathbf{x}_i = \mathbf{x}_i^\top \mathbf{A}\mathbf{x}_j$  since

$$\mathbf{x}_j^\top \mathbf{A}\mathbf{x}_i = \sum_{m=1}^n \sum_{\ell=1}^n x_{\ell j} a_{\ell m} x_{m i} = \sum_{\ell=1}^n \sum_{m=1}^n x_{m i} a_{m \ell} x_{\ell j} = \mathbf{x}_i^\top \mathbf{A}\mathbf{x}_j.$$

- Altogether it follows that  $\lambda_i\mathbf{x}_j^\top \mathbf{x}_i = \lambda_j\mathbf{x}_i^\top \mathbf{x}_j$  and  $(\lambda_i - \lambda_j)\mathbf{x}_j^\top \mathbf{x}_i = 0$ .
- As  $\lambda_i - \lambda_j \neq 0$ , it holds that  $\mathbf{x}_j^\top \mathbf{x}_i = 0$ . □

### 1.1.3 Diagonalization Method

- Now, let  $\mathbf{A}$  be an invertible symmetric  $n \times n$  matrix.
- In Lemma 1.4 we have shown that all eigenvalues  $\lambda_1, \dots, \lambda_n$  of  $\mathbf{A}$  are real numbers (where it is possible that one number occurs more than once in this sequence).
- Since  $\det \mathbf{A} \neq 0$ , we get that  $\lambda = 0$  is no solution of (3), i.e., all eigenvalues  $\lambda_1, \dots, \lambda_n$  of  $\mathbf{A}$  are different from zero.
- Furthermore, it can be proved that there are *orthonormal* (basis) vectors  $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$ , i.e.,

$$\mathbf{v}_i^\top \mathbf{v}_i = 1, \quad \mathbf{v}_i^\top \mathbf{v}_j = 0, \quad \forall i, j \in \{1, \dots, n\} \text{ with } i \neq j, \quad (6)$$

such that  $\mathbf{v}_i$  is an eigenvector that corresponds to  $\lambda_i$ ;  $i = 1, \dots, n$ .

- If all eigenvalues  $\lambda_1, \dots, \lambda_n$  differ from each other, then this is an immediate consequence of part 2 of Lemma 1.4.
- As a consequence, the following *diagonalization method* for invertible symmetric matrices is obtained.



**Lemma 1.5**

- Let  $\mathbf{A}$  be an invertible symmetric  $n \times n$  matrix and let  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$  be the  $n \times n$  matrix that consists of the orthonormal eigenvalues  $\mathbf{v}_1, \dots, \mathbf{v}_n$ .
- Then

$$\mathbf{V}^\top \mathbf{A} \mathbf{V} = \mathbf{\Lambda}, \quad (7)$$

where  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$  denotes the  $n \times n$  diagonal matrix which consists of the eigenvalues  $\lambda_1, \dots, \lambda_n$ .

**Proof**

- Equation (2) in the definition of eigenvalues and eigenvectors implies that  $\mathbf{A} \mathbf{v}_i = \lambda_i \mathbf{v}_i$  for each  $i = 1, \dots, n$ .
- This means that  $\mathbf{A} \mathbf{V} = (\lambda_1 \mathbf{v}_1, \dots, \lambda_n \mathbf{v}_n)$  and with (6) it follows that  $\mathbf{V}^\top \mathbf{A} \mathbf{V} = \mathbf{V}^\top (\lambda_1 \mathbf{v}_1, \dots, \lambda_n \mathbf{v}_n) = \mathbf{\Lambda}$ .  $\square$

**1.1.4 Symmetric and Definite Matrices; Factorization**

**Lemma 1.6** Let  $\mathbf{A}$  be a symmetric and positive definite  $n \times n$  matrix, i.e.,  $\mathbf{A} = \mathbf{A}^\top$  and  $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$  for each vector  $\mathbf{x} = (x_1, \dots, x_n)^\top \in \mathbb{R}^n$  with  $\mathbf{x} \neq \mathbf{o}$ . Then  $\mathbf{A}$  is invertible and there is an invertible  $n \times n$  matrix  $\mathbf{H}$ , such that

$$\mathbf{A} = \mathbf{H} \mathbf{H}^\top. \quad (8)$$

**Proof** We only prove the second part of Lemma 1.6.

- Lemma 1.5 implies that  $\mathbf{V}^\top \mathbf{A} \mathbf{V} = \mathbf{\Lambda}$  and

$$\mathbf{A} = (\mathbf{V}^\top)^{-1} \mathbf{\Lambda} \mathbf{V}^{-1}, \quad (9)$$

- where  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$  is the  $n \times n$  matrix which consists of the orthonormal eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$ ,
- and  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$  denotes the  $n \times n$  diagonal matrix which consists of the (positive) eigenvalues  $\lambda_1, \dots, \lambda_n$ .
- Now, let  $\mathbf{\Lambda}^{1/2}$  be the  $n \times n$  diagonal matrix  $\mathbf{\Lambda}^{1/2} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$  and let

$$\mathbf{H} = (\mathbf{V}^\top)^{-1} \mathbf{\Lambda}^{1/2} \mathbf{V}^\top. \quad (10)$$

- It is obvious that the matrix  $\mathbf{H}$ , given in (10), is invertible. Because of  $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$  it also holds that

$$\begin{aligned} \mathbf{H} \mathbf{H}^\top &= (\mathbf{V}^\top)^{-1} \mathbf{\Lambda}^{1/2} \mathbf{V}^\top \left( (\mathbf{V}^\top)^{-1} \mathbf{\Lambda}^{1/2} \mathbf{V}^\top \right)^\top = (\mathbf{V}^\top)^{-1} \mathbf{\Lambda}^{1/2} \mathbf{V}^\top \mathbf{V} \mathbf{\Lambda}^{1/2} \mathbf{V}^{-1} \\ &= (\mathbf{V}^\top)^{-1} \mathbf{\Lambda}^{1/2} \mathbf{\Lambda}^{1/2} \mathbf{V}^{-1} = (\mathbf{V}^\top)^{-1} \mathbf{\Lambda} \mathbf{V}^{-1} = \mathbf{A}, \end{aligned}$$

where the last equality follows from (9).  $\square$

**Remark**

- Each invertible  $n \times n$  matrix  $\mathbf{H}$  with  $\mathbf{A} = \mathbf{H} \mathbf{H}^\top$  is called a *square root* of  $\mathbf{A}$  and is denoted by  $\mathbf{A}^{1/2}$ .
- Using the *Cholesky decomposition* for symmetric and positive definite matrices, one can show that there exists a (uniquely determined) lower triangular matrix  $\mathbf{H}$  with  $\mathbf{A} = \mathbf{H} \mathbf{H}^\top$ .

The following property of symmetric matrices is a generalization of Lemma 1.6.

**Lemma 1.7** Let  $\mathbf{A}$  be a symmetric and positive semidefinite  $n \times n$  matrix, i.e., it holds that  $\mathbf{A} = \mathbf{A}^\top$  and  $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$  for each vector  $\mathbf{x} = (x_1, \dots, x_n)^\top \in \mathbb{R}^n$ . Now, let  $\text{rk}(\mathbf{A}) = r (\leq n)$ . Then there exists an  $n \times r$  matrix  $\mathbf{H}$  with  $\text{rk}(\mathbf{H}) = r$ , such that  $\mathbf{A} = \mathbf{H}\mathbf{H}^\top$ .

The *proof* of Lemma 1.7 is similar to the proof of Lemma 1.6.

**Lemma 1.8**

- Let  $m, r \in \mathbb{N}$  be arbitrary natural numbers with  $1 \leq r \leq m$ . Let  $\mathbf{A}$  be a symmetric and positive definite  $m \times m$  matrix and let  $\mathbf{B}$  be an  $r \times m$  matrix with full rank  $\text{rk}(\mathbf{B}) = r$ .
- Then also the matrices  $\mathbf{B}\mathbf{A}\mathbf{B}^\top$  and  $\mathbf{A}^{-1}$  are positive definite.

**Proof**

- Because of the full rank of  $\mathbf{B}^\top$  it holds that  $\mathbf{B}^\top \mathbf{x} \neq \mathbf{o}$  for each  $\mathbf{x} \in \mathbb{R}^r$  with  $\mathbf{x} \neq \mathbf{o}$ .
- Since  $\mathbf{A}$  is positive definite, it also holds that

$$\mathbf{x}^\top (\mathbf{B}\mathbf{A}\mathbf{B}^\top) \mathbf{x} = (\mathbf{B}^\top \mathbf{x})^\top \mathbf{A} (\mathbf{B}^\top \mathbf{x}) > 0$$

for each  $\mathbf{x} \in \mathbb{R}^r$  with  $\mathbf{x} \neq \mathbf{o}$ , i.e.,  $\mathbf{B}\mathbf{A}\mathbf{B}^\top$  is positive definite.

- Therefore, we get for  $\mathbf{B} = \mathbf{A}^{-1}$  that

$$\mathbf{A}^{-1} = \mathbf{A}^{-1} (\mathbf{A}\mathbf{A}^{-1}) = \mathbf{A}^{-1} \mathbf{A} (\mathbf{A}^{-1})^\top$$

is positive definite. □

## 1.2 Multivariate Normal Distribution

In this section we recall the notion of a multivariate normal distribution and discuss some fundamental properties of this family of distributions.

### 1.2.1 Definition and Fundamental Properties

- Let  $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$  be independent and (identically) normally distributed random variables, i.e.,

$$X_i \sim N(\mu, \sigma^2), \quad \forall i = 1, \dots, n, \quad (11)$$

where  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ .

- The assumption of normality in (11) and the independence of the sampling variables  $X_1, \dots, X_n$  mean in vector notation that the distribution of the random sample  $\mathbf{X} = (X_1, \dots, X_n)^\top$  is given by

$$\mathbf{X} \sim N(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_n), \quad (12)$$

where  $\boldsymbol{\mu} = (\mu, \dots, \mu)^\top$  and  $N(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_n)$  denotes the  $n$ -dimensional normal distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\sigma^2 \mathbf{I}_n$ .

- *Recall* (cf. Section WR-4.3.4): In general, the  $n$ -dimensional normal distribution is defined as follows.
  - Let  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top \in \mathbb{R}^n$  be an arbitrary vector and let  $\mathbf{K}$  be a symmetric and positive definite  $n \times n$ -matrix.

- Let  $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$  be an absolutely continuous random vector, where the joint density of  $\mathbf{Z}$  is given by

$$f(\mathbf{z}) = \left(\frac{1}{\sqrt{2\pi}}\right)^n \frac{1}{\sqrt{\det \mathbf{K}}} \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^\top \mathbf{K}^{-1}(\mathbf{z} - \boldsymbol{\mu})\right) \quad (13)$$

for each  $\mathbf{z} = (z_1, \dots, z_n)^\top \in \mathbb{R}^n$ .

- Then the random vector  $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$  is called (nondegenerate) *normally distributed*.
- Notation:  $\mathbf{Z} \sim \mathbf{N}(\boldsymbol{\mu}, \mathbf{K})$

Now, we show that the function given in (13) is an ( $n$ -dimensional) probability density.

**Theorem 1.1** *Let  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top \in \mathbb{R}^n$  be an arbitrary vector and let  $\mathbf{K}$  be a symmetric and positive definite  $n \times n$ -matrix. Then it holds that*

$$\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{K}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) dx_1 \dots dx_n = (2\pi)^{n/2} (\det \mathbf{K})^{1/2}. \quad (14)$$

### Proof

- Since  $\mathbf{K}$  is symmetric and positive definite (and therefore invertible), Lemma 1.5 implies that there exists an  $n \times n$  matrix  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$  consisting of the orthogonal eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of  $\mathbf{K}$ , such that

$$\mathbf{V}^\top \mathbf{K} \mathbf{V} = \boldsymbol{\Lambda}, \quad (15)$$

where  $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$  denotes the  $n \times n$  diagonal matrix which is built up of the eigenvalues  $\lambda_1, \dots, \lambda_n$  of  $\mathbf{K}$ .

- Since  $\mathbf{K}$  is positive definite, it holds that  $\mathbf{v}_i^\top \mathbf{K} \mathbf{v}_i = \lambda_i > 0$  for each  $i = 1, \dots, n$ , i.e., all eigenvalues  $\lambda_1, \dots, \lambda_n$  of  $\mathbf{K}$  are positive.
- Because of  $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$ , it holds  $\mathbf{V}^\top = \mathbf{V}^{-1}$  and  $\mathbf{V} \mathbf{V}^\top = \mathbf{I}$ , respectively.
- Due to the fact that  $(\mathbf{A} \mathbf{B})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$  and due to (15), it follows that

$$(\mathbf{V}^\top \mathbf{K} \mathbf{V})^{-1} = \mathbf{V}^\top \mathbf{K}^{-1} \mathbf{V} = \text{diag}(\lambda_1^{-1}, \dots, \lambda_n^{-1}).$$

- The mapping  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^n$  with  $\mathbf{y} = \varphi(\mathbf{x}) = \mathbf{V}^\top(\mathbf{x} - \boldsymbol{\mu})$ , i.e.,  $\mathbf{x} - \boldsymbol{\mu} = \mathbf{V} \mathbf{y}$ , maps  $\mathbb{R}^n$  bijectively onto itself and for the Jacobian determinant of  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^n$  it holds that

$$\det\left(\frac{\partial \varphi_i}{\partial x_j}(x_1, \dots, x_n)\right) = \det \mathbf{V} = \pm 1,$$

where the last equality follows from the fact that  $1 = \det(\mathbf{V}^\top \mathbf{V}) = (\det \mathbf{V})^2$ .

- Therefore, the integral on the left-hand side of (14) can be written as

$$\begin{aligned} & \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{K}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) dx_1 \dots dx_n \\ &= \int_{\mathbb{R}^n} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{K}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) d(x_1, \dots, x_n) = \int_{\mathbb{R}^n} \exp\left(-\frac{1}{2} \sum_{i=1}^n \frac{y_i^2}{\lambda_i}\right) d(y_1, \dots, y_n) \\ &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2} \sum_{i=1}^n \frac{y_i^2}{\lambda_i}\right) dy_1 \dots dy_n = \prod_{i=1}^n (2\pi \lambda_i)^{1/2}. \end{aligned}$$

- This implies (14) since

$$\prod_{i=1}^n \lambda_i = \det \boldsymbol{\Lambda} = \det(\mathbf{V}^\top \mathbf{K} \mathbf{V}) = \det(\mathbf{V}^\top \mathbf{V}) \det \mathbf{K} = \det \mathbf{K}. \quad \square$$

### 1.2.2 Characteristics of the Multivariate Normal Distribution

- Let  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top \in \mathbb{R}^n$  be an arbitrary vector and let  $\mathbf{K} = (k_{ij})$  be a symmetric and positive definite  $n \times n$  matrix.
- First, we determine the characteristic function of a normally distributed random vector.
- *Recall:* The *characteristic function*  $\varphi : \mathbb{R}^n \rightarrow \mathbb{C}$  of an arbitrary  $n$ -dimensional random vector  $\mathbf{X} = (X_1, \dots, X_n)^\top : \Omega \rightarrow \mathbb{R}^n$  is given by

$$\varphi(\mathbf{t}) = \mathbb{E} \exp(i \mathbf{t}^\top \mathbf{X}) = \mathbb{E} \exp\left(i \sum_{\ell=1}^n t_\ell X_\ell\right), \quad \forall \mathbf{t} = (t_1, \dots, t_n)^\top \in \mathbb{R}^n. \quad (16)$$

#### Theorem 1.2

- Let the random vector  $\mathbf{X} = (X_1, \dots, X_n)^\top : \Omega \rightarrow \mathbb{R}^n$  be normally distributed with  $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{K})$ .
- Then the characteristic function  $\varphi : \mathbb{R}^n \rightarrow \mathbb{C}$  of  $\mathbf{X}$  fulfills

$$\varphi(\mathbf{t}) = \exp\left(i \mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2} \mathbf{t}^\top \mathbf{K} \mathbf{t}\right), \quad \forall \mathbf{t} \in \mathbb{R}^n. \quad (17)$$

#### Proof

- Equations (13) and (16) imply

$$\begin{aligned} \varphi(\mathbf{t}) &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(i \sum_{\ell=1}^n t_\ell x_\ell\right) f(x_1, \dots, x_n) dx_1 \dots dx_n \\ &= \frac{1}{(2\pi)^{n/2} (\det \mathbf{K})^{1/2}} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(i \mathbf{t}^\top \mathbf{x} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{K}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) dx_1 \dots dx_n \\ &= \frac{\exp(i \mathbf{t}^\top \boldsymbol{\mu})}{(2\pi)^{n/2} (\det \mathbf{K})^{1/2}} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(i \mathbf{t}^\top \mathbf{y} - \frac{1}{2} \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}\right) dy_1 \dots dy_n, \end{aligned}$$

where the last equality holds due to the substitution  $\mathbf{y} = \mathbf{x} - \boldsymbol{\mu}$ , for which the matrix of the partial derivatives is the identity matrix and therefore the Jacobian determinant is equal to 1.

- Similar to the proof of Theorem 1.1 it follows with the help of the substitutions  $\mathbf{y} = \mathbf{V}\mathbf{x}$  and  $\mathbf{t} = \mathbf{V}\mathbf{s}$  that

$$\begin{aligned} \varphi(\mathbf{t}) &= \frac{\exp(i \mathbf{t}^\top \boldsymbol{\mu})}{(2\pi)^{n/2} (\det \mathbf{K})^{1/2}} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(i \mathbf{s}^\top \mathbf{x} - \frac{1}{2} \mathbf{x}^\top \mathbf{V}^\top \mathbf{K}^{-1} \mathbf{V} \mathbf{x}\right) dx_1 \dots dx_n \\ &= \frac{\exp(i \mathbf{t}^\top \boldsymbol{\mu})}{(2\pi)^{n/2} (\det \mathbf{K})^{1/2}} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(\sum_{\ell=1}^n \left(i s_\ell x_\ell - \frac{x_\ell^2}{2\lambda_\ell}\right)\right) dx_1 \dots dx_n \end{aligned}$$

and thus

$$\begin{aligned} \varphi(\mathbf{t}) &= \frac{\exp(i \mathbf{t}^\top \boldsymbol{\mu})}{(2\pi)^{n/2} (\det \mathbf{K})^{1/2}} \prod_{\ell=1}^n \int_{-\infty}^{\infty} \exp\left(i s_\ell x_\ell - \frac{x_\ell^2}{2\lambda_\ell}\right) dx_\ell \\ &= \exp(i \mathbf{t}^\top \boldsymbol{\mu}) \prod_{\ell=1}^n \frac{1}{\sqrt{2\pi\lambda_\ell}} \int_{-\infty}^{\infty} \exp\left(i s_\ell x_\ell - \frac{x_\ell^2}{2\lambda_\ell}\right) dx_\ell, \end{aligned}$$

where the matrix  $\mathbf{V}$  consists of the orthonormal eigenvectors of  $\mathbf{K}$  and  $\lambda_1, \dots, \lambda_n > 0$  are the eigenvalues of  $\mathbf{K}$  with  $\det \mathbf{K} = \lambda_1 \cdot \dots \cdot \lambda_n$ .

- Now, it is sufficient to consider that  $\varphi_\ell : \mathbb{R} \rightarrow \mathbb{C}$  with

$$\varphi_\ell(s) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\lambda_\ell}} \exp\left(i s x - \frac{x^2}{2\lambda_\ell}\right) dx$$

is the characteristic function of the (one dimensional)  $N(0, \lambda_\ell)$ -distribution.

- In Section WR-5.3.3 we already have seen that  $\varphi_\ell(s) = \exp(-\lambda_\ell s^2/2)$ .
- Hence, we get

$$\begin{aligned} \varphi(\mathbf{t}) &= \exp(i \mathbf{t}^\top \boldsymbol{\mu}) \prod_{\ell=1}^n \exp\left(-\frac{\lambda_\ell s_\ell^2}{2}\right) = \exp(i \mathbf{t}^\top \boldsymbol{\mu}) \exp\left(-\frac{\sum_{\ell=1}^n \lambda_\ell s_\ell^2}{2}\right) \\ &= \exp(i \mathbf{t}^\top \boldsymbol{\mu}) \exp\left(-\frac{\mathbf{t}^\top \mathbf{K} \mathbf{t}}{2}\right). \end{aligned} \quad \square$$

Using (17) for the characteristic function we are able to determine expectation and covariance matrix of a normally distributed random vector.

**Corollary 1.1** *If  $\mathbf{X} = (X_1, \dots, X_n)^\top \sim N(\boldsymbol{\mu}, \mathbf{K})$ , it holds for arbitrary  $i, j = 1, \dots, n$  that*

$$\mathbb{E} X_i = \mu_i, \quad \text{and} \quad \text{Cov}(X_i, X_j) = k_{ij}. \quad (18)$$

**Proof**

- From (17) it follows that

$$\frac{\partial \varphi(\mathbf{t})}{\partial t_i} = \left(i \mu_i - \sum_{\ell=1}^n k_{i\ell} t_\ell\right) \varphi(\mathbf{t}) \quad (19)$$

and

$$\frac{\partial^2 \varphi(\mathbf{t})}{\partial t_i \partial t_j} = -k_{ij} \varphi(\mathbf{t}) + \left(i \mu_i - \sum_{\ell=1}^n k_{i\ell} t_\ell\right) \left(i \mu_j - \sum_{\ell=1}^n k_{j\ell} t_\ell\right) \varphi(\mathbf{t}). \quad (20)$$

- It is easy to see that

$$\mathbb{E} X_i = i^{-1} \left. \frac{\partial \varphi(\mathbf{t})}{\partial t_i} \right|_{\mathbf{t}=\mathbf{o}}.$$

Because of  $\varphi(\mathbf{o}) = 1$  and (19), it follows that  $\mathbb{E} X_i = \mu_i$ .

- Furthermore,

$$\mathbb{E}(X_i X_j) = - \left. \frac{\partial^2 \varphi(\mathbf{t})}{\partial t_i \partial t_j} \right|_{\mathbf{t}=\mathbf{o}}.$$

This equation and (20) imply  $\text{Cov}(X_i, X_j) = k_{ij}$ . □

**Remark**

- In Theorem WR-4.14 we have shown that the covariance matrix  $\mathbf{K} = \mathbf{K}_{\mathbf{X}}$  of an arbitrary random vector  $\mathbf{X} = (X_1, \dots, X_n)^\top$  is always symmetric and positive semidefinite.
- In (13), where the density of the nondegenerate multivariate normal distribution is defined, it is additionally required that the covariance matrix  $\mathbf{K}$  is positive definite.
- Here,  $\mathbf{K}$  being positive definite is not only sufficient but also necessary to ensure that the matrix  $\mathbf{K}$  is invertible, i.e.,  $\det \mathbf{K} \neq 0$  or  $\mathbf{K}$  has full rank.

### 1.2.3 Marginal Distributions and Independence of Subvectors; Convolution Properties

- In this section it is shown how to derive further interesting properties of the multivariate normal distribution using Theorem 1.2.
- For this purpose we need a *vectorial version* of the uniqueness theorem for characteristic functions (cf. Corollary WR-5.5), which we will state without proof.

**Lemma 1.9** *Let  $\mathbf{X}, \mathbf{Y} : \Omega \rightarrow \mathbb{R}^n$  be arbitrary random vectors;  $\mathbf{X} = (X_1, \dots, X_n)^\top$ ,  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ . Then it holds that*

$$\mathbf{X} \stackrel{d}{=} \mathbf{Y} \quad \text{if and only if} \quad \varphi_{\mathbf{X}}(\mathbf{t}) = \varphi_{\mathbf{Y}}(\mathbf{t}) \quad \forall \mathbf{t} = (t_1, \dots, t_n)^\top \in \mathbb{R}^n, \quad (21)$$

where

$$\varphi_{\mathbf{X}}(\mathbf{t}) = \mathbb{E} \exp\left(i \sum_{j=1}^n t_j X_j\right) \quad \text{and} \quad \varphi_{\mathbf{Y}}(\mathbf{t}) = \mathbb{E} \exp\left(i \sum_{j=1}^n t_j Y_j\right)$$

are the characteristic functions of  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively.

First, we show that arbitrary subvectors of normally distributed random vectors are also normally distributed.

- We assume  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top \in \mathbb{R}^n$  to be an arbitrary vector and  $\mathbf{K} = (k_{ij})$  to be a symmetric and positive definite  $n \times n$ -matrix.
- It is obvious that the random vector  $(X_{\pi_1}, \dots, X_{\pi_n})^\top$  is normally distributed for each permutation  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_n)^\top$  of the natural numbers  $1, \dots, n$  if  $\mathbf{X} = (X_1, \dots, X_n)^\top$  is normally distributed.
- Therefore, we can w.l.o.g restrict the examination of the distribution of subvectors of normally distributed random vectors to the examination of the first components.

**Corollary 1.2** *Let  $\mathbf{X} = (X_1, \dots, X_n)^\top \sim N(\boldsymbol{\mu}, \mathbf{K})$ , where  $\mathbf{K}$  is positive definite. Then it holds that*

$$(X_1, \dots, X_m)^\top \sim N(\boldsymbol{\mu}_m, \mathbf{K}_m) \quad \forall m = 1, \dots, n,$$

where  $\boldsymbol{\mu}_m = (\mu_1, \dots, \mu_m)^\top$  and  $\mathbf{K}_m$  denotes the  $m \times m$  matrix which consists of the first  $m$  rows and columns of  $\mathbf{K}$ .

**Proof**

- Let  $\varphi : \mathbb{R}^n \rightarrow \mathbb{C}$  be the characteristic function of  $(X_1, \dots, X_n)^\top$ .
- Now, the characteristic function  $\varphi_m : \mathbb{R}^m \rightarrow \mathbb{C}$  of  $(X_1, \dots, X_m)^\top$  fulfills

$$\varphi_m(\mathbf{t}_m) = \varphi\left(\underbrace{(\mathbf{t}_m, 0, \dots, 0)}_{n-m}\right), \quad \forall \mathbf{t}_m = (t_1, \dots, t_m)^\top \in \mathbb{R}^m.$$

- This result and (17) imply that

$$\varphi_m(\mathbf{t}_m) = \exp\left(i \mathbf{t}_m^\top \boldsymbol{\mu}_m - \frac{1}{2} \mathbf{t}_m^\top \mathbf{K}_m \mathbf{t}_m\right), \quad \forall \mathbf{t}_m \in \mathbb{R}^m.$$

- Since  $\mathbf{K}$  is symmetric and positive definite, we know that also the  $m \times m$  matrix  $\mathbf{K}_m$  is symmetric and positive definite. From this fact and from Theorem 1.2 it follows that the characteristic function of the subvector  $(X_1, \dots, X_m)^\top$  is identical with the characteristic function of the  $N(\boldsymbol{\mu}_m, \mathbf{K}_m)$ -distribution.
- The statement follows because of the one-to-one correspondence of characteristic functions and distributions of random vectors (cf. Lemma 1.9).  $\square$

There is a simple criterion for two subvectors  $(X_1, \dots, X_m)^\top$  and  $(X_{m+1}, \dots, X_n)^\top$ , with  $1 \leq m < n$ , of the normally distributed random vector  $\mathbf{X} = (X_1, \dots, X_n)^\top$  being independent.

**Corollary 1.3** *Let  $\mathbf{X} = (X_1, \dots, X_n)^\top$  be a normally distributed random vector with  $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{K})$ ;  $\mathbf{K} = (k_{ij})$ . The subvectors  $(X_1, \dots, X_m)^\top$  and  $(X_{m+1}, \dots, X_n)^\top$  are independent if and only if  $k_{ij} = 0$  for arbitrary  $i \in \{1, \dots, m\}$  and  $j \in \{m+1, \dots, n\}$ .*

**Proof**

- If the subvectors  $(X_1, \dots, X_m)^\top$  and  $(X_{m+1}, \dots, X_n)^\top$  are independent, then the (one-dimensional) random variables  $X_i$  and  $X_j$  are independent for arbitrary  $i \in \{1, \dots, m\}$  and  $j \in \{m+1, \dots, n\}$ .
- Thus, it holds that  $\text{Cov}(X_i, X_j) = 0$  and Corollary 1.1 implies that  $k_{ij} = 0$ .
- Let us now assume that  $k_{ij} = 0$  for arbitrary  $i \in \{1, \dots, m\}$  and  $j \in \{m+1, \dots, n\}$ .
- Then Theorem 1.2 implies that the characteristic function  $\varphi(\mathbf{t})$  of  $\mathbf{X} = (X_1, \dots, X_n)^\top$  has the following factorization.
- For each  $\mathbf{t} = (t_1, \dots, t_n)^\top \in \mathbb{R}^n$  it holds that

$$\begin{aligned} \varphi(\mathbf{t}) &= \exp\left(\mathbf{i} \mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2} \mathbf{t}^\top \mathbf{K} \mathbf{t}\right) = \exp\left(\mathbf{i} \sum_{i=1}^n t_i \mu_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n t_i k_{ij} t_j\right) \\ &= \exp\left(\mathbf{i} \sum_{i=1}^m t_i \mu_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m t_i k_{ij} t_j\right) \exp\left(\mathbf{i} \sum_{i=m+1}^n t_i \mu_i - \frac{1}{2} \sum_{i=m+1}^n \sum_{j=m+1}^n t_i k_{ij} t_j\right), \end{aligned}$$

where the factors of the last term are the characteristic functions of  $(X_1, \dots, X_m)^\top$  and  $(X_{m+1}, \dots, X_n)^\top$ .

- The statement follows because of the one-to-one correspondence of characteristic functions and distributions of random vectors (cf. Lemma 1.9).  $\square$

**Remark**

- Finally, we show that the family of multivariate normal distributions is closed under convolution. In the following we call this property briefly "convolution stability" of the multivariate normal distribution. In Corollary WR-3.2 we already have proved the convolution stability of one-dimensional normal distributions.
- The following formula for the characteristic function of sums of independent random vectors is useful in this context. The proof is analog to the proof of the one-dimensional case (cf. Theorem WR-5.18).

**Lemma 1.10** *Let  $\mathbf{Z}_1, \mathbf{Z}_2 : \Omega \rightarrow \mathbb{R}^n$  be independent random vectors. The characteristic function  $\varphi_{\mathbf{Z}_1 + \mathbf{Z}_2} : \mathbb{R}^n \rightarrow \mathbb{C}$  of the sum  $\mathbf{Z}_1 + \mathbf{Z}_2$  can then be written as*

$$\varphi_{\mathbf{Z}_1 + \mathbf{Z}_2}(\mathbf{t}) = \varphi_{\mathbf{Z}_1}(\mathbf{t}) \varphi_{\mathbf{Z}_2}(\mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^n, \quad (22)$$

where  $\varphi_{\mathbf{Z}_i}$  denotes the characteristic function of  $\mathbf{Z}_i$ ;  $i = 1, 2$ .

The following statement is called *convolution stability* of the multivariate normal distribution.

**Corollary 1.4** *Let  $\mathbf{Z}_1, \mathbf{Z}_2 : \Omega \rightarrow \mathbb{R}^n$  be independent random vectors with  $\mathbf{Z}_i \sim N(\boldsymbol{\mu}_i, \mathbf{K}_i)$  for  $i = 1, 2$ . Then it holds that  $\mathbf{Z}_1 + \mathbf{Z}_2 \sim N(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2, \mathbf{K}_1 + \mathbf{K}_2)$ .*

**Proof**

- Equations (17) and (22) imply that

$$\begin{aligned}\varphi_{\mathbf{Z}_1+\mathbf{Z}_2}(\mathbf{t}) &= \varphi_{\mathbf{Z}_1}(\mathbf{t}) \varphi_{\mathbf{Z}_2}(\mathbf{t}) \\ &= \exp\left(i\mathbf{t}^\top \boldsymbol{\mu}_1 - \frac{1}{2}\mathbf{t}^\top \mathbf{K}_1 \mathbf{t}\right) \exp\left(i\mathbf{t}^\top \boldsymbol{\mu}_2 - \frac{1}{2}\mathbf{t}^\top \mathbf{K}_2 \mathbf{t}\right) \\ &= \exp\left(i\mathbf{t}^\top (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2) - \frac{1}{2}\mathbf{t}^\top (\mathbf{K}_1 + \mathbf{K}_2) \mathbf{t}\right).\end{aligned}$$

- This result and the uniqueness theorem for characteristic functions (cf. Lemma 1.9) imply the statement.  $\square$

**1.2.4 Linear Transformation of Normally Distributed Random Vectors**

Now, we show that the linear transformation of a normally distributed random vector again is a normally distributed random vector.

**Theorem 1.3**

- Let  $\mathbf{Y} \sim \mathbf{N}(\boldsymbol{\mu}, \mathbf{K})$  be an  $n$ -dimensional normally distributed random vector with mean vector  $\boldsymbol{\mu} \in \mathbb{R}^n$  and (positive definite) covariance matrix  $\mathbf{K}$ .
- Moreover, let  $m \leq n$ , let  $\mathbf{A}$  be an arbitrary  $m \times n$  matrix having full rank  $\text{rk}(\mathbf{A}) = m$  and let  $\mathbf{c} \in \mathbb{R}^m$  be an arbitrary  $m$ -dimensional vector.
- Then it holds that  $\mathbf{Z} = \mathbf{A}\mathbf{Y} + \mathbf{c}$  is an ( $m$ -dimensional) normally distributed random vector with

$$\mathbf{Z} \sim \mathbf{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{c}, \mathbf{A}\mathbf{K}\mathbf{A}^\top). \quad (23)$$

**Proof**

- For each  $\mathbf{a} \in \mathbb{R}^m$  it holds that

$$\varphi_{\mathbf{Z}}(\mathbf{t}) = \exp(i\mathbf{t}^\top \mathbf{a}) \varphi_{\mathbf{Z}-\mathbf{a}}(\mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^m.$$

- From (17) derived in Theorem 1.2 and from the uniqueness theorem for the characteristic function of normally distributed random vectors it follows that

$$\mathbf{Z} \sim \mathbf{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{c}, \mathbf{A}\mathbf{K}\mathbf{A}^\top) \quad \text{if and only if} \quad \mathbf{Z} - (\mathbf{A}\boldsymbol{\mu} + \mathbf{c}) \sim \mathbf{N}(\mathbf{o}, \mathbf{A}\mathbf{K}\mathbf{A}^\top).$$

- Therefore, we will w.l.o.g. assume that  $\mathbf{Y} \sim \mathbf{N}(\mathbf{o}, \mathbf{K})$  and  $\mathbf{c} = \mathbf{o}$ .
- Then the characteristic function  $\varphi_{\mathbf{Z}}(\mathbf{t})$  of  $\mathbf{Z} = \mathbf{A}\mathbf{Y}$  fulfills

$$\begin{aligned}\varphi_{\mathbf{Z}}(\mathbf{t}) &= \mathbb{E} e^{i\mathbf{t}^\top \mathbf{Z}} \\ &= \mathbb{E} e^{i\mathbf{t}^\top \mathbf{A}\mathbf{Y}} = \mathbb{E} e^{i(\mathbf{A}^\top \mathbf{t})^\top \mathbf{Y}} \\ &= \varphi_{\mathbf{Y}}(\mathbf{A}^\top \mathbf{t}),\end{aligned}$$

for each  $\mathbf{t} \in \mathbb{R}^m$ , where  $\varphi_{\mathbf{Y}}(\mathbf{A}^\top \mathbf{t})$  denotes the value of the characteristic function of the normally distributed random vector  $\mathbf{Y}$  at  $\mathbf{A}^\top \mathbf{t} \in \mathbb{R}^n$ .



- Now, formula (17) for the characteristic function of normally distributed random vectors implies

$$\begin{aligned}\varphi_{\mathbf{Z}}(\mathbf{t}) &= \varphi_{\mathbf{Y}}(\mathbf{A}^\top \mathbf{t}) \\ &= \exp\left(-\frac{1}{2}(\mathbf{A}^\top \mathbf{t})^\top \mathbf{K}(\mathbf{A}^\top \mathbf{t})\right) \\ &= \exp\left(-\frac{1}{2}\mathbf{t}^\top (\mathbf{A}\mathbf{K}\mathbf{A}^\top)\mathbf{t}\right).\end{aligned}$$

- In other words: The characteristic function of  $\mathbf{Z}$  is equal to the characteristic function of  $N(\mathbf{o}, \mathbf{A}\mathbf{K}\mathbf{A}^\top)$ .
- The uniqueness theorem for characteristic functions of random vectors implies  $\mathbf{Z} \sim N(\mathbf{o}, \mathbf{A}\mathbf{K}\mathbf{A}^\top)$ .  $\square$

By using Theorem 1.3 it follows in particular that it is possible to create normally distributed random vectors by a linear transformation of vectors whose components are independent  $N(0, 1)$ -distributed random variables.

### Corollary 1.5

- Let  $Y_1, \dots, Y_n : \Omega \rightarrow \mathbb{R}$  be independent random variables with  $Y_i \sim N(0, 1)$  for each  $i = 1, \dots, n$ , i.e.,  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top \sim N(\mathbf{o}, \mathbf{I})$ .
- Let  $\mathbf{K}$  be a symmetric and positive definite  $n \times n$  matrix and let  $\boldsymbol{\mu} \in \mathbb{R}^n$ .
- Then the random vector  $\mathbf{Z} = \mathbf{K}^{1/2}\mathbf{Y} + \boldsymbol{\mu}$  satisfies  $\mathbf{Z} \sim N(\boldsymbol{\mu}, \mathbf{K})$ , where  $\mathbf{K}^{1/2}$  is the square root of  $\mathbf{K}$ .

### Proof

- With the help of Theorem 1.3 it follows that

$$\mathbf{Z} \sim N(\boldsymbol{\mu}, \mathbf{K}^{1/2}(\mathbf{K}^{1/2})^\top).$$

- Now, this result and Lemma 1.6 imply the statement.  $\square$

### 1.2.5 Degenerate Multivariate Normal Distribution

In the following, we will give a generalization of the notation of (nondegenerate) multivariate normal distributions, which was introduced in Section 1.2.1.

- A factorization property of covariance matrices which has already been mentioned in Lemma 1.7 is useful in this context.
- *Recall:* Let  $\mathbf{K}$  be a symmetric and positive semidefinite  $n \times n$  matrix with  $\text{rk}(\mathbf{K}) = r \leq n$ . Then there is an  $n \times r$  matrix  $\mathbf{B}$  with  $\text{rk}(\mathbf{B}) = r$ , such that

$$\mathbf{K} = \mathbf{B}\mathbf{B}^\top. \tag{24}$$

### Definition

- Let  $\mathbf{Y}$  be an  $n$ -dimensional random vector with mean vector  $\boldsymbol{\mu} = \mathbb{E}\mathbf{Y}$  and covariance matrix  $\mathbf{K} = \text{Cov}(\mathbf{Y})$ , such that  $\text{rk}(\mathbf{K}) = r$  with  $r \leq n$ .
- Then  $\mathbf{Y}$  is called normally distributed if  $\mathbf{Y} \stackrel{d}{=} \boldsymbol{\mu} + \mathbf{B}\mathbf{Z}$ , where  $\mathbf{B}$  is an  $n \times r$  matrix with  $\text{rk}(\mathbf{B}) = r$  fulfilling (24) and where  $\mathbf{Z}$  is an  $r$ -dimensional random vector with  $\mathbf{Z} \sim N(\mathbf{o}, \mathbf{I}_r)$ .

- We say that  $\mathbf{Y} \sim \mathbf{N}(\boldsymbol{\mu}, \mathbf{K})$  follows a *degenerate normal distribution* if  $\text{rk}(\mathbf{K}) < n$ . (Notation:  $\mathbf{Y} \sim \mathbf{N}(\boldsymbol{\mu}, \mathbf{K})$ )

### Remark

- If  $\text{rk}(\mathbf{K}) = r < n$ , then the random vector  $\mathbf{Y} \sim \mathbf{N}(\boldsymbol{\mu}, \mathbf{K})$  is *not* absolutely continuous
  - because the values of  $\mathbf{Y} \stackrel{d}{=} \boldsymbol{\mu} + \mathbf{B}\mathbf{Z}$  are almost surely (with probability 1) elements of the  $r$ -dimensional subset  $\{\boldsymbol{\mu} + \mathbf{B}\mathbf{x} : \mathbf{x} \in \mathbb{R}^r\}$  of  $\mathbb{R}^n$ ,
  - i.e., the distribution of  $\mathbf{Y}$  has no density with respect to the  $n$ -dimensional Lebesgue measure.
  - An example for this is the random vector  $\mathbf{Y} = (Z, Z)^\top = \mathbf{B}\mathbf{Z}$  with  $Z \sim \mathbf{N}(0, \sigma^2)$  and  $\mathbf{B} = (1, 1)^\top$ , which only takes values on the diagonal  $\{(z_1, z_2) \in \mathbb{R}^2 : z_1 = z_2\}$ .
- The distribution of the random vector  $\boldsymbol{\mu} + \mathbf{B}\mathbf{Z}$  does *not* depend on the choice of matrix  $\mathbf{B}$  of the factorization (24).
- This is an immediate consequence of both of the following criteria for (degenerate and nondegenerate) multivariate normal distributions.

### Theorem 1.4

- Let  $\mathbf{Y}$  be an  $n$ -dimensional random vector with mean vector  $\boldsymbol{\mu} = \mathbb{E}\mathbf{Y}$  and covariance matrix  $\mathbf{K} = \text{Cov}(\mathbf{Y})$ , such that  $\text{rk}(\mathbf{K}) = r$  with  $r \leq n$ .
- The random vector  $\mathbf{Y}$  is normally distributed if and only if one of the following conditions is fulfilled:

1. The characteristic function  $\varphi(\mathbf{t}) = \mathbb{E} \exp\left(i \sum_{j=1}^n t_j Y_j\right)$  of  $\mathbf{Y}$  is given by

$$\varphi(\mathbf{t}) = \exp\left(i \mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2} \mathbf{t}^\top \mathbf{K} \mathbf{t}\right), \quad \forall \mathbf{t} = (t_1, \dots, t_n)^\top \in \mathbb{R}^n. \quad (25)$$

2. The linear function  $\mathbf{c}^\top \mathbf{Y}$  of  $\mathbf{Y}$  is normally distributed for each  $\mathbf{c} \in \mathbb{R}^n$  with  $\mathbf{c} \neq \mathbf{o}$  and

$$\mathbf{c}^\top \mathbf{Y} \sim \mathbf{N}(\mathbf{c}^\top \boldsymbol{\mu}, \mathbf{c}^\top \mathbf{K} \mathbf{c}).$$

The *proof* of Theorem 1.4 is omitted (and left as an exercise).

## 1.3 Linear and Quadratic Forms of Normally Distributed Random Vectors

### 1.3.1 Definition, Expectation and Covariance

#### Definition

- Let  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  and  $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$  be arbitrary  $n$ -dimensional random vectors and let  $\mathbf{A}$  be a symmetric  $n \times n$  matrix with real-valued entries.
- Then the (real-valued) random variable  $\mathbf{Y}^\top \mathbf{A} \mathbf{Y} : \Omega \rightarrow \mathbb{R}$  is called a *quadratic form* of  $\mathbf{Y}$  (with respect to  $\mathbf{A}$ ).
- The random variable  $\mathbf{Y}^\top \mathbf{A} \mathbf{Z} : \Omega \rightarrow \mathbb{R}$  is called a *bilinear form* of  $\mathbf{Y}$  and  $\mathbf{Z}$  (with respect to  $\mathbf{A}$ ).

First, we derive the expectation of quadratic or bilinear forms.

**Theorem 1.5** Let  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  and  $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$  be arbitrary  $n$ -dimensional random vectors and let  $\mathbf{A}$  be a symmetric  $n \times n$  matrix with real-valued entries. Furthermore, let the mean vectors  $\boldsymbol{\mu}_\mathbf{Y} = \mathbb{E}\mathbf{Y}$  and  $\boldsymbol{\mu}_\mathbf{Z} = \mathbb{E}\mathbf{Z}$  as well as the covariance matrices  $\mathbf{K}_{\mathbf{Y}\mathbf{Y}} = (\text{Cov}(Y_i, Y_j))$  and  $\mathbf{K}_{\mathbf{Z}\mathbf{Y}} = (\text{Cov}(Z_i, Y_j))$  be well-defined. Then it holds that

$$\mathbb{E}(\mathbf{Y}^\top \mathbf{A} \mathbf{Y}) = \text{tr}(\mathbf{A} \mathbf{K}_{\mathbf{Y}\mathbf{Y}}) + \boldsymbol{\mu}_\mathbf{Y}^\top \mathbf{A} \boldsymbol{\mu}_\mathbf{Y} \quad \text{and} \quad \mathbb{E}(\mathbf{Y}^\top \mathbf{A} \mathbf{Z}) = \text{tr}(\mathbf{A} \mathbf{K}_{\mathbf{Z}\mathbf{Y}}) + \boldsymbol{\mu}_\mathbf{Y}^\top \mathbf{A} \boldsymbol{\mu}_\mathbf{Z}. \quad (26)$$

**Proof**

- We only prove the second formula in (26) since the first formula follows as a special case for  $\mathbf{Z} = \mathbf{Y}$ .
- It obviously holds that  $\mathbf{Y}^\top \mathbf{A} \mathbf{Z} = \text{tr}(\mathbf{Y}^\top \mathbf{A} \mathbf{Z})$ . Moreover, from Lemma 1.1 it follows that  $\text{tr}(\mathbf{Y}^\top \mathbf{A} \mathbf{Z}) = \text{tr}(\mathbf{A} \mathbf{Z} \mathbf{Y}^\top)$ .
- Altogether we get

$$\begin{aligned} \mathbb{E}(\mathbf{Y}^\top \mathbf{A} \mathbf{Z}) &= \mathbb{E} \text{tr}(\mathbf{Y}^\top \mathbf{A} \mathbf{Z}) = \mathbb{E} \text{tr}(\mathbf{A} \mathbf{Z} \mathbf{Y}^\top) = \text{tr}(\mathbf{A} \mathbb{E}(\mathbf{Z} \mathbf{Y}^\top)) \\ &= \text{tr}(\mathbf{A}(\mathbf{K}_{\mathbf{Z}\mathbf{Y}} + \boldsymbol{\mu}_{\mathbf{Z}} \boldsymbol{\mu}_{\mathbf{Y}}^\top)) = \text{tr}(\mathbf{A} \mathbf{K}_{\mathbf{Z}\mathbf{Y}}) + \boldsymbol{\mu}_{\mathbf{Y}}^\top \mathbf{A} \boldsymbol{\mu}_{\mathbf{Z}}. \end{aligned} \quad \square$$

In a similar way it is possible to derive a formula for the covariance of quadratic forms of normally distributed random vectors. The following formulas for the third and fourth mixed moments of the components of centered normally distributed random vectors are useful in this context.

**Lemma 1.11** *Let  $\mathbf{Z} = (Z_1, \dots, Z_n)^\top \sim \mathcal{N}(\mathbf{o}, \mathbf{K})$  be a normally distributed random vector with mean vector  $\boldsymbol{\mu} = \mathbf{o}$  and with an arbitrary covariance matrix  $\mathbf{K} = (k_{ij})$ . Then it holds that*

$$\mathbb{E}(Z_i Z_j Z_\ell) = 0 \quad \text{and} \quad \mathbb{E}(Z_i Z_j Z_\ell Z_m) = k_{ij} k_{\ell m} + k_{i\ell} k_{jm} + k_{j\ell} k_{im} \quad \forall i, j, \ell, m \in \{1, \dots, n\}. \quad (27)$$

The *proof* of Lemma 1.11 is omitted. It is an immediate consequence of Theorems 1.2 and 1.4, cf. the proof of Corollary 1.1.

**Theorem 1.6**

- Let  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  be an  $n$ -dimensional random vector with  $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$  and let  $\mathbf{A} = (a_{ij})$ ,  $\mathbf{B} = (b_{ij})$  be arbitrary symmetric  $n \times n$  matrices.
- Then
 
$$\text{Cov}(\mathbf{Y}^\top \mathbf{A} \mathbf{Y}, \mathbf{Y}^\top \mathbf{B} \mathbf{Y}) = 2 \text{tr}(\mathbf{A} \mathbf{K} \mathbf{B} \mathbf{K}) + 4 \boldsymbol{\mu}^\top \mathbf{A} \mathbf{K} \mathbf{B} \boldsymbol{\mu}. \quad (28)$$
- In particular, it holds that  $\text{Var}(\mathbf{Y}^\top \mathbf{A} \mathbf{Y}) = 2 \text{tr}((\mathbf{A} \mathbf{K})^2) + 4 \boldsymbol{\mu}^\top \mathbf{A} \mathbf{K} \mathbf{A} \boldsymbol{\mu}$ .

**Proof**

- From the definition of the covariance and from Theorem 1.5 it follows that

$$\begin{aligned} \text{Cov}(\mathbf{Y}^\top \mathbf{A} \mathbf{Y}, \mathbf{Y}^\top \mathbf{B} \mathbf{Y}) &= \mathbb{E}((\mathbf{Y}^\top \mathbf{A} \mathbf{Y} - \mathbb{E}(\mathbf{Y}^\top \mathbf{A} \mathbf{Y}))(\mathbf{Y}^\top \mathbf{B} \mathbf{Y} - \mathbb{E}(\mathbf{Y}^\top \mathbf{B} \mathbf{Y}))) \\ &= \mathbb{E}((\mathbf{Y}^\top \mathbf{A} \mathbf{Y} - \text{tr}(\mathbf{A} \mathbf{K}) - \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu})(\mathbf{Y}^\top \mathbf{B} \mathbf{Y} - \text{tr}(\mathbf{B} \mathbf{K}) - \boldsymbol{\mu}^\top \mathbf{B} \boldsymbol{\mu})). \end{aligned}$$

- With the substitution  $\mathbf{Z} = \mathbf{Y} - \boldsymbol{\mu}$  or  $\mathbf{Y} = \mathbf{Z} + \boldsymbol{\mu}$  it follows that

$$\begin{aligned} \text{Cov}(\mathbf{Y}^\top \mathbf{A} \mathbf{Y}, \mathbf{Y}^\top \mathbf{B} \mathbf{Y}) &= \mathbb{E}((\mathbf{Z}^\top \mathbf{A} \mathbf{Z} + 2 \boldsymbol{\mu}^\top \mathbf{A} \mathbf{Z} - \text{tr}(\mathbf{A} \mathbf{K}))(\mathbf{Z}^\top \mathbf{B} \mathbf{Z} + 2 \boldsymbol{\mu}^\top \mathbf{B} \mathbf{Z} - \text{tr}(\mathbf{B} \mathbf{K}))) \\ &= \mathbb{E}(\mathbf{Z}^\top \mathbf{A} \mathbf{Z} \mathbf{Z}^\top \mathbf{B} \mathbf{Z}) + 2 \boldsymbol{\mu}^\top \mathbf{A} \mathbb{E}(\mathbf{Z} \mathbf{Z}^\top \mathbf{B} \mathbf{Z}) + 2 \boldsymbol{\mu}^\top \mathbf{B} \mathbb{E}(\mathbf{Z} \mathbf{Z}^\top \mathbf{A} \mathbf{Z}) \\ &\quad - \mathbb{E}(\mathbf{Z}^\top \mathbf{A} \mathbf{Z}) \text{tr}(\mathbf{B} \mathbf{K}) - \mathbb{E}(\mathbf{Z}^\top \mathbf{B} \mathbf{Z}) \text{tr}(\mathbf{A} \mathbf{K}) \\ &\quad + 4 \boldsymbol{\mu}^\top \mathbf{A} \mathbf{K} \mathbf{B} \boldsymbol{\mu} + \text{tr}(\mathbf{A} \mathbf{K}) \text{tr}(\mathbf{B} \mathbf{K}) \\ &= \mathbb{E}(\mathbf{Z}^\top \mathbf{A} \mathbf{Z} \mathbf{Z}^\top \mathbf{B} \mathbf{Z}) + 2 \boldsymbol{\mu}^\top \mathbf{A} \mathbb{E}(\mathbf{Z} \mathbf{Z}^\top \mathbf{B} \mathbf{Z}) + 2 \boldsymbol{\mu}^\top \mathbf{B} \mathbb{E}(\mathbf{Z} \mathbf{Z}^\top \mathbf{A} \mathbf{Z}) \\ &\quad + 4 \boldsymbol{\mu}^\top \mathbf{A} \mathbf{K} \mathbf{B} \boldsymbol{\mu} - \text{tr}(\mathbf{A} \mathbf{K}) \text{tr}(\mathbf{B} \mathbf{K}), \end{aligned}$$

where the last equality is a result of Theorem 1.5 because  $\mathbf{Z} \sim \mathcal{N}(\mathbf{o}, \mathbf{K})$ , which implies  $\mathbb{E}(\mathbf{Z}^\top \mathbf{A} \mathbf{Z}) = \text{tr}(\mathbf{A} \mathbf{K})$ .

- Since the matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{K}$  are symmetric, it follows from Lemma 1.11 that

$$\begin{aligned}
\mathbb{E}(\mathbf{Z}^\top \mathbf{A} \mathbf{Z} \mathbf{Z}^\top \mathbf{B} \mathbf{Z}) &= \mathbb{E}(\mathbf{Z}^\top \mathbf{A} \mathbf{Z} \cdot \mathbf{Z}^\top \mathbf{B} \mathbf{Z}) \\
&= \sum_{i=1}^n \sum_{j=1}^n \sum_{\ell=1}^n \sum_{m=1}^n a_{ij} b_{\ell m} \mathbb{E}(Z_i Z_j Z_\ell Z_m) \\
&= \sum_{i=1}^n \sum_{j=1}^n \sum_{\ell=1}^n \sum_{m=1}^n (a_{ij} k_{ji} b_{\ell m} k_{m\ell} + a_{ji} k_{i\ell} b_{\ell m} k_{mj} + a_{ij} k_{j\ell} b_{\ell m} k_{mi}) \\
&= \text{tr}(\mathbf{A} \mathbf{K}) \text{tr}(\mathbf{B} \mathbf{K}) + 2 \text{tr}(\mathbf{A} \mathbf{K} \mathbf{B} \mathbf{K}).
\end{aligned}$$

- Furthermore, Lemma 1.11 implies that

$$\mathbb{E}(\mathbf{Z} \mathbf{Z}^\top \mathbf{A} \mathbf{Z}) = \left( \sum_{i=1}^n \sum_{j=1}^n a_{ij} \mathbb{E}(Z_i Z_j Z_\ell) \right)_\ell = \mathbf{o} \quad (29)$$

and analogously  $\mathbb{E}(\mathbf{Z} \mathbf{Z}^\top \mathbf{B} \mathbf{Z}) = \mathbf{o}$ .

- This result and the above derived expression for  $\text{Cov}(\mathbf{Y}^\top \mathbf{A} \mathbf{Y}, \mathbf{Y}^\top \mathbf{B} \mathbf{Y})$  imply the statement.  $\square$

Now, we derive the following formula for the covariance vector of linear or quadratic forms of normally distributed random vectors.

**Theorem 1.7** *Let  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  be an  $n$ -dimensional random vector with  $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$  and let  $\mathbf{A} = (a_{ij})$ ,  $\mathbf{B} = (b_{ij})$  be arbitrary symmetric  $n \times n$  matrices. Then it holds*

$$\text{Cov}(\mathbf{A} \mathbf{Y}, \mathbf{Y}^\top \mathbf{B} \mathbf{Y}) = 2 \mathbf{A} \mathbf{K} \mathbf{B} \boldsymbol{\mu}. \quad (30)$$

### Proof

- As  $\mathbb{E}(\mathbf{A} \mathbf{Y}) = \mathbf{A} \boldsymbol{\mu}$  and as it has been shown in Theorem 1.5 that

$$\mathbb{E}(\mathbf{Y}^\top \mathbf{B} \mathbf{Y}) = \text{tr}(\mathbf{B} \mathbf{K}) + \boldsymbol{\mu}^\top \mathbf{B} \boldsymbol{\mu},$$

it follows that

$$\begin{aligned}
\text{Cov}(\mathbf{A} \mathbf{Y}, \mathbf{Y}^\top \mathbf{B} \mathbf{Y}) &= \mathbb{E}((\mathbf{A} \mathbf{Y} - \mathbf{A} \boldsymbol{\mu})(\mathbf{Y}^\top \mathbf{B} \mathbf{Y} - \boldsymbol{\mu}^\top \mathbf{B} \boldsymbol{\mu} - \text{tr}(\mathbf{B} \mathbf{K}))) \\
&= \mathbb{E}((\mathbf{A} \mathbf{Y} - \mathbf{A} \boldsymbol{\mu})((\mathbf{Y} - \boldsymbol{\mu})^\top \mathbf{B} (\mathbf{Y} - \boldsymbol{\mu}) + 2(\mathbf{Y} - \boldsymbol{\mu})^\top \mathbf{B} \boldsymbol{\mu} - \text{tr}(\mathbf{B} \mathbf{K}))).
\end{aligned}$$

- Moreover, it holds that  $\mathbb{E}(\mathbf{A} \mathbf{Y} - \mathbf{A} \boldsymbol{\mu}) = \mathbf{o}$  and from (29) it follows with  $\mathbf{Z} = \mathbf{Y} - \boldsymbol{\mu}$  that

$$\mathbb{E}((\mathbf{A} \mathbf{Y} - \mathbf{A} \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})^\top \mathbf{B} (\mathbf{Y} - \boldsymbol{\mu})) = \mathbf{A} \mathbb{E}((\mathbf{Y} - \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})^\top \mathbf{B} (\mathbf{Y} - \boldsymbol{\mu})) = \mathbf{o}.$$

- Therefore, we get

$$\begin{aligned}
\text{Cov}(\mathbf{A} \mathbf{Y}, \mathbf{Y}^\top \mathbf{B} \mathbf{Y}) &= 2 \mathbb{E}((\mathbf{A} \mathbf{Y} - \mathbf{A} \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})^\top \mathbf{B} \boldsymbol{\mu}) \\
&= 2 \mathbf{A} \mathbb{E}((\mathbf{Y} - \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})^\top) \mathbf{B} \boldsymbol{\mu} \\
&= 2 \mathbf{A} \mathbf{K} \mathbf{B} \boldsymbol{\mu}.
\end{aligned}$$

$\square$

### 1.3.2 Noncentral $\chi^2$ -Distribution

To determine the distribution of quadratic forms of normally distributed random vectors we introduce the (parametric) family of the noncentral  $\chi^2$ -distribution.

**Definition** Let  $\boldsymbol{\mu} \in \mathbb{R}^n$  and  $(X_1, \dots, X_n)^\top \sim N(\boldsymbol{\mu}, \mathbf{I})$ . Then the random variable

$$Z = (X_1, \dots, X_n)(X_1, \dots, X_n)^\top = \sum_{i=1}^n X_i^2$$

is distributed according to a *noncentral  $\chi^2$ -distribution* with  $n$  degrees of freedom and the *noncentrality parameter*  $\lambda = \boldsymbol{\mu}^\top \boldsymbol{\mu}$ . (Notation:  $Z \sim \chi_{n,\lambda}^2$ )

#### Remark

- For  $\boldsymbol{\mu} = \mathbf{o}$  we obtain the (central)  $\chi^2$ -distribution  $\chi_n^2$  with  $n$  degrees of freedom, which has already been introduced in Section I-1.3.1, as a special case.
- To derive a formula for the density of the noncentral  $\chi^2$ -distribution we consider (in addition to the characteristic function) still another *integral transform* of probability densities.

#### Definition

- Let  $f : \mathbb{R} \rightarrow [0, \infty)$  be the density of a real-valued random variable, such that the integral  $\int_{-\infty}^{\infty} e^{tx} f(x) dx$  is well-defined for each  $t \in (a, b)$  in a certain interval  $(a, b)$  with  $a < b$ .
- Then the mapping  $\psi : (a, b) \rightarrow \mathbb{R}$  with

$$\psi(t) = \int_{-\infty}^{\infty} e^{tx} f(x) dx, \quad \forall t \in (a, b) \quad (31)$$

is called the *moment generating function* of the density  $f$ .

The following *uniqueness theorem* for moment generating functions is true, which we state without proof.

#### Lemma 1.12

- Let  $f, f' : \mathbb{R} \rightarrow [0, \infty)$  be densities of real-valued random variables and let the corresponding moment generating functions  $\psi : (a, b) \rightarrow \mathbb{R}$  and  $\psi' : (a, b) \rightarrow \mathbb{R}$  be well-defined in a (common) interval  $(a, b)$  with  $a < b$ .
- It holds that  $\psi(t) = \psi'(t)$  for each  $t \in (a, b)$  if and only if  $f(x) = f'(x)$  for almost all  $x \in \mathbb{R}$ .

By using Lemma 1.12 we are now able to identify the density of the noncentral  $\chi^2$ -distribution.

#### Theorem 1.8

- Let the random variable  $Z_{n,\lambda} : \Omega \rightarrow \mathbb{R}$  be distributed according to the  $\chi_{n,\lambda}^2$ -distribution with  $n$  degrees of freedom and noncentrality parameter  $\lambda$ .
- Then the density of  $Z_{n,\lambda}$  is given by

$$f_{Z_{n,\lambda}}(z) = \begin{cases} \exp\left(-\frac{\lambda+z}{2}\right) \sum_{j=0}^{\infty} \frac{\left(\frac{\lambda}{2}\right)^j z^{\frac{n}{2}+j-1}}{j! 2^{\frac{n}{2}+j} \Gamma\left(\frac{n}{2}+j\right)}, & \text{if } z > 0, \\ 0 & \text{otherwise.} \end{cases} \quad (32)$$

**Proof**

- Let  $\boldsymbol{\mu} \in \mathbb{R}^n$  and  $(X_1, \dots, X_n)^\top \sim N(\boldsymbol{\mu}, \mathbf{I})$ .
- The moment generating function  $\psi_Z(t)$  of  $Z = (X_1, \dots, X_n)(X_1, \dots, X_n)^\top = \sum_{j=1}^n X_j^2$  is well-defined for  $t \in (-\infty, 1/2)$  and for each  $t < 1/2$  it holds that

$$\begin{aligned} \psi_Z(t) &= \mathbb{E} \exp\left(t \sum_{j=1}^n X_j^2\right) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(t \sum_{j=1}^n x_j^2\right) \prod_{j=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_j - \mu_j)^2\right) dx_1 \dots dx_n \\ &= \left(\frac{1}{2\pi}\right)^{n/2} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp\left(t \sum_{j=1}^n x_j^2 - \frac{1}{2} \sum_{j=1}^n (x_j - \mu_j)^2\right) dx_1 \dots dx_n \\ &= \prod_{j=1}^n \int_{-\infty}^{\infty} (2\pi)^{-1/2} \exp\left(tx_j^2 - \frac{1}{2}(x_j - \mu_j)^2\right) dx_j. \end{aligned}$$

- It is possible to rewrite the exponent of the last term as follows:

$$\begin{aligned} tx_j^2 - \frac{1}{2}(x_j - \mu_j)^2 &= -\frac{1}{2}(-2tx_j^2 + x_j^2 - 2x_j\mu_j + \mu_j^2) \\ &= -\frac{1}{2}\left(x_j^2(1-2t) - 2x_j\mu_j + \mu_j^2(1-2t)^{-1} + \mu_j^2 - \mu_j^2(1-2t)^{-1}\right) \\ &= -\frac{1}{2}\left((x_j - \mu_j(1-2t)^{-1})^2(1-2t) + \mu_j^2(1 - (1-2t)^{-1})\right). \end{aligned}$$

- Hence, it holds that

$$\begin{aligned} \psi_Z(t) &= \exp\left(-\frac{1}{2}(1 - (1-2t)^{-1}) \sum_{j=1}^n \mu_j^2\right) \prod_{j=1}^n \int_{-\infty}^{\infty} (2\pi)^{-1/2} \exp\left(-\frac{(x_j - \mu_j(1-2t)^{-1})^2}{2(1-2t)^{-1}}\right) dx_j \\ &= (1-2t)^{-n/2} \exp\left(-\frac{\lambda}{2}(1 - (1-2t)^{-1})\right) \end{aligned}$$

as the integrand represents the density of the one-dimensional normal distribution (except for the constant factor  $(1-2t)^{1/2}$ );  $\lambda = \boldsymbol{\mu}^\top \boldsymbol{\mu}$ .

- On the other hand, the moment generating function  $\psi(t)$  of the density  $f_{Z_{n,\lambda}}(z)$  given in (32) can be written as

$$\psi(t) = \sum_{j=0}^{\infty} \frac{e^{-\lambda/2}(\lambda/2)^j}{j!} \int_0^{\infty} e^{tz} \frac{z^{n/2+j-1} e^{-z/2}}{2^{\frac{n}{2}+j} \Gamma\left(\frac{n}{2}+j\right)} dz,$$

where the integral is the moment generating function of the (central)  $\chi^2$ -distribution  $\chi_{n+2j}^2$  with  $n+2j$  degrees of freedom.

- Similar to the way the characteristic function (cf. Theorem I-1.5) is defined, the moment generating function of this distribution is given by

$$\psi_{\chi_{n+2j}^2}(t) = \frac{1}{(1-2t)^{n/2+j}}.$$

- Therefore, it holds that

$$\int_0^{\infty} e^{tz} \frac{z^{n/2+j-1} e^{-z/2}}{2^{\frac{n}{2}+j} \Gamma\left(\frac{n}{2}+j\right)} dz = \frac{1}{(1-2t)^{n/2+j}},$$

and

$$\begin{aligned}\psi(t) &= e^{-\lambda/2}(1-2t)^{-n/2} \sum_{j=0}^{\infty} \frac{1}{j!} \left(\frac{\lambda}{2}(1-2t)^{-1}\right)^j \\ &= (1-2t)^{-n/2} \exp\left(-\frac{\lambda}{2}(1-(1-2t)^{-1})\right).\end{aligned}$$

- Hence,  $\psi(t) = \psi_Z(t)$  for each  $t < 1/2$  and the statement follows from Lemma 1.12.  $\square$

### 1.3.3 Distributional Properties of Linear and Quadratic Forms

- *Recall:* The definition of the noncentral  $\chi^2$ -distribution in Section 1.3.2 considers the sum of squares of the components of  $N(\boldsymbol{\mu}, \mathbf{I})$ -distributed random vectors.
- One can show that the (adequately modified) sum of squares is distributed according to the noncentral  $\chi^2$ -distribution even if the considered normally distributed random vector has an *arbitrary* positive definite covariance matrix.
- Indeed, let  $\boldsymbol{\mu} \in \mathbb{R}^n$  and let  $\mathbf{K}$  be a symmetric and positive definite  $n \times n$  matrix.
- If  $\mathbf{Z} = (Z_1, \dots, Z_n)^\top \sim N(\boldsymbol{\mu}, \mathbf{K})$ , Theorem 1.3 implies that

$$\mathbf{K}^{-1/2}\mathbf{Z} \sim N(\mathbf{K}^{-1/2}\boldsymbol{\mu}, \mathbf{I}).$$

- Therefore, by the definition of the noncentral  $\chi^2$ -distribution it follows that

$$\mathbf{Z}^\top \mathbf{K}^{-1} \mathbf{Z} = (\mathbf{K}^{-1/2}\mathbf{Z})^\top \mathbf{K}^{-1/2}\mathbf{Z} \sim \chi_{n,\lambda}^2, \quad (33)$$

where  $\lambda = (\mathbf{K}^{-1/2}\boldsymbol{\mu})^\top \mathbf{K}^{-1/2}\boldsymbol{\mu} = \boldsymbol{\mu}^\top \mathbf{K}^{-1}\boldsymbol{\mu}$ .

The distributional property (33) for quadratic forms of normally distributed random vectors has the following generalization. In this context Lemma 1.7 about the factorization of symmetric and positive semidefinite matrices is useful.

#### Theorem 1.9

- Let  $\mathbf{Z} = (Z_1, \dots, Z_n)^\top \sim N(\boldsymbol{\mu}, \mathbf{K})$ , where the covariance matrix  $\mathbf{K}$  be positive definite. Moreover, let  $\mathbf{A}$  be a symmetric  $n \times n$  matrix with  $\text{rk}(\mathbf{A}) = r \leq n$ .
- If the matrix  $\mathbf{AK}$  is idempotent, i.e., if  $\mathbf{AK} = (\mathbf{AK})^2$ , it holds that  $\mathbf{Z}^\top \mathbf{AZ} \sim \chi_{r,\lambda}^2$ , where  $\lambda = \boldsymbol{\mu}^\top \mathbf{A}\boldsymbol{\mu}$ .

#### Proof

- Let the matrix  $\mathbf{AK}$  be idempotent. Then it holds that

$$\mathbf{AK} = \mathbf{AKAK}.$$

- Since  $\mathbf{K}$  is nondegenerate, it is allowed to multiply both sides of the above equation from the right by  $\mathbf{K}^{-1}$ . It follows

$$\mathbf{A} = \mathbf{AKA} \quad (34)$$

or

$$\mathbf{x}^\top \mathbf{Ax} = \mathbf{x}^\top \mathbf{AKAx} = (\mathbf{Ax})^\top \mathbf{K}(\mathbf{Ax}) \geq 0$$

for each  $\mathbf{x} \in \mathbb{R}^n$ , i.e.,  $\mathbf{A}$  is positive semidefinite.

- According to Lemma 1.7 there exists a decomposition

$$\mathbf{A} = \mathbf{H}\mathbf{H}^\top, \quad (35)$$

such that the  $n \times r$  matrix  $\mathbf{H}$  has full column rank  $r$ .

- Now, Lemma 1.2 implies that the inverse matrix  $(\mathbf{H}^\top\mathbf{H})^{-1}$  exists.
- From Theorem 1.3 about the linear transformation of normally distributed random vectors it follows for the  $r$ -dimensional vector  $\mathbf{Z}' = \mathbf{H}^\top\mathbf{Z}$  that

$$\mathbf{Z}' \sim N(\mathbf{H}^\top\boldsymbol{\mu}, \mathbf{I}_r) \quad (36)$$

because

$$\begin{aligned} \mathbf{H}^\top\mathbf{K}\mathbf{H} &= (\mathbf{H}^\top\mathbf{H})^{-1}(\mathbf{H}^\top\mathbf{H})(\mathbf{H}^\top\mathbf{K}\mathbf{H})(\mathbf{H}^\top\mathbf{H})(\mathbf{H}^\top\mathbf{H})^{-1} \\ &= (\mathbf{H}^\top\mathbf{H})^{-1}\mathbf{H}^\top(\mathbf{A}\mathbf{K}\mathbf{A})\mathbf{H}(\mathbf{H}^\top\mathbf{H})^{-1} \\ &= (\mathbf{H}^\top\mathbf{H})^{-1}\mathbf{H}^\top\mathbf{A}\mathbf{H}(\mathbf{H}^\top\mathbf{H})^{-1} = \mathbf{I}_r, \end{aligned}$$

where the last three equalities follow from (34) and (35).

- As on the other hand

$$\mathbf{Z}^\top\mathbf{A}\mathbf{Z} = \mathbf{Z}^\top\mathbf{H}\mathbf{H}^\top\mathbf{Z} = (\mathbf{H}^\top\mathbf{Z})^\top\mathbf{H}^\top\mathbf{Z} = (\mathbf{Z}')^\top\mathbf{Z}'$$

and since

$$(\mathbf{H}^\top\boldsymbol{\mu})^\top\mathbf{H}^\top\boldsymbol{\mu} = \boldsymbol{\mu}^\top\mathbf{H}\mathbf{H}^\top\boldsymbol{\mu} = \boldsymbol{\mu}^\top\mathbf{A}\boldsymbol{\mu},$$

the statement follows from (36) and from the definition of the noncentral  $\chi^2$ -distribution.  $\square$

Furthermore, the following criterion for the independence of linear and quadratic forms of normally distributed random vectors is useful. It can be considered as a (vectorial) generalization of Lemma I-5.3.

### Theorem 1.10

- Let  $\mathbf{Z} = (Z_1, \dots, Z_n)^\top \sim N(\boldsymbol{\mu}, \mathbf{K})$ , where  $\mathbf{K}$  is an arbitrary (symmetric and positive semidefinite) covariance matrix.
- Moreover, let  $\mathbf{A}$ ,  $\mathbf{B}$  be arbitrary  $r_1 \times n$  and  $r_2 \times n$  matrices with  $r_1, r_2 \leq n$  and let  $\mathbf{C}$  be a symmetric and positive semidefinite  $n \times n$  matrix.
- If the additional condition

$$\mathbf{A}\mathbf{K}\mathbf{B}^\top = \mathbf{0} \quad \text{or} \quad \mathbf{A}\mathbf{K}\mathbf{C} = \mathbf{0} \quad (37)$$

is fulfilled, the random variables  $\mathbf{A}\mathbf{Z}$  and  $\mathbf{B}\mathbf{Z}$  or  $\mathbf{A}\mathbf{Z}$  and  $\mathbf{Z}^\top\mathbf{C}\mathbf{Z}$ , respectively, are independent.

### Proof

- First, we show that (37) implies the independence of the linear forms  $\mathbf{A}\mathbf{Z}$  and  $\mathbf{B}\mathbf{Z}$ .
- Because of the uniqueness theorem for characteristic functions of random vectors (cf. Lemma 1.9), it suffices to show that  $\mathbf{t}_2 \in \mathbb{R}^{r_2}$

$$\mathbb{E} \exp(i(\mathbf{t}_1^\top\mathbf{A}\mathbf{Z} + \mathbf{t}_2^\top\mathbf{B}\mathbf{Z})) = \mathbb{E} \exp(i\mathbf{t}_1^\top\mathbf{A}\mathbf{Z})\mathbb{E} \exp(i\mathbf{t}_2^\top\mathbf{B}\mathbf{Z})$$

for arbitrary  $\mathbf{t}_1 \in \mathbb{R}^{r_1}$ .



- From (37) it follows that

$$\mathbf{BKA}^\top = \left( (\mathbf{BKA}^\top)^\top \right)^\top = \left( \mathbf{AKB}^\top \right)^\top = \mathbf{0}.$$

- Therefore, it holds for arbitrary  $\mathbf{t}_1 \in \mathbb{R}^{r_1}$ ,  $\mathbf{t}_2 \in \mathbb{R}^{r_2}$  that

$$(\mathbf{t}_1^\top \mathbf{A})\mathbf{K}(\mathbf{t}_2^\top \mathbf{B})^\top = \mathbf{t}_1^\top \mathbf{AKB}^\top \mathbf{t}_2 = \mathbf{0}, \quad (\mathbf{t}_2^\top \mathbf{B})\mathbf{K}(\mathbf{t}_1^\top \mathbf{A})^\top = \mathbf{t}_2^\top \mathbf{BKA}^\top \mathbf{t}_1 = \mathbf{0}. \quad (38)$$

- Then the representation formula (25) for the characteristic function of normally distributed random vectors derived in Theorem 1.4 and (38) imply that

$$\begin{aligned} \mathbb{E} \exp(i(\mathbf{t}_1^\top \mathbf{AZ} + \mathbf{t}_2^\top \mathbf{BZ})) &= \mathbb{E} \exp(i(\mathbf{t}_1^\top \mathbf{A} + \mathbf{t}_2^\top \mathbf{B})\mathbf{Z}) \\ &= \exp\left(i(\mathbf{t}_1^\top \mathbf{A} + \mathbf{t}_2^\top \mathbf{B})\boldsymbol{\mu} - \frac{1}{2}(\mathbf{t}_1^\top \mathbf{A} + \mathbf{t}_2^\top \mathbf{B})\mathbf{K}(\mathbf{t}_1^\top \mathbf{A} + \mathbf{t}_2^\top \mathbf{B})^\top\right) \\ &= \exp\left(i(\mathbf{t}_1^\top \mathbf{A} + \mathbf{t}_2^\top \mathbf{B})\boldsymbol{\mu} - \frac{1}{2}(\mathbf{t}_1^\top \mathbf{A})\mathbf{K}(\mathbf{t}_1^\top \mathbf{A})^\top - \frac{1}{2}(\mathbf{t}_2^\top \mathbf{B})\mathbf{K}(\mathbf{t}_2^\top \mathbf{B})^\top\right) \\ &= \exp\left(i(\mathbf{t}_1^\top \mathbf{A})\boldsymbol{\mu} - \frac{1}{2}(\mathbf{t}_1^\top \mathbf{A})\mathbf{K}(\mathbf{t}_1^\top \mathbf{A})^\top\right) \exp\left(i(\mathbf{t}_2^\top \mathbf{B})\boldsymbol{\mu} - \frac{1}{2}(\mathbf{t}_2^\top \mathbf{B})\mathbf{K}(\mathbf{t}_2^\top \mathbf{B})^\top\right) \\ &= \mathbb{E} \exp(i\mathbf{t}_1^\top \mathbf{AZ}) \mathbb{E} \exp(i\mathbf{t}_2^\top \mathbf{BZ}). \end{aligned}$$

- Now, it remains to show that the independence of  $\mathbf{AZ}$  of  $\mathbf{Z}^\top \mathbf{CZ}$  is a result of the second condition of (37).
- Let  $\text{rk}(\mathbf{C}) = r \leq n$ . According to Lemma 1.7 there is an  $n \times r$  matrix  $\mathbf{H}$  with  $\text{rk}(\mathbf{H}) = r$ , such that  $\mathbf{C} = \mathbf{HH}^\top$ .
- Then it follows from (37) that  $\mathbf{AKHH}^\top = \mathbf{0}$  and  $\mathbf{AKHH}^\top \mathbf{H} = \mathbf{0}$ .
- Because of Lemma 1.2, the  $r \times r$  matrix  $\mathbf{H}^\top \mathbf{H}$  has (full) rank  $\text{rk}(\mathbf{H}) = r$ . Hence,  $\mathbf{H}^\top \mathbf{H}$  is invertible.
- Finally, it follows that  $\mathbf{AKH} = \mathbf{0}$ .
- Therefore, the first part of the proof implies that the linear forms  $\mathbf{AZ}$  and  $\mathbf{H}^\top \mathbf{Z}$  are independent.
- Because of

$$\mathbf{Z}^\top \mathbf{CZ} = \mathbf{Z}^\top \mathbf{HH}^\top \mathbf{Z} = (\mathbf{H}^\top \mathbf{Z})^\top \mathbf{H}^\top \mathbf{Z},$$

the transformation theorem for independent random vectors (cf. Theorem I-1.8) implies that also  $\mathbf{AZ}$  and  $\mathbf{Z}^\top \mathbf{CZ}$  are independent.  $\square$

## 2 Linear Models; Design Matrix with Full Rank

Recall (cf. Chapter 5 of the lecture “Stochastik I”):

- In *simple linear regression* one considers two datasets  $(x_1, \dots, x_n) \in \mathbb{R}^n$  and  $(y_1, \dots, y_n) \in \mathbb{R}^n$ , which shall be modeled stochastically.
- In doing so, we perceive the vectors  $(x_1, y_1), \dots, (x_n, y_n)$  as realizations of  $n$  random vectors, say  $(X_1, Y_1), \dots, (X_n, Y_n)$ , which are typically *not* identically distributed.
- We interpret the random variables  $Y_1, \dots, Y_n$  as *response variables* and assume that they depend on the *predictor variables*  $X_1, \dots, X_n$  in the following way:

$$Y_i = \varphi(X_i) + \varepsilon_i, \quad \forall i = 1, \dots, n, \quad (1)$$

where

- $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  is an arbitrary Borel measurable function, the so-called *regression function*, and
- $\varepsilon_1, \dots, \varepsilon_n : \Omega \rightarrow \mathbb{R}$  are random variables, so-called *error terms*, which can be used to model random errors, e.g., errors in measurement.
- We are dealing with an important special case if the regression function  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  is a linear function, the so-called *regression line*, i.e., if there are real numbers  $\beta_1, \beta_2 \in \mathbb{R}$ , such that

$$\varphi(x) = \beta_1 + \beta_2 x, \quad \forall x \in \mathbb{R}, \quad (2)$$

where  $\beta_1$  is called the *intercept* and  $\beta_2$  is called the *regression coefficient*.

- The quantities  $\beta_1, \beta_2 \in \mathbb{R}$  are unknown model parameters, which are to be estimated from the observed data  $(x_1, \dots, x_n) \in \mathbb{R}^n$  and  $(y_1, \dots, y_n) \in \mathbb{R}^n$ .

We now consider the following *multivariate* generalization of the simple linear regression model and let  $m, n \geq 2$  be arbitrary natural numbers, such that  $m \leq n$ .

- We assume that the response variables  $Y_1, \dots, Y_n$  depend on *vectorial*  $m$ -dimensional predictor variables  $(X_{11}, \dots, X_{1m})^\top, \dots, (X_{n1}, \dots, X_{nm})^\top$ , i.e.,

$$Y_i = \varphi(X_{i1}, \dots, X_{im}) + \varepsilon_i, \quad \forall i = 1, \dots, n, \quad (3)$$

where

- the regression function  $\varphi : \mathbb{R}^m \rightarrow \mathbb{R}$  is given by

$$\varphi(x_1, \dots, x_m) = \beta_1 x_1 + \dots + \beta_m x_m, \quad \forall (x_1, \dots, x_m)^\top \in \mathbb{R}^m \quad (4)$$

with (unknown) regression coefficients  $\beta_1, \dots, \beta_m \in \mathbb{R}$  and

- the random error terms  $\varepsilon_1, \dots, \varepsilon_n : \Omega \rightarrow \mathbb{R}$  satisfy the following requirements:

$$\mathbb{E} \varepsilon_i = 0, \quad \text{Var} \varepsilon_i = \sigma^2, \quad \text{Cov}(\varepsilon_i, \varepsilon_j) = 0, \quad \forall i, j = 1, \dots, n \text{ with } i \neq j \quad (5)$$

for a certain (unknown)  $\sigma^2 > 0$ .

- Here we just consider the case of deterministic predictor variables  $(X_{11}, \dots, X_{1m})^\top, \dots, (X_{n1}, \dots, X_{nm})^\top$ , i.e., we put

$$(X_{11}, \dots, X_{1m})^\top = (x_{11}, \dots, x_{1m})^\top, \dots, (X_{n1}, \dots, X_{nm})^\top = (x_{n1}, \dots, x_{nm})^\top$$

for certain vectors  $(x_{11}, \dots, x_{1m})^\top, \dots, (x_{n1}, \dots, x_{nm})^\top \in \mathbb{R}^m$ .

**Remark**

- In matrix notation the model given in (3) and (4) can be expressed as follows:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (6)$$

where

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} x_{11} & \dots & x_{1m} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nm} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}. \quad (7)$$

- Here  $\mathbf{X}$  is called the *design matrix* of the regression model.

**2.1 Method of Least Squares**

The goal of this section consists in estimating the unknown model parameters  $\beta_1, \dots, \beta_m$  and  $\sigma^2$  from the observed data  $(x_{11}, \dots, x_{1m})^\top, \dots, (x_{n1}, \dots, x_{nm})^\top \in \mathbb{R}^m$  and  $(y_1, \dots, y_n)^\top \in \mathbb{R}^n$ .

- Similar to the way this is done in Section I-5.1 we consider the *method of least squares* in order to determine estimators  $\hat{\beta}_1, \dots, \hat{\beta}_m$  for the unknown regression coefficients  $\beta_1, \dots, \beta_m$ .
- In detail, this means that a random vector  $\hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \dots, \hat{\beta}_m)^\top$  is to be determined, such that the *mean squared error*

$$e(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n (Y_i - (\beta_1 x_{i1} + \dots + \beta_m x_{im}))^2 \quad (8)$$

is minimal for  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$ .

**Remark** Except the model assumptions made in (5) no further preconditions concerning the distribution of the random error terms  $\varepsilon_1, \dots, \varepsilon_n : \Omega \rightarrow \mathbb{R}$  are required up to now.

**2.1.1 Normal Equation**

It can easily be shown that the function  $e(\boldsymbol{\beta})$  considered in (8) has a uniquely determined minimum if the design matrix  $\mathbf{X}$  has full (column) rank, i.e.,  $\text{rk}(\mathbf{X}) = m$ .

**Theorem 2.1** Let  $\text{rk}(\mathbf{X}) = m$ .

- The mean squared error  $e(\boldsymbol{\beta})$  in (8) is minimal if and only if  $\boldsymbol{\beta}$  is a solution of the so-called normal equation:

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{Y}. \quad (9)$$

- The normal equation (9) has the uniquely determined solution

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}. \quad (10)$$

**Proof**

- The function  $e(\boldsymbol{\beta})$  given in (8) is differentiable, where

$$e'(\boldsymbol{\beta}) = \left( \frac{\partial e(\boldsymbol{\beta})}{\partial \beta_1}, \dots, \frac{\partial e(\boldsymbol{\beta})}{\partial \beta_m} \right)^\top = \frac{2}{n} \left( \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} - \mathbf{X}^\top \mathbf{Y} \right)$$

and

$$e''(\boldsymbol{\beta}) = \left( \frac{\partial^2 e(\boldsymbol{\beta})}{\partial \beta_i \partial \beta_j} \right) = \frac{2}{n} \mathbf{X}^\top \mathbf{X}.$$

- Setting  $e'(\boldsymbol{\beta}) = \mathbf{0}$  results in the normal equation (9).
- Moreover, it follows from Lemma 1.2 that  $\text{rk}(\mathbf{X}^\top \mathbf{X}) = m$ .
  - Hence, the  $m \times m$  matrix  $\mathbf{X}^\top \mathbf{X}$  (and consequently also  $e''(\boldsymbol{\beta})$ ) is invertible and positive definite.
  - Therefore,  $e(\boldsymbol{\beta})$  is minimal if and only if  $\boldsymbol{\beta}$  is a solution of (9).
- As the  $m \times m$  matrix  $\mathbf{X}^\top \mathbf{X}$  is invertible, we know that (9) has a uniquely determined solution  $\widehat{\boldsymbol{\beta}}$ , which is given by (10).  $\square$

**Remark** The estimator  $\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$  for  $\boldsymbol{\beta}$  is a linear transformation of the random sample  $\mathbf{Y}$ , i.e.,  $\widehat{\boldsymbol{\beta}}$  is a *linear estimator* for  $\boldsymbol{\beta}$ .

**Examples** (simple and multiple linear regression model)

- For  $m = 2$  and

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad (11)$$

we obtain the *simple linear regression model* already considered in Section I-5.1 as a special case.

- The design matrix  $\mathbf{X}$  in (11) has full rank  $\text{rk}(\mathbf{X}) = 2$  if and only if not all  $x_1, \dots, x_n$  are equal.
- The estimators  $\widehat{\boldsymbol{\beta}} = (\widehat{\beta}_1, \widehat{\beta}_2)$  for the intercept  $\beta_1$  and the regression coefficient  $\beta_2$ , considered in (10) (see also Theorem I-5.1), are then given by

$$\widehat{\beta}_2 = \frac{s_{xy}^2}{s_{xx}^2} \quad \text{and} \quad \widehat{\beta}_1 = \bar{y}_n - \widehat{\beta}_2 \bar{x}_n, \quad (12)$$

respectively, where  $\bar{x}_n, \bar{y}_n$  denote the sample means, i.e.,

$$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y}_n = \frac{1}{n} \sum_{i=1}^n y_i,$$

and the sample variances  $s_{xx}^2, s_{yy}^2$  and sample covariance  $s_{xy}^2$  are given by

$$s_{xx}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2, \quad s_{xy}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n) \quad \text{and} \quad s_{yy}^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y}_n)^2.$$

- For  $m > 2$  and

$$\mathbf{X} = \begin{pmatrix} 1 & x_{12} & \dots & x_{1m} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n2} & \dots & x_{nm} \end{pmatrix} \quad (13)$$

we obtain the so-called *multiple linear regression model*.

### 2.1.2 Properties of the LS–Estimator $\widehat{\boldsymbol{\beta}}$

From now on in Section 2.1, we always assume that the design matrix  $\mathbf{X}$  has full (column) rank and derive three different *properties* of the LS–estimator  $\widehat{\boldsymbol{\beta}} = (\widehat{\beta}_1, \dots, \widehat{\beta}_m)^\top$  given in (10).

**Theorem 2.2** *The estimator  $\widehat{\boldsymbol{\beta}}$  is unbiased for  $\boldsymbol{\beta}$ , i.e.,  $\mathbb{E}\widehat{\boldsymbol{\beta}} = \boldsymbol{\beta}$  for all  $\boldsymbol{\beta} \in \mathbb{R}^m$ .*

**Proof**

Due to  $\mathbb{E}\boldsymbol{\varepsilon} = \mathbf{o}$  it follows from (6) and (10) that

$$\begin{aligned} \mathbb{E}\widehat{\boldsymbol{\beta}} &\stackrel{(10)}{=} \mathbb{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}) \stackrel{(6)}{=} \mathbb{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})) = \boldsymbol{\beta} + \mathbb{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}) \\ &= \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}\boldsymbol{\varepsilon} = \boldsymbol{\beta}. \end{aligned} \quad \square$$

The LS–estimator  $\widehat{\boldsymbol{\beta}}$  additionally has the following *minimum variance* property. We denote by  $\mathcal{L}$  the family of all unbiased linear estimators  $\widetilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y} + \mathbf{a}$  for  $\boldsymbol{\beta}$ , where  $\mathbf{A}$  is an  $(m \times n)$ –dimensional matrix and  $\mathbf{a} = (a_1, \dots, a_m)^\top \in \mathbb{R}^m$ .

**Theorem 2.3** *For all  $\widetilde{\boldsymbol{\beta}} = (\widetilde{\beta}_1, \dots, \widetilde{\beta}_m) \in \mathcal{L}$  it holds that*

$$\text{Var } \widehat{\beta}_i \leq \text{Var } \widetilde{\beta}_i, \quad \forall i = 1, \dots, m, \quad (14)$$

where the equality in (14) is true for all  $i = 1, \dots, m$  if and only if  $\widetilde{\boldsymbol{\beta}} = \widehat{\boldsymbol{\beta}}$ .

**Proof**

- As it is assumed that the linear estimator  $\widetilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y} + \mathbf{a}$  is unbiased for  $\boldsymbol{\beta}$ , one has

$$\boldsymbol{\beta} = \mathbb{E}\widetilde{\boldsymbol{\beta}} = \mathbb{E}(\mathbf{A}\mathbf{Y} + \mathbf{a}) \stackrel{(6)}{=} \mathbb{E}(\mathbf{A}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})) + \mathbf{a} = \mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbf{A}\mathbb{E}\boldsymbol{\varepsilon} + \mathbf{a} = \mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbf{a}$$

for all  $\boldsymbol{\beta} \in \mathbb{R}^m$ , where the last equality arises from  $\mathbb{E}\boldsymbol{\varepsilon} = \mathbf{o}$ .

- Herefrom, it follows that

$$\mathbf{A}\mathbf{X} = \mathbf{I} \quad \text{and} \quad \mathbf{a} = \mathbf{o}. \quad (15)$$

- Hence, one has

$$\widetilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y} = \mathbf{A}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbf{A}\boldsymbol{\varepsilon} = \boldsymbol{\beta} + \mathbf{A}\boldsymbol{\varepsilon},$$

i.e., each linear unbiased estimator  $\widetilde{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$  is of the form

$$\widetilde{\boldsymbol{\beta}} = \boldsymbol{\beta} + \mathbf{A}\boldsymbol{\varepsilon}. \quad (16)$$

- For the covariance matrix  $\text{Cov}(\widetilde{\boldsymbol{\beta}})$  of the random vector  $\widetilde{\boldsymbol{\beta}}$  it thus holds that

$$\text{Cov}(\widetilde{\boldsymbol{\beta}}) = \mathbb{E}((\widetilde{\boldsymbol{\beta}} - \boldsymbol{\beta})(\widetilde{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top) = \mathbb{E}((\mathbf{A}\boldsymbol{\varepsilon})(\mathbf{A}\boldsymbol{\varepsilon})^\top) = \mathbf{A}\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top)\mathbf{A}^\top = \sigma^2 \mathbf{A}\mathbf{A}^\top,$$

i.e.,

$$\text{Cov}(\widetilde{\boldsymbol{\beta}}) = \sigma^2 \mathbf{A}\mathbf{A}^\top. \quad (17)$$

- Furthermore, it results from (17) with  $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  that the covariance matrix  $\text{Cov}(\widehat{\boldsymbol{\beta}})$  of the LS-estimator  $\widehat{\boldsymbol{\beta}}$  is given by

$$\text{Cov}(\widehat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \quad (18)$$

because

$$\begin{aligned} \text{Cov}(\widehat{\boldsymbol{\beta}}) &= \sigma^2 ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}. \end{aligned}$$

- So in order to prove the validity of (14), it has to be shown that

$$((\mathbf{X}^\top \mathbf{X})^{-1})_{ii} \leq (\mathbf{A} \mathbf{A}^\top)_{ii}, \quad \forall i = 1, \dots, m. \quad (19)$$

- For  $\mathbf{D} = \mathbf{A} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  one has

$$\begin{aligned} \mathbf{A} \mathbf{A}^\top &= (\mathbf{D} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) (\mathbf{D} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top \\ &= \mathbf{D} \mathbf{D}^\top + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{D}^\top + \mathbf{D} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} + (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \mathbf{D} \mathbf{D}^\top + (\mathbf{X}^\top \mathbf{X})^{-1} \end{aligned}$$

because due to (15) it holds that

$$\mathbf{D} \mathbf{X} = (\mathbf{A} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{X} = \mathbf{A} \mathbf{X} - \mathbf{I} = \mathbf{I} - \mathbf{I} = \mathbf{0},$$

where  $\mathbf{0}$  denotes the zero matrix.

- As for  $\mathbf{D} = (d_{ij})$  the inequality  $(\mathbf{D} \mathbf{D}^\top)_{ii} = \sum_{j=1}^m d_{ij}^2 \geq 0$  is fulfilled, this gives (19).
- Moreover, it becomes clear that the equality in (19) for each  $i = 1, \dots, m$  holds if and only if  $\mathbf{D} = \mathbf{0}$ , i.e.,  $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ .  $\square$

### Remark

- It follows from Theorems 2.1 and 2.2 that  $\widehat{\boldsymbol{\beta}} \in \mathcal{L}$ . Moreover, it arises from Theorem 2.3 that  $\widehat{\boldsymbol{\beta}}$  is the *best unbiased linear estimator* for  $\boldsymbol{\beta}$  in terms of (14).
- We now derive a sufficient condition for  $\widehat{\boldsymbol{\beta}}$  to be a weakly consistent estimator for  $\boldsymbol{\beta}$ , where the sample size  $n$ , i.e., the number of rows of the design matrix  $\mathbf{X} = \mathbf{X}_n$  tends to  $\infty$ .
- *Recall:* An estimator  $\widetilde{\boldsymbol{\beta}}_n = \widetilde{\boldsymbol{\beta}}(Y_1, \dots, Y_n)$  for  $\boldsymbol{\beta}$  is called *weakly consistent* if

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\boldsymbol{\beta}}(|\widetilde{\boldsymbol{\beta}}_n - \boldsymbol{\beta}| > \varepsilon) = 0, \quad \forall \varepsilon > 0, \boldsymbol{\beta} \in \mathbb{R}^m.$$

- Under similar conditions, it can also be shown that  $\widehat{\boldsymbol{\beta}}_n$  is asymptotically normally distributed if  $n \rightarrow \infty$  (cf. Section III.3.2 in Pruscha (2000)).

**Theorem 2.4** *Let  $f : \mathbb{N} \rightarrow \mathbb{R} \setminus \{0\}$  be a function satisfying  $\lim_{n \rightarrow \infty} f(n) = 0$ , such that the limit*

$$\mathbf{Q} = \lim_{n \rightarrow \infty} (f(n) \mathbf{X}_n^\top \mathbf{X}_n) \quad (20)$$

*exists and the  $m \times m$  matrix  $\mathbf{Q}$  is invertible. Then  $\widehat{\boldsymbol{\beta}}_n$  is a weakly consistent estimator for  $\boldsymbol{\beta}$ .*

**Proof**

- As  $\widehat{\boldsymbol{\beta}}_n$  is unbiased (cf. Theorem 2.2), it holds for each  $n \geq m$  that

$$\begin{aligned} \mathbb{P}_{\boldsymbol{\beta}}(|\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}| > \varepsilon) &= \mathbb{P}_{\boldsymbol{\beta}}(|\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}|^2 > \varepsilon^2) = \mathbb{P}_{\boldsymbol{\beta}}\left(\sum_{i=1}^m (\widehat{\beta}_{in} - \beta_i)^2 > \varepsilon^2\right) \\ &\leq \mathbb{P}_{\boldsymbol{\beta}}\left(\bigcup_{i=1}^m \{(\widehat{\beta}_{in} - \beta_i)^2 > \frac{\varepsilon^2}{m}\}\right) \leq \sum_{i=1}^m \mathbb{P}_{\boldsymbol{\beta}}\left((\widehat{\beta}_{in} - \beta_i)^2 > \frac{\varepsilon^2}{m}\right) \\ &\leq \frac{m}{\varepsilon^2} \sum_{i=1}^m \text{Var } \widehat{\beta}_{in}, \end{aligned}$$

where the last inequality results from the Chebyshev inequality (cf. Theorem WR-4.18).

- Hence, it suffices to show that

$$\lim_{n \rightarrow \infty} \text{Var } \widehat{\beta}_{in} = 0, \quad \forall i = 1, \dots, m. \quad (21)$$

- The matrix  $\mathbf{Q}^{-1}$  is well-defined because we assume that the (limiting) matrix  $\mathbf{Q}$  is invertible. Moreover, it follows from (20) that

$$\mathbf{Q}^{-1} = \lim_{n \rightarrow \infty} (f(n) \mathbf{X}_n^{\top} \mathbf{X}_n)^{-1}.$$

- From the formula for the covariance matrix of the random vector  $\widehat{\boldsymbol{\beta}}_n$  derived in (18), it now results that

$$\lim_{n \rightarrow \infty} \text{Cov}(\widehat{\boldsymbol{\beta}}_n) = \sigma^2 \lim_{n \rightarrow \infty} (\mathbf{X}_n^{\top} \mathbf{X}_n)^{-1} = \sigma^2 \lim_{n \rightarrow \infty} f(n) \lim_{n \rightarrow \infty} (f(n) \mathbf{X}_n^{\top} \mathbf{X}_n)^{-1} = \left(\sigma^2 \lim_{n \rightarrow \infty} f(n)\right) \mathbf{Q}^{-1} = \mathbf{0}.$$

- This particularly implies (21). □

**2.1.3 Unbiased Estimation of the Variance  $\sigma^2$  of the Error Terms**

- Besides the conditions on the error terms  $\varepsilon_1, \dots, \varepsilon_n$  formulated in (5), we now assume that  $n > m$ . Furthermore, we again assume that the design matrix has full rank, i.e.,  $\text{rk}(\mathbf{X}) = m$ .
- By generalizing the approach we considered in Section I-5.1.3 for the estimation of  $\sigma^2$  in the simple linear regression model, we now consider

$$S^2 = \frac{1}{n - m} (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})^{\top} (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}). \quad (22)$$

- For normally distributed error terms,  $S^2$  can be regarded as a modified version of a maximum-likelihood estimator for  $\sigma^2$ ; cf. Section 2.2.

We show that the random variable  $S^2$  defined in (22) gives an unbiased estimator for  $\sigma^2$ . Here, the following lemmata are useful.

**Lemma 2.1** *The  $n \times n$  matrix*

$$\mathbf{G} = \mathbf{I} - \mathbf{X}(\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \quad (23)$$

*is idempotent and symmetric, i.e.,*

$$\mathbf{G} = \mathbf{G}^2 \quad \text{and} \quad \mathbf{G} = \mathbf{G}^{\top}. \quad (24)$$

**Proof**

- The second part of the statement in (24) follows directly from the definition of  $\mathbf{G}$  and the computation rules for transposed matrices because

$$\mathbf{G}^\top = \left( \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right)^\top = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{G}.$$

- Furthermore, it holds that

$$\begin{aligned} \mathbf{G}^2 &= \left( \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) \left( \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) \\ &= \mathbf{I} - 2\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \\ &= \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{G}. \end{aligned}$$

□

**Lemma 2.2** For the  $n \times n$  matrix  $\mathbf{G}$  given in (23) it holds that  $\text{tr}(\mathbf{G}) = n - m$ .

**Proof**

- One can easily see (cf. Lemmas 1.1 and 1.3) that
  - $\text{tr}(\mathbf{A} - \mathbf{B}) = \text{tr}(\mathbf{A}) - \text{tr}(\mathbf{B})$  for arbitrary  $n \times n$  matrices  $\mathbf{A}$  and  $\mathbf{B}$ ,
  - $\text{tr}(\mathbf{CD}) = \text{tr}(\mathbf{DC})$  for arbitrary  $n \times m$  matrices  $\mathbf{C}$  and arbitrary  $m \times n$  matrices  $\mathbf{D}$ .
- Herefrom and from the definition of  $\mathbf{G}$  in (23) it follows that

$$\begin{aligned} \text{tr}(\mathbf{G}) &= \text{tr}\left(\mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top\right) = \text{tr}(\mathbf{I}_n) - \text{tr}\left(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top\right) \\ &= \text{tr}(\mathbf{I}_n) - \text{tr}\left(\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}\right) = \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{I}_m) = n - m, \end{aligned}$$

where  $\mathbf{I}_\ell$  denotes the  $(\ell \times \ell)$ -dimensional identity matrix. □

**Theorem 2.5** It holds that  $\mathbb{E}S^2 = \sigma^2$  for any  $\sigma^2 > 0$ , i.e.,  $S^2$  is an unbiased estimator for  $\sigma^2$ .

**Proof**

- It obviously holds that

$$\mathbf{GX} = (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{X} = \mathbf{0}. \quad (25)$$

- Herefrom, it follows with the aid of (10) and (23) that

$$\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} \stackrel{(10)}{=} \mathbf{Y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \stackrel{(23)}{=} \mathbf{GY} = \mathbf{GX}\boldsymbol{\beta} + \mathbf{G}\boldsymbol{\varepsilon} \stackrel{(25)}{=} \mathbf{G}\boldsymbol{\varepsilon}.$$

- Hence, for the estimator  $S^2$  introduced in (22) it is true that

$$\begin{aligned} S^2 &= \frac{1}{n-m} (\mathbf{G}\boldsymbol{\varepsilon})^\top (\mathbf{G}\boldsymbol{\varepsilon}) = \frac{1}{n-m} \boldsymbol{\varepsilon}^\top \mathbf{G}^\top \mathbf{G}\boldsymbol{\varepsilon} = \frac{1}{n-m} \boldsymbol{\varepsilon}^\top \mathbf{G}\boldsymbol{\varepsilon} \\ &= \frac{1}{n-m} \text{tr}(\boldsymbol{\varepsilon}^\top \mathbf{G}\boldsymbol{\varepsilon}) = \frac{1}{n-m} \text{tr}(\mathbf{G}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top), \end{aligned}$$

on account of  $\mathbf{G}^\top \mathbf{G} = \mathbf{G}^2 = \mathbf{G}$  (cf. Lemma 2.1).

- Due to  $\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top) = \sigma^2 \mathbf{I}_n$  this leads to

$$\mathbb{E}S^2 = \frac{1}{n-m} \text{tr}(\mathbf{G}\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top)) = \frac{1}{n-m} \text{tr}(\mathbf{G}\sigma^2 \mathbf{I}_n) = \frac{\sigma^2}{n-m} \text{tr}(\mathbf{G}) = \sigma^2,$$

where the last equality results from Lemma 2.2. □



## 2.2 Normally Distributed Error Terms

- In addition to the model assumptions that were made at the beginning of Chapter 2, we now assume that the random error terms  $\varepsilon_1, \dots, \varepsilon_n : \Omega \rightarrow \mathbb{R}$  are independent and normally distributed, i.e.,  $\varepsilon_i \sim N(0, \sigma^2)$  for each  $i = 1, \dots, n$ .
- Moreover, let  $\text{rk}(\mathbf{X}) = m$  and  $n > m$ .
- According to Theorem 1.3, the distributions of the vector  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  of response variables and of the LS-estimator  $\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$  are given by

$$\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}) \quad (26)$$

and

$$\widehat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}), \quad (27)$$

respectively.

### 2.2.1 Maximum-Likelihood Estimation

- A *parametric model* for the distribution of the vector  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  of the sampling variables  $Y_1, \dots, Y_n$  is given by (26).
- This means that aside from the method of least squares, which was discussed in Section 2.1, we can now also use maximum-likelihood estimation in order to construct estimators for the unknown model parameters  $\boldsymbol{\beta}$  and  $\sigma^2$ .
- It follows from (1.13) and (26) that

$$f_{\mathbf{Y}}(\mathbf{y}) = \left( \frac{1}{\sigma\sqrt{2\pi}} \right)^n \exp\left( -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) \quad (28)$$

for each  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$ .

- Hence, we consider the *likelihood function*

$$L(\mathbf{y}; \boldsymbol{\beta}, \sigma^2) = \left( \frac{1}{\sigma\sqrt{2\pi}} \right)^n \exp\left( -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) \quad (29)$$

or the *loglikelihood function*

$$\log L(\mathbf{y}; \boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} |\mathbf{y} - \mathbf{X}\boldsymbol{\beta}|^2. \quad (30)$$

- We want to find estimators  $\widehat{\boldsymbol{\beta}}, \widehat{\sigma}^2$  for  $\boldsymbol{\beta}, \sigma^2$ , such that

$$L(\mathbf{Y}; \widehat{\boldsymbol{\beta}}, \widehat{\sigma}^2) = \sup_{\boldsymbol{\beta} \in \mathbb{R}^m, \sigma^2 > 0} L(\mathbf{Y}; \boldsymbol{\beta}, \sigma^2) \quad (31)$$

with probability 1 or equivalently such that

$$\log L(\mathbf{Y}; \widehat{\boldsymbol{\beta}}, \widehat{\sigma}^2) = \sup_{\boldsymbol{\beta} \in \mathbb{R}^m, \sigma^2 > 0} \log L(\mathbf{Y}; \boldsymbol{\beta}, \sigma^2) \quad (32)$$

with probability 1.

**Remark** The maximization in (31) or (32) can be carried out in two steps: first with respect to  $\boldsymbol{\beta}$  and then with respect to  $\sigma^2$ . Due to (30), the first step is identical with the minimization method considered in Section 2.1.1.

**Theorem 2.6** *The solutions of the maximization problems (31) and (32) are uniquely determined and given by*

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \quad (33)$$

and

$$\widehat{\sigma}^2 = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}), \quad (34)$$

respectively.

**Proof**

- For arbitrary but fixed  $\mathbf{y} \in \mathbb{R}^n$  and  $\sigma^2 > 0$ , we first consider the mapping

$$\mathbb{R}^m \ni \boldsymbol{\beta} \mapsto \log L(\mathbf{y}; \boldsymbol{\beta}, \sigma^2). \quad (35)$$

- In Theorem 2.1 we have shown that the mapping given in (35) has the uniquely determined global maximum  $\widehat{\boldsymbol{\beta}}(\mathbf{y}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ , which does not depend on  $\sigma^2$ .
- For each (fixed)  $\mathbf{y} \in \mathbb{R}^n$ , we now consider the mapping

$$(0, \infty) \ni \sigma^2 \mapsto \log L(\mathbf{y}; \widehat{\boldsymbol{\beta}}(\mathbf{y}), \sigma^2). \quad (36)$$

- This mapping is continuous and it obviously holds that

$$\lim_{\sigma^2 \rightarrow \infty} \log L(\mathbf{y}; \widehat{\boldsymbol{\beta}}(\mathbf{y}), \sigma^2) = -\infty.$$

- As  $n > m$  is assumed, the  $n$ -dimensional absolutely continuous random vector  $\mathbf{Y}$  only takes values in the  $m$ -dimensional subset  $\{\mathbf{X}\mathbf{z} : \mathbf{z} \in \mathbb{R}^m\}$  of  $\mathbb{R}^n$  with probability 0.
- Therefore, we have that  $|\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}|^2 > 0$  with probability 1 and

$$\lim_{\sigma^2 \rightarrow 0} \log L(\mathbf{y}; \widehat{\boldsymbol{\beta}}(\mathbf{y}), \sigma^2) = -\infty$$

for almost every  $\mathbf{y} \in \mathbb{R}^n$ .

- Thus, for almost every  $\mathbf{y} \in \mathbb{R}^n$ , the mapping given in (36) has at least one global maximum in  $(0, \infty)$ .
- For each of these maxima, it holds that

$$\frac{\partial \log L(\mathbf{y}; \widehat{\boldsymbol{\beta}}(\mathbf{y}), \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}(\mathbf{y}))^\top (\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}(\mathbf{y})) = 0.$$

- The (uniquely determined) solution of this equation is

$$\widehat{\sigma}^2(\mathbf{y}) = \frac{1}{n} (\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}(\mathbf{y}))^\top (\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}(\mathbf{y})). \quad \square$$

**Remark**

- The ML-estimator  $\widehat{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$  derived in Theorem 2.6 coincides with the LS-estimator derived in Theorem 2.1.
- In contrast, the ML-estimator  $\widehat{\sigma}^2$  for  $\sigma^2$  differs from the (unbiased) estimator  $S^2$  for  $\sigma^2$  considered in Section 2.1.3 in a constant proportionality factor because

$$\widehat{\sigma}^2 = \frac{n-m}{n} S^2.$$

### 2.2.2 Distributional Properties of $\widehat{\boldsymbol{\beta}}$ and $S^2$

- Apart from the fact that  $\widehat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$  is normally distributed, which was already mentioned in (27), it is also possible to determine the distribution of the estimator

$$S^2 = \frac{1}{n-m} (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}). \quad (37)$$

for the variance  $\sigma^2$  of the error terms.

- For this purpose we use the representation formula

$$\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}} = \mathbf{G}\boldsymbol{\varepsilon}, \quad (38)$$

which we have shown in the proof of Theorem 2.5, where  $\mathbf{G} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top$ .

From the condition derived in Theorem 1.9, under which a quadratic form of normally distributed random vectors follows a  $\chi^2$ -distribution, we obtain the following result.

**Theorem 2.7** *It holds that*

$$\frac{(n-m)S^2}{\sigma^2} \sim \chi_{n-m}^2, \quad (39)$$

*i. e., the random variable  $(n-m)S^2/\sigma^2$  is distributed according to the (central)  $\chi^2$ -distribution with  $n-m$  degrees of freedom.*

**Proof**

- In Lemma 2.1 we have shown that the matrix  $\mathbf{G} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top$  is idempotent and symmetric.
- Herefrom and from (38) it follows that

$$\begin{aligned} \frac{(n-m)S^2}{\sigma^2} &= \frac{1}{\sigma^2} (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}) = \frac{1}{\sigma^2} (\mathbf{G}\boldsymbol{\varepsilon})^\top \mathbf{G}\boldsymbol{\varepsilon} = \frac{1}{\sigma^2} \boldsymbol{\varepsilon}^\top \mathbf{G}^\top \mathbf{G}\boldsymbol{\varepsilon} \\ &= (\sigma^{-1}\boldsymbol{\varepsilon})^\top \mathbf{G}(\sigma^{-1}\boldsymbol{\varepsilon}). \end{aligned}$$

- As  $\sigma^{-1}\boldsymbol{\varepsilon} \sim N(\mathbf{o}, \mathbf{I})$  and as the matrix  $\mathbf{G}\mathbf{I} = \mathbf{G}$  is idempotent, it suffices to show that  $\text{rk}(\mathbf{G}) = n-m$  due to Theorem 1.9.
- This is a result of Lemma 1.3 and 2.2 because

$$\text{rk}(\mathbf{G}) \stackrel{\text{Lemma 1.3}}{=} \text{sp}(\mathbf{G}) \stackrel{\text{Lemma 2.2}}{=} n-m. \quad \square$$

Moreover, we use the criterion for the independence of linear and quadratic forms of normally distributed random vectors, which has been derived in Theorem 1.10 in order to show the following result.

**Theorem 2.8** *The estimators  $\widehat{\boldsymbol{\beta}}$  and  $S^2$  for  $\boldsymbol{\beta}$  and  $\sigma^2$ , respectively, are independent.*

**Proof**

- From  $\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$  and  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  it follows that

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon} + \boldsymbol{\beta}.$$

- Furthermore, we have shown in the proof of Theorem 2.7 that the estimator

$$S^2 = \frac{1}{n-m} (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})$$

can be written as a quadratic form of  $\boldsymbol{\varepsilon}$ :

$$S^2 = \frac{1}{n-m} \boldsymbol{\varepsilon}^\top \mathbf{G} \boldsymbol{\varepsilon}, \quad \text{where } \mathbf{G} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top.$$

- Due to  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$  and

$$\left( (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) \left( \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) = \mathbf{0},$$

it follows from Theorem 1.10 that the linear form  $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}$  and the quadratic form  $\boldsymbol{\varepsilon}^\top \mathbf{G} \boldsymbol{\varepsilon}$  are independent.

- This implies that also the random variables  $\widehat{\boldsymbol{\beta}}$  and  $S^2$  are independent. □

**2.2.3 Statistical Tests for the Regression Coefficients**

- By use of the distributional properties of linear and quadratic forms of normally distributed random vectors which have been derived in Sections 1.3.3 and 2.2.2, we are able to construct statistical t-tests and F-tests for the verification of hypotheses about the regression coefficients  $\beta_1, \dots, \beta_m$ .
- In doing so, we still consider the (independent) estimators  $\widehat{\boldsymbol{\beta}}$  and  $S^2$  for  $\boldsymbol{\beta}$  and  $\sigma^2$ , where

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}) \quad (40)$$

and

$$\frac{(n-m)S^2}{\sigma^2} = \frac{1}{\sigma^2} (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}) \sim \chi_{n-m}^2. \quad (41)$$

We first discuss the following F-test, which is also called a *test of model significance*.

- Here, the null hypothesis  $H_0 : \beta_1 = \dots = \beta_m = 0$  is tested (against the alternative  $H_1 : \beta_j \neq 0$  for at least one  $j \in \{1, \dots, m\}$ ).
- The choice of the test statistic is motivated by the following decomposition.

**Theorem 2.9** *With the notation  $\widehat{\mathbf{Y}} = \mathbf{X}\widehat{\boldsymbol{\beta}}$ , it holds that*

$$\mathbf{Y}^\top \mathbf{Y} = \widehat{\mathbf{Y}}^\top \widehat{\mathbf{Y}} + (\mathbf{Y} - \widehat{\mathbf{Y}})^\top (\mathbf{Y} - \widehat{\mathbf{Y}}). \quad (42)$$

**Proof**

- We have that

$$\begin{aligned} \mathbf{Y}^\top \mathbf{Y} &= \sum_{i=1}^n Y_i^2 = \sum_{i=1}^n \left( (Y_i - \hat{Y}_i) + \hat{Y}_i \right)^2 = \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + 2 \underbrace{\sum_{i=1}^n (Y_i - \hat{Y}_i) \hat{Y}_i}_{=0} \\ &= \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \hat{\mathbf{Y}}^\top \hat{\mathbf{Y}} + (\mathbf{Y} - \hat{\mathbf{Y}})^\top (\mathbf{Y} - \hat{\mathbf{Y}}). \end{aligned}$$

- Here, the last but one equality holds due to the following consideration:

$$\begin{aligned} \sum_{i=1}^n (Y_i - \hat{Y}_i) \hat{Y}_i &= (\mathbf{Y} - \hat{\mathbf{Y}})^\top \hat{\mathbf{Y}} = (\mathbf{Y}^\top - \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top) \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{Y}^\top \mathbf{X} \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^\top \underbrace{\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}}}_{=\mathbf{X}^\top \mathbf{Y}} \\ &= \mathbf{Y}^\top \mathbf{X} \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{Y} = (\mathbf{Y}^\top \mathbf{X} \hat{\boldsymbol{\beta}})^\top - \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{Y} = 0. \end{aligned} \quad \square$$

**Remark**

- The first summand  $\hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}$  on the right side of (42) is the squared length of the vector  $\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}}$  of the estimated target values  $\hat{Y}_1, \dots, \hat{Y}_n$ .
- The second component of the decomposition (42), i.e., the sum of the deviation squares  $(\mathbf{Y} - \hat{\mathbf{Y}})^\top (\mathbf{Y} - \hat{\mathbf{Y}})$ , is called *residual variance*.
- Sometimes the so-called *coefficient of determination*  $R^2$  is considered as well, which is given by

$$R^2 = 1 - \frac{(\mathbf{Y} - \hat{\mathbf{Y}})^\top (\mathbf{Y} - \hat{\mathbf{Y}})}{\sum_{i=1}^n (Y_i - \bar{Y})^2}, \quad \text{where} \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

Our model assumption that the design matrix  $\mathbf{X}$  has full rank, i.e.,  $\text{rk}(\mathbf{X}) = m$ , implies that the inequality  $(\mathbf{X}\boldsymbol{\beta})^\top (\mathbf{X}\boldsymbol{\beta}) = \boldsymbol{\beta}^\top (\mathbf{X}^\top \mathbf{X}) \boldsymbol{\beta} > 0$  holds if the hypothesis  $H_0 : \beta_1 = \dots = \beta_m = 0$  is not true.

- Therefore, it is natural to reject the hypothesis  $H_0$  if the squared length  $\hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}$  of the random vector  $\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}}$  is sufficiently large.
- In order to decide what “sufficiently large” means in this context, we also consider the variability  $\sigma^2$  of the data.

– Assuming that  $H_0 : \boldsymbol{\beta} = \mathbf{o}$  is true, it holds that

$$\mathbb{E}(\mathbf{Y}^\top \mathbf{Y}) = \mathbb{E}(\boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}) = \mathbb{E}\left(\sum_{i=1}^n \varepsilon_i^2\right) = \sum_{i=1}^n \mathbb{E} \varepsilon_i^2 = n\sigma^2.$$

– In this case, due to Theorem 2.9, one has

$$n\sigma^2 = \mathbb{E}(\hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}) + \mathbb{E}((\mathbf{Y} - \hat{\mathbf{Y}})^\top (\mathbf{Y} - \hat{\mathbf{Y}})),$$

– which is the reason why the quotient of  $\hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}$  and the sum of deviation squares  $(\mathbf{Y} - \hat{\mathbf{Y}})^\top (\mathbf{Y} - \hat{\mathbf{Y}})$  is considered for testing the hypothesis  $H_0 : \boldsymbol{\beta} = \mathbf{o}$ .

- More precisely, we consider the following test statistic

$$T_{\text{mod}} = \frac{\widehat{\boldsymbol{\beta}}^\top (\mathbf{X}^\top \mathbf{X}) \widehat{\boldsymbol{\beta}}}{mS^2}. \quad (43)$$

For being able to construct a test of the hypothesis  $H_0 : \beta_1 = \dots = \beta_m = 0$  based on  $T_{\text{mod}}$ , the distribution of the test statistic  $T_{\text{mod}}$  has to be determined.

**Theorem 2.10** *Assuming that  $H_0 : \beta_1 = \dots = \beta_m = 0$  is true, it holds that*

$$T_{\text{mod}} \sim F_{m, n-m}, \quad (44)$$

*i. e., the test statistic  $T_{\text{mod}}$  given in (43) has an F-distribution with  $(m, n - m)$  degrees of freedom.*

### Proof

- Assuming that  $H_0 : \beta_1 = \dots = \beta_m = 0$  is true, it holds that  $\widehat{\boldsymbol{\beta}} \sim N(\mathbf{o}, \mathbf{K})$  with  $\mathbf{K} = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$ .
- This implies that  $(\sigma^{-1}\mathbf{X})^\top (\sigma^{-1}\mathbf{X})\mathbf{K} = (\sigma^{-1}\mathbf{X})^\top (\sigma^{-1}\mathbf{X})\sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{I}$ , i. e., in particular, that the matrix  $(\sigma^{-1}\mathbf{X})^\top (\sigma^{-1}\mathbf{X})\mathbf{K}$  is idempotent.
- Now it follows from Theorem 1.9 that the quadratic form  $\sigma^{-2}\widehat{\boldsymbol{\beta}}^\top (\mathbf{X}^\top \mathbf{X}) \widehat{\boldsymbol{\beta}}$  has a (noncentral)  $\chi^2$ -distribution with  $m$  degrees of freedom.
- Moreover, we have shown in Theorem 2.7 that the random variable  $(n - m)S^2/\sigma^2$  has a (central)  $\chi^2$ -distribution with  $n - m$  degrees of freedom.
- In Theorem 2.8 we have shown that  $\widehat{\boldsymbol{\beta}}$  and  $S^2$  are independent.
- Thus, it follows from the transformation theorem for independent random vectors (cf. Theorem I-1.8) that the random variables  $\sigma^{-2}\widehat{\boldsymbol{\beta}}^\top (\mathbf{X}^\top \mathbf{X}) \widehat{\boldsymbol{\beta}}$  and  $(n - m)S^2/\sigma^2$  are independent as well.
- Now the statement follows from the definition of the F-distribution, cf. Section I-3.1.3.  $\square$

### Remark

- When testing the hypothesis  $H_0 : \beta_1 = \dots = \beta_m = 0$  with a significance level of  $\alpha \in (0, 1)$  (against the alternative  $H_1 : \beta_j \neq 0$  for at least one  $j \in \{1, \dots, m\}$ ), the null hypothesis  $H_0$  is rejected if

$$T_{\text{mod}} > F_{m, n-m, 1-\alpha}, \quad (45)$$

where  $F_{m, n-m, 1-\alpha}$  denotes the  $1 - \alpha$  quantile of the F-distribution with  $(m, n - m)$  degrees of freedom.

- In a similar way, an F-test for the verification of the hypothesis  $H_0 : \boldsymbol{\beta} = \boldsymbol{\beta}_0$  with significance level  $\alpha \in (0, 1)$  (against the alternative  $H_1 : \boldsymbol{\beta} \neq \boldsymbol{\beta}_0$ ) for an arbitrary hypothetical parameter vector  $\boldsymbol{\beta}_0 = (\beta_{01}, \dots, \beta_{0m})$  can be constructed.
- Proceeding as in the proof of Theorem 2.10, one can show that if  $H_0 : \boldsymbol{\beta} = \boldsymbol{\beta}_0$  is true, the test statistic

$$T_{\boldsymbol{\beta}_0} = \frac{(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top (\mathbf{X}^\top \mathbf{X}) (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)}{mS^2} \quad (46)$$

has an F-distribution with  $(m, n - m)$  degrees of freedom.

- Thus, the null hypothesis  $H_0 : \boldsymbol{\beta} = \boldsymbol{\beta}_0$  is rejected if

$$T_{\boldsymbol{\beta}_0} > F_{m, n-m, 1-\alpha}. \quad (47)$$

For the verification of hypotheses about *single* components of  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)^\top$ , t-tests are used instead.

- Let  $j \in \{1, \dots, m\}$ . In order to test a hypothetical value  $\beta_{0,j}$  of the  $j$ -th component  $\beta_j$  of the parameter vector  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)^\top$ , we consider the test statistic

$$T_j = \frac{\widehat{\beta}_j - \beta_j}{S\sqrt{x^{jj}}}, \quad (48)$$

where  $x^{ij}$  denotes the entry of the (inverse) matrix  $(\mathbf{X}^\top \mathbf{X})^{-1}$  at position  $(i, j)$ .

- From (40) – (41) and from the independence of  $\widehat{\boldsymbol{\beta}}$  and  $S^2$  it follows that  $T_j \sim t_{n-m}$ .
- When testing the hypothesis  $H_0 : \beta_j = \beta_{0,j}$  with a significance level of  $\alpha \in (0, 1)$  (against the alternative  $H_1 : \beta_j \neq \beta_{0,j}$ ), the null hypothesis  $H_0$  is rejected if

$$\frac{|\widehat{\beta}_j - \beta_{0,j}|}{S\sqrt{x^{jj}}} > t_{n-m, 1-\alpha/2}, \quad (49)$$

where  $t_{n-m, 1-\alpha/2}$  denotes the  $(1 - \alpha/2)$ -quantile of the  $t$ -distribution with  $n - m$  degrees of freedom.

### Remark

- The test of the hypothesis  $H_0 : \beta_j = 0$  (against the alternative  $H_1 : \beta_j \neq 0$ ) is particularly interesting because by using it, one can verify how far the response variables  $Y_1, \dots, Y_n$  depend on the  $j$ -th predictor at all.
- In this test of significance of the  $j$ -th predictor, the null hypothesis  $H_0 : \beta_j = 0$  is rejected if

$$\frac{|\widehat{\beta}_j|}{S\sqrt{x^{jj}}} > t_{n-m, 1-\alpha/2}. \quad (50)$$

The tests we considered up to now in this section are special cases of the following *ubiquitous test*. Here, an *arbitrary* part of the components of the parameter vector  $\boldsymbol{\beta}$  is tested.

- For  $\ell \in \{1, \dots, m\}$  and  $\beta_{0\ell}, \dots, \beta_{0m} \in \mathbb{R}$  the hypothesis

$$H_0 : \beta_\ell = \beta_{0\ell}, \dots, \beta_m = \beta_{0m} \quad \text{versus} \quad H_1 : \beta_j \neq \beta_{0j} \text{ for at least one } j \in \{\ell, \dots, m\} \quad (51)$$

shall be tested.

- For this purpose, we consider the following  $(m - \ell + 1) \times (m - \ell + 1)$ -dimensional submatrix  $\mathbf{K}_{\text{uni}}$  of the matrix  $(\mathbf{X}^\top \mathbf{X})^{-1} = (x^{ij})$  with

$$\mathbf{K}_{\text{uni}} = \begin{pmatrix} x^{\ell\ell} & \dots & x^{\ell m} \\ \vdots & & \vdots \\ x^{m\ell} & \dots & x^{mm} \end{pmatrix}.$$

- One can show that the inverse matrix  $\mathbf{K}_{\text{uni}}^{-1}$  is well-defined because  $\mathbf{K}_{\text{uni}} = \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top$ , where  $\mathbf{H} = (\mathbf{0}, \mathbf{I})$ , the null matrix  $\mathbf{0}$  has the dimension  $(m - \ell + 1) \times (\ell - 1)$  and the identity matrix  $\mathbf{I}$  has the dimension  $(m - \ell + 1) \times (m - \ell + 1)$ .
- Herefrom and from Lemma 1.8 it follows that the matrix  $\mathbf{K}_{\text{uni}}$  is positive definite and thus invertible.
- A possible approach to solve the testing problem given in (51) is then given by the test statistic

$$T_{\text{uni}} = \frac{(\widehat{\boldsymbol{\beta}}_{\text{uni}} - \boldsymbol{\beta}_{\text{uni}})^\top \mathbf{K}_{\text{uni}}^{-1} (\widehat{\boldsymbol{\beta}}_{\text{uni}} - \boldsymbol{\beta}_{\text{uni}})}{(m - \ell + 1)S^2}, \quad (52)$$

where  $\widehat{\boldsymbol{\beta}}_{\text{uni}} = (\widehat{\beta}_\ell, \dots, \widehat{\beta}_m)$  and  $\boldsymbol{\beta}_{\text{uni}} = (\beta_{0\ell}, \dots, \beta_{0m})$ .

- The following Theorem 2.11 implies that, assuming that the null hypothesis  $H_0$  formulated in (51) is true,

$$T_{\text{uni}} \sim F_{m-\ell+1, n-m}. \quad (53)$$

- Hence, the hypothesis  $H_0 : \beta_\ell = \beta_{0,\ell}, \dots, \beta_m = \beta_{0m}$  is rejected if

$$T_{\text{uni}} > F_{m-\ell+1, n-m, 1-\alpha}. \quad (54)$$

We now discuss one further test, a general *test for linear forms* of the parameter vector  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)$ .

- Let  $r \in \{1, \dots, m\}$ , let  $\mathbf{H}$  be an  $r \times m$  matrix with full rank  $\text{rk}(\mathbf{H}) = r$ , and let  $\mathbf{c} \in \mathbb{R}^r$ .
- The hypothesis to be tested is

$$H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{c} \quad \text{versus} \quad H_1 : \mathbf{H}\boldsymbol{\beta} \neq \mathbf{c}, \quad (55)$$

where the following test statistic  $T_{\mathbf{H}}$  is considered:

$$T_{\mathbf{H}} = \frac{(\mathbf{H}\widehat{\boldsymbol{\beta}} - \mathbf{c})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\widehat{\boldsymbol{\beta}} - \mathbf{c})}{rS^2}. \quad (56)$$

**Theorem 2.11** *Assuming that  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{c}$  is true, it holds that*

$$T_{\mathbf{H}} \sim F_{r, n-m}, \quad (57)$$

*i.e., the test statistic  $T_{\mathbf{H}}$  given in (56) has an F-distribution with  $(r, n - m)$  degrees of freedom.*

### Proof

- As the design matrix  $\mathbf{X}$  has full rank, the symmetric matrix  $\mathbf{X}^\top \mathbf{X}$  is positive definite.
  - Due to Lemma 1.8 the matrices  $(\mathbf{X}^\top \mathbf{X})^{-1}$  and  $\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top$  are then positive definite as well,
  - i.e., in particular, the matrix  $\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top$  has full rank and thus is invertible.
- Therefore, the quantity  $\mathbf{Z}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} \mathbf{Z}$  considered in (56) is well-defined, where

$$\mathbf{Z} = \mathbf{H}\widehat{\boldsymbol{\beta}} - \mathbf{c} \quad \text{with} \quad \widehat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}).$$

- Theorem 1.3 implies that, assuming that  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{c}$  is true, it holds that

$$\mathbf{Z} \sim N(\mathbf{0}, \sigma^2 \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top).$$

- Furthermore, the  $r \times r$  matrix  $\mathbf{A} = (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1}$  is symmetric because

$$\begin{aligned} \mathbf{A}^\top &= ((\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1})^\top = ((\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^\top)^{-1} = (\mathbf{H}((\mathbf{X}^\top \mathbf{X})^{-1})^\top \mathbf{H}^\top)^{-1} \\ &= (\mathbf{H}((\mathbf{X}^\top \mathbf{X})^\top)^{-1} \mathbf{H}^\top)^{-1} = (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} = \mathbf{A}. \end{aligned}$$

- As the matrix  $(\sigma^{-2} \mathbf{A})(\sigma^2 \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top) = \mathbf{I}$  obviously is idempotent, Theorem 1.9 implies that  $\sigma^{-2} \mathbf{Z}^\top \mathbf{A} \mathbf{Z}$  is a random variable having a  $\chi_r^2$ -distribution.
- The rest of the proof continues in the same way as the proof of Theorem 2.10.  $\square$

**Remark** The null hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{c}$  is rejected if  $T_{\mathbf{H}} > F_{r, n-m, 1-\alpha}$ , where  $T_{\mathbf{H}}$  is the test statistic given in (56).



### 2.2.4 Confidence Regions; Prediction of Response Variables

- *Recall:* In Section 2.2.3 we have considered the test statistic  $T_j = (\hat{\beta}_j - \beta_j)/(S\sqrt{x^{jj}})$ , where  $x^{jj}$  denotes the entry of the (inverse) matrix  $(\mathbf{X}^\top \mathbf{X})^{-1}$  at position  $(i, j)$ .
- In doing so, we have shown that  $T_j \sim t_{n-m}$  for each  $j \in \{1, \dots, m\}$ .
- This leads to the following confidence intervals with confidence level  $1 - \alpha \in (0, 1)$  for each single regression coefficient  $\beta_j$ .
- It holds for each  $j \in \{1, \dots, m\}$  with probability  $1 - \alpha$  that

$$\hat{\beta}_j - t_{n-m, 1-\alpha/2} S\sqrt{x^{jj}} < \beta_j < \hat{\beta}_j + t_{n-m, 1-\alpha/2} S\sqrt{x^{jj}}. \quad (58)$$

#### Remark

- In the same way as in the proof of Theorem I-5.8 a *common confidence region* with confidence level  $1 - \alpha \in (0, 1)$  for all  $m$  regression coefficients  $\beta_1, \dots, \beta_m$  can be derived by use of the Bonferroni inequality (cf. Lemma I-5.4).
- Indeed, the probability that

$$\hat{\beta}_j - t_{n-m, 1-\alpha/2m} S\sqrt{x^{jj}} < \beta_j < \hat{\beta}_j + t_{n-m, 1-\alpha/2m} S\sqrt{x^{jj}} \quad (59)$$

for all  $j = 1, \dots, m$  at the same time is at least equal to  $1 - \alpha$ .

- Moreover, Theorem 2.10 leads to an *exact* common confidence region with confidence level  $1 - \alpha$  for all  $m$  regression coefficients  $\beta_1, \dots, \beta_m$ .
  - It holds (cf. (46) – (47)) that

$$\mathbb{P}_\beta \left( \frac{(\hat{\beta} - \beta)^\top (\mathbf{X}^\top \mathbf{X}) (\hat{\beta} - \beta)}{mS^2} < F_{m, n-m, 1-\alpha} \right) = 1 - \alpha.$$

- Here, the confidence region  $E$  with

$$E = \left\{ \beta = (\beta_1, \dots, \beta_m) : \frac{(\hat{\beta} - \beta)^\top (\mathbf{X}^\top \mathbf{X}) (\hat{\beta} - \beta)}{mS^2} < F_{m, n-m, 1-\alpha} \right\}$$

forms a (random) *ellipsoid* with center  $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_m)$ .

- One can show that the ellipsoid  $E$  can be embedded into an  $m$ -dimensional paraxial cuboid  $E' \supset E$ , where

$$E' = \prod_{j=1}^m \left( \hat{\beta}_j - S\sqrt{m x^{jj} F_{m, n-m, 1-\alpha}}, \hat{\beta}_j + S\sqrt{m x^{jj} F_{m, n-m, 1-\alpha}} \right).$$

- The confidence region  $E'$  has a simpler form than  $E$ . However, due to  $E' \supset E$  it is clear that  $E'$  is an estimation, which is less accurate than  $E$ .

In a similar way, a *confidence interval for the expected target value*

$$\varphi(x_{01}, \dots, x_{0m}) = \beta_1 x_{01} + \dots + \beta_m x_{0m}$$

corresponding to a given vector  $\mathbf{x}_0 = (x_{01}, \dots, x_{0m})^\top \in \mathbb{R}^m$  of values  $x_{01}, \dots, x_{0m}$  of the  $m$  predictor variables can be derived.

- For this purpose, we consider the  $1 \times m$  matrix  $\mathbf{H} = (x_{01}, \dots, x_{0m}) (= \mathbf{x}_0^\top)$ .

- Then, Theorem 2.11 implies that

$$\sqrt{T_{\mathbf{H}}} = \frac{|\widehat{\boldsymbol{\beta}}^{\top} \mathbf{x}_0 - \varphi(\mathbf{x}_0)|}{S \sqrt{\mathbf{x}_0^{\top} (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{x}_0}} \stackrel{d}{=} |T|,$$

where  $T$  is a random variable having a  $t$ -distribution with  $n - m$  degrees of freedom.

- Hence, it holds with probability  $1 - \alpha$  that

$$\widehat{\boldsymbol{\beta}}^{\top} \mathbf{x}_0 - Z_0 < \varphi(\mathbf{x}_0) < \widehat{\boldsymbol{\beta}}^{\top} \mathbf{x}_0 + Z_0, \quad (60)$$

where

$$Z_0 = t_{n-m, 1-\alpha/2} S \sqrt{\mathbf{x}_0^{\top} (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{x}_0}.$$

### Remark

- Analogously, one can derive a *prediction interval* for the response variable  $Y_0 = \beta_1 x_{01} + \dots + \beta_m x_{0m} + \varepsilon_0$ , where the error term  $\varepsilon_0$  is normally distributed and independent of the error terms  $\varepsilon_1, \dots, \varepsilon_n$ ;  $\varepsilon_0 \sim N(0, \sigma^2)$ .
- Indeed, it is  $\widehat{\boldsymbol{\beta}}^{\top} \mathbf{x}_0 - Y_0 \sim N(0, \sigma^2(1 + \mathbf{x}_0^{\top} (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{x}_0))$  and thus

$$\widehat{\boldsymbol{\beta}}^{\top} \mathbf{x}_0 - Z'_0 < Y_0 < \widehat{\boldsymbol{\beta}}^{\top} \mathbf{x}_0 + Z'_0, \quad (61)$$

with probability  $1 - \alpha$ , where  $Z'_0 = t_{n-m, 1-\alpha/2} S \sqrt{1 + \mathbf{x}_0^{\top} (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{x}_0}$ .

### 2.2.5 Confidence Band

In this section we assume that the design matrix  $\mathbf{X}$  has the form

$$\mathbf{X} = \begin{pmatrix} 1 & x_{12} & \dots & x_{1m} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n2} & \dots & x_{nm} \end{pmatrix}, \quad (62)$$

i.e., we consider the (multiple) linear regression model.

- In the definition of the regression function  $\varphi(x_1, \dots, x_m) = \beta_1 x_1 + \dots + \beta_m x_m$  in (4), we now set  $x_1 = 1$  and determine a *confidence band for the regression hyperplane*

$$y = \varphi(1, x_2, \dots, x_m) = \beta_1 + \beta_2 x_2 + \dots + \beta_m x_m, \quad \forall x_2, \dots, x_m \in \mathbb{R}.$$

- This means that we need to find a number  $a_\gamma > 0$ , such that with the given (coverage) probability  $\gamma = 1 - \alpha \in (0, 1)$  it holds that

$$\widehat{\beta}_1 + \widehat{\beta}_2 x_2 + \dots + \widehat{\beta}_m x_m - a_\gamma Z_{\mathbf{x}} < \varphi(1, x_2, \dots, x_m) < \widehat{\beta}_1 + \widehat{\beta}_2 x_2 + \dots + \widehat{\beta}_m x_m + a_\gamma Z_{\mathbf{x}} \quad (63)$$

for each  $\mathbf{x} = (1, x_2, \dots, x_m) \in \mathbb{R}^m$  simultaneously, where

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{Y} \quad \text{and} \quad Z_{\mathbf{x}} = S \sqrt{\mathbf{x}^{\top} (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{x}}.$$

For solving this problem, the following lemma is useful.

**Lemma 2.3** *It holds with probability 1 that*

$$\max_{\mathbf{x} \in \mathbb{R}_1^{m-1}} \frac{((\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x})^2}{\mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} = (\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \boldsymbol{\varepsilon}), \quad (64)$$

where  $\mathbb{R}_1^{m-1}$  denotes the set of all vectors  $\mathbf{x} \in \mathbb{R}^m$  with  $\mathbf{x} = (1, x_2, \dots, x_m)^\top$ .

**Proof**

- From Lemmas 1.6 and 1.8 it follows that  $(\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{H}\mathbf{H}^\top$  for an invertible  $m \times m$  matrix  $\mathbf{H}$ .
- Therefore, the expression

$$(\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x} = ((\mathbf{X}\mathbf{H})^\top \boldsymbol{\varepsilon})^\top \mathbf{H}^\top \mathbf{x}$$

can be perceived as the scalar product of the  $m$ -dimensional vectors  $(\mathbf{X}\mathbf{H})^\top \boldsymbol{\varepsilon}$  and  $\mathbf{H}^\top \mathbf{x}$ .

- Analogously, it holds that

$$(\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon} = ((\mathbf{X}\mathbf{H})^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}\mathbf{H})^\top \boldsymbol{\varepsilon} \quad \text{and} \quad \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x} = (\mathbf{H}^\top \mathbf{x})^\top \mathbf{H}^\top \mathbf{x}.$$

- From this result and from the inequality

$$|\mathbf{y}^\top \mathbf{z}| \leq \sqrt{\mathbf{y}^\top \mathbf{y}} \sqrt{\mathbf{z}^\top \mathbf{z}} \quad \forall \mathbf{y}, \mathbf{z} \in \mathbb{R}^m \quad (65)$$

together with  $\mathbf{y} = (\mathbf{X}\mathbf{H})^\top \boldsymbol{\varepsilon}$  and  $\mathbf{z} = \mathbf{H}^\top \mathbf{x}$  it follows that

$$|(\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}| \leq \sqrt{(\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \boldsymbol{\varepsilon})} \sqrt{\mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}}$$

and thus

$$\frac{((\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x})^2}{\mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} \leq (\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \boldsymbol{\varepsilon}). \quad (66)$$

- As the random vector  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$  has independent and absolutely continuous components, we get that  $\sum_{i=1}^n \varepsilon_i \neq 0$  with probability 1.
- Now let  $\sum_{i=1}^n \varepsilon_i \neq 0$ . Then it follows from the form of the design matrix  $\mathbf{X}$  considered in (62) that the vector  $\mathbf{x} = \mathbf{X}^\top \boldsymbol{\varepsilon} / \sum_{i=1}^n \varepsilon_i$  belongs to  $\mathbb{R}_1^{m-1}$  and that in this case the equality in (66) holds.  $\square$

The following result, which is a vectorial generalization of Theorem I-5.9, leads to the desired confidence band.

**Theorem 2.12** *Let  $a_\gamma = \sqrt{m F_{m, n-m, \gamma}}$ . Then it holds that*

$$\mathbb{P}_\beta \left( \max_{\mathbf{x} \in \mathbb{R}_1^{m-1}} \frac{(\hat{\boldsymbol{\beta}}^\top \mathbf{x} - \varphi(\mathbf{x}))^2}{S^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} \leq a_\gamma^2 \right) = \gamma. \quad (67)$$

**Proof**

- For each  $\mathbf{x} \in \mathbb{R}_1^{m-1}$  it holds that

$$\begin{aligned} \hat{\boldsymbol{\beta}}^\top \mathbf{x} - \varphi(\mathbf{x}) &= \hat{\boldsymbol{\beta}}^\top \mathbf{x} - \boldsymbol{\beta}^\top \mathbf{x} = ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y})^\top \mathbf{x} - \boldsymbol{\beta}^\top \mathbf{x} = (\boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon})^\top \mathbf{x} - \boldsymbol{\beta}^\top \mathbf{x} \\ &= (\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x} \end{aligned}$$

and thus

$$\begin{aligned} \max_{\mathbf{x} \in \mathbb{R}_1^{m-1}} \frac{(\hat{\boldsymbol{\beta}}^\top \mathbf{x} - \varphi(\mathbf{x}))^2}{S^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} &= \max_{\mathbf{x} \in \mathbb{R}_1^{m-1}} \frac{((\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x})^2}{S^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} = \frac{1}{S^2} \max_{\mathbf{x} \in \mathbb{R}_1^{m-1}} \frac{((\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x})^2}{\mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} \\ &= \frac{(\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \boldsymbol{\varepsilon})}{S^2}, \end{aligned}$$

where the last equality follows from Lemma 2.3.

- Therefore, one has

$$\max_{\mathbf{x} \in \mathbb{R}_1^{m-1}} \frac{(\widehat{\boldsymbol{\beta}}^\top \mathbf{x} - \varphi(\mathbf{x}))^2}{S^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} = \frac{(\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \boldsymbol{\varepsilon})}{S^2}. \quad (68)$$

- Due to  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$  and

$$\mathbf{X}^\top (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \mathbf{0},$$

Theorem 1.10 implies that  $\mathbf{X}^\top \boldsymbol{\varepsilon}$  and  $\boldsymbol{\varepsilon}^\top (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \boldsymbol{\varepsilon}$  are independent.

- Hence, it follows from the representation formula

$$S^2 = \frac{1}{n-m} \boldsymbol{\varepsilon}^\top (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \boldsymbol{\varepsilon},$$

which has already been derived in Theorem 2.7, that also the random variables  $(\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \boldsymbol{\varepsilon})$  and  $S^2$  are independent.

- In Theorem 2.7 we have shown that

$$(n-m)S^2/\sigma^2 \sim \chi_{n-m}^2.$$

- Moreover, Theorem 1.9 implies that

$$(\mathbf{X}^\top \boldsymbol{\varepsilon})^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \boldsymbol{\varepsilon})/\sigma^2 \sim \chi_m^2$$

since the  $m \times m$  (covariance) matrix  $\mathbf{X}^\top \mathbf{X}$  of the normally distributed random vector  $\mathbf{X}^\top \boldsymbol{\varepsilon}$  has full rank and since the matrix  $(\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{X}) = \mathbf{I}$  is idempotent.

- Due to (68) we have altogether shown that

$$\frac{1}{m} \max_{\mathbf{x} \in \mathbb{R}_1^{m-1}} \frac{(\widehat{\boldsymbol{\beta}}^\top \mathbf{x} - \varphi(\mathbf{x}))^2}{S^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} \sim F_{m, n-m}.$$

- For the threshold value considered in (63) and (67), we hence obtain  $a_\gamma = \sqrt{m \overline{F}_{m, n-m, \gamma}}$ .  $\square$

### 3 Arbitrary Design Matrix; Generalized Inverse

- We now consider the following generalization of the linear model discussed in Chapter 2,

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (1)$$

for which we have assumed so far that the design matrix

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{pmatrix} \quad (2)$$

is an  $(n \times m)$ -dimensional matrix with full (column) rank  $\text{rk}(\mathbf{X}) = m$ , where  $n \geq m$ .

- In contrast to this, in this chapter we will consider the case that  $\text{rk}(\mathbf{X}) \leq m$ , i.e., we allow that the design matrix  $\mathbf{X}$  does *not* have full rank.
- As we did in Section 2.1, we first assume for the random vector  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$  that

$$\mathbb{E}\varepsilon_i = 0, \quad \text{Var}\varepsilon_i = \sigma^2, \quad \text{Cov}(\varepsilon_i, \varepsilon_j) = 0, \quad \forall i, j = 1, \dots, n \text{ with } i \neq j \quad (3)$$

for a certain (unknown)  $\sigma^2 > 0$ .

#### 3.1 Analysis of Variance as a Linear Model

To begin with, we discuss two examples of problems leading to linear models whose design matrix does not have full rank, cf. Section 3.4.

The term „analysis of variance” does not mean that variances of random variables are analyzed in this context, but refers to the analysis of the *variability of expectations*. In literature, ANOVA is typically used as an abbreviation.

##### 3.1.1 One-Factor Analysis of Variance; ANOVA Null Hypothesis

- In a one-factor analysis of variance, we assume that the random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  can be partitioned into  $k$  classes of subsamples  $(Y_{ij}, j = 1, \dots, n_i)$ , where

- $n_i > 1$  for each  $i = 1, \dots, k$  and  $\sum_{i=1}^k n_i = n$
- and the sampling variables belonging to the same class have the same expectation  $\theta_i$ .

- In other words: We assume that

$$Y_{ij} = \theta_i + \varepsilon_{ij}, \quad \forall i = 1, \dots, k, j = 1, \dots, n_i, \quad (4)$$

where  $\theta_1, \dots, \theta_k \in \mathbb{R}$  are (unknown) parameters and the error terms  $\varepsilon_{ij} : \Omega \rightarrow \mathbb{R}$  are uncorrelated with

$$\mathbb{E}\varepsilon_{ij} = 0, \quad \text{Var}\varepsilon_{ij} = \sigma^2, \quad \forall i = 1, \dots, k, j = 1, \dots, n_i. \quad (5)$$

#### Remark

- The numbers  $i = 1, \dots, k$  of the classes  $(Y_{ij}, j = 1, \dots, n_i)$  are interpreted as *levels of a predictor variable*.

- The model assumptions made above imply in particular that the observed values  $y_1, \dots, y_n$  of the response variables  $Y_1, \dots, Y_n$  can be structured in table form as follows:

level	1	2	3	...	$k$
	$y_{11}$	$y_{21}$	$y_{31}$	...	$y_{k1}$
	$y_{12}$	$y_{22}$	$y_{32}$	...	$y_{k2}$
	$\vdots$	$\vdots$	$\vdots$	...	$y_{k3}$
			$y_{3n_3}$		$\vdots$
	$y_{1n_1}$				
		$y_{2n_2}$			$y_{kn_k}$

We show that the classical ANOVA null hypothesis  $H_0 : \theta_1 = \dots = \theta_k$  can be expressed by use of so-called contrasts.

- For this purpose, we consider the following set  $\mathcal{A} \subset \mathbb{R}^k$  with

$$\mathcal{A} = \left\{ \mathbf{a} = (a_1, \dots, a_k)^\top : \mathbf{a} \neq \mathbf{0}, \sum_{i=1}^k a_i = 0 \right\}.$$

- Let  $\mathbf{t} = (t_1, \dots, t_k)^\top \in \mathbb{R}^k$  be an arbitrary vector of variables and let  $\mathbf{a} = (a_1, \dots, a_k)^\top \in \mathcal{A}$  be a vector of (known) constants. The mapping  $\mathbf{t} \rightarrow \sum_{i=1}^k a_i t_i$  is then called a *contrast*.

**Lemma 3.1** *Let  $\theta_1, \dots, \theta_k \in \mathbb{R}$  be arbitrary real numbers. For the validity of  $\theta_1 = \dots = \theta_k$  it is then necessary and sufficient that*

$$\sum_{i=1}^k a_i \theta_i = 0 \quad \forall \mathbf{a} \in \mathcal{A}. \quad (6)$$

**Proof**

- If  $\theta_1 = \dots = \theta_k = \theta$  is true, we get for each  $\mathbf{a} \in \mathcal{A}$  that

$$\sum_{i=1}^k a_i \theta_i = \theta \sum_{i=1}^k a_i = 0.$$

- In order to show the sufficiency of the condition, we consider the vectors  $\mathbf{a}_1, \dots, \mathbf{a}_{k-1} \in \mathcal{A}$  with

$$\mathbf{a}_1 = (1, -1, 0, \dots, 0)^\top, \quad \mathbf{a}_2 = (0, 1, -1, 0, \dots, 0)^\top, \quad \dots, \quad \mathbf{a}_{k-1} = (0, \dots, 0, 1, -1)^\top.$$

- For each  $i \in \{1, \dots, k-1\}$  the validity of condition (6) for  $\mathbf{a}_i$  implies that  $-\theta_i + \theta_{i+1} = 0$ , i.e.,  $\theta_i = \theta_{i+1}$ . Therefore, it follows that  $\theta_1 = \dots = \theta_k$ .  $\square$

**Remark**

- Due to Lemma 3.1, the classical ANOVA null hypothesis  $H_0 : \theta_1 = \dots = \theta_k$  is equivalent to the hypothesis  $H_0 : \sum_{i=1}^k a_i \theta_i = 0$  for each  $\mathbf{a} = (a_1, \dots, a_k)^\top \in \mathcal{A}$ .

- Moreover, assuming that  $H_0$  is true, it is obvious that
  - $\sum_{i=1}^k a_i \bar{Y}_{i\cdot}$  with  $\bar{Y}_{i\cdot} = \sum_{j=1}^{n_i} Y_{ij}/n_i$  is an unbiased estimator for  $\sum_{i=1}^k a_i \theta_i = 0$  for each  $\mathbf{a} \in \mathcal{A}$ ,
  - the variance of  $\sum_{i=1}^k a_i \bar{Y}_{i\cdot}$  is given by  $\text{Var} \sum_{i=1}^k a_i \bar{Y}_{i\cdot} = \sigma^2 \sum_{i=1}^k a_i^2/n_i$
  - and

$$S_p^2 = \frac{1}{n-k} \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\cdot})^2 \quad (7)$$

is an unbiased estimator for  $\sigma^2$ , the so-called *pooled sample variance*.

- Hence, it is reasonable to reject  $H_0 : \theta_1 = \dots = \theta_k$  if the supremum over  $\mathbf{a} \in \mathcal{A}$  of the (suitably normed) absolute values of  $\sum_{i=1}^k a_i \bar{Y}_{i\cdot}$  exceeds a certain threshold value, where the test statistic  $\sup_{\mathbf{a} \in \mathcal{A}} T_{\mathbf{a}}^2$  is considered with  $T_{\mathbf{a}} = \left( \sum_{i=1}^k a_i \bar{Y}_{i\cdot} \right) / \sqrt{S_p^2 \sum_{i=1}^k a_i^2/n_i}$ .
- In a similar way as in the proof of Lemma 2.3 one can show that, assuming that  $H_0 : \theta_1 = \dots = \theta_k$  is true, it holds that

$$\sup_{\mathbf{a} \in \mathcal{A}} T_{\mathbf{a}}^2 = \frac{\sum_{i=1}^k n_i (\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot})^2}{S_p^2}, \quad (8)$$

where  $\bar{Y}_{\cdot\cdot} = \sum_{i=1}^k n_i \bar{Y}_{i\cdot} / \sum_{i=1}^k n_i$ .

The following *decomposition* implies an intuitive interpretation of numerator and denominator of the test statistic  $\sup_{\mathbf{a} \in \mathcal{A}} T_{\mathbf{a}}^2$  considered in (8), cf. also Theorem 2.9.

**Theorem 3.1** *It holds that*

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{\cdot\cdot})^2 = \sum_{i=1}^k n_i (\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot})^2 + \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\cdot})^2. \quad (9)$$

**Proof** By expanding the left-hand side of (9), one obtains that

$$\begin{aligned} \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{\cdot\cdot})^2 &= \sum_{i=1}^k \sum_{j=1}^{n_i} ((Y_{ij} - \bar{Y}_{i\cdot}) + (\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot}))^2 \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} ((Y_{ij} - \bar{Y}_{i\cdot})^2 + 2(Y_{ij} - \bar{Y}_{i\cdot})(\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot}) + (\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot})^2) \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\cdot})^2 + 2 \sum_{i=1}^k (\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot}) \underbrace{\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\cdot})}_{=0} + \sum_{i=1}^k n_i (\bar{Y}_{i\cdot} - \bar{Y}_{\cdot\cdot})^2. \end{aligned} \quad \square$$

**Remark**

- The double sum on the left-hand side of (9) can be interpreted as a measure for the (total) *variability of the sampling variables*  $\{Y_{ij}, i = 1, \dots, k, j = 1, \dots, n_i\}$ .
- The first sum on the right-hand side of (9) is a measure for the variability *between* the levels of the predictor variable, while the double sum on the right-hand side of (9) is a measure for the variability *within* the levels of the predictor variable.
- So due to the definition of  $S_p^2$  given in (7), the test statistic considered in (8) is proportional to the ratio of the variability between the levels and the variability within the levels of the predictor variable.
- Therefore, the ANOVA null hypothesis  $H_0 : \theta_1 = \dots = \theta_k$  is rejected if the variability between the levels is significantly higher than the variability within the levels of the predictor variable.

### 3.1.2 Reparametrization of the Expectations

The model of one-factor analysis of variance considered in Section 3.1.1 can be represented as a linear model in two different ways.

- In both cases, the random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  is “structured”, i.e., we use the notation  $\mathbf{Y} = (Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2}, \dots, Y_{k1}, \dots, Y_{kn_k})^\top$ , where  $n_1 + \dots + n_k = n$ .
- The random vector  $\mathbf{Y}$  is represented in the form  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , where the design matrix  $\mathbf{X}$  and the parameter vector  $\boldsymbol{\beta}$  are chosen differently in each case.
  - While  $\mathbf{X}$  has full rank in the first case, it does not have full rank in the second case.
  - The second (reparametrized) representation is especially useful for the application of general estimation and test methods, which are discussed in Sections 3.2 and 3.3.
  - If the error terms are normally distributed, we can thus determine the distribution of the test statistic  $\sup_{\mathbf{a} \in \mathcal{A}} T_{\mathbf{a}}^2$  considered in (8), assuming that  $H_0 : \theta_1 = \dots = \theta_k$  is true, cf. formula (89) in Section 3.4.1.

#### Case 1

- In this case the design matrix  $\mathbf{X}$  is given by the  $n \times k$  matrix

$$\mathbf{X} = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 1 \end{pmatrix}, \quad (10)$$

and the parameter vector  $\boldsymbol{\beta}$  is given by  $\boldsymbol{\beta} = (\theta_1, \dots, \theta_k)^\top$ .

#### Case 2

- We consider the following *reparametrization* of the expectations  $\theta_1, \dots, \theta_k$ , which correspond to the levels of the predictor variable.
- Let  $\mu \in \mathbb{R}$  and  $\alpha_1, \dots, \alpha_k \in \mathbb{R}$  be real numbers, such that

$$\theta_i = \mu + \alpha_i, \quad \forall i = 1, \dots, k \quad (11)$$

and

$$\sum_{i=1}^k n_i \alpha_i = 0. \quad (12)$$



- Then the random sample  $\mathbf{Y}$  of the model of one-factor analysis of variance can also be written in the form  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , where the design matrix  $\mathbf{X}$ , however, is now given by the  $n \times (k + 1)$  matrix

$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 1 & 0 & 0 & \dots & 0 & 0 \\ 1 & 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 0 & 1 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & 0 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 0 & 0 & 0 & \dots & 0 & 1 \end{pmatrix}, \quad (13)$$

and the parameter vector  $\boldsymbol{\beta}$  is given by  $\boldsymbol{\beta} = (\mu, \alpha_1, \dots, \alpha_k)^\top$ .

#### Remark

- The linear additional condition (12) for the components  $\alpha_1, \dots, \alpha_k$  of the parameter vector  $\boldsymbol{\beta}$  ensures that the representation (11) – (12) of the expectations  $\theta_1, \dots, \theta_k$  is unique.
- Furthermore, (11) and (12) imply that

$$\frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} \mathbb{E} Y_{ij} = \mu,$$

where

- the parameter  $\mu$  can be interpreted as *general mean* of the expectations  $\mathbb{E} Y_{ij}$  of the sampling variables  $Y_{ij}$  and
- the (deviation) parameter  $\alpha_i$  is called the *effect* of the  $i$ -th level of the predictor variable.
- For the design matrix  $\mathbf{X}$  given in (13) it holds that  $\text{rk}(\mathbf{X}) = k$ , i.e., the  $n \times (k + 1)$ -dimensional matrix  $\mathbf{X}$  does not have full column rank.

**Theorem 3.2** *It holds that*

$$\mathbb{E} \bar{Y}_{..} = \mu \quad \text{and} \quad \mathbb{E} (\bar{Y}_{i.} - \bar{Y}_{..}) = \alpha_i \quad (14)$$

for each  $i = 1, \dots, k$ , i.e.,  $\bar{Y}_{..}$  and  $\bar{Y}_{i.} - \bar{Y}_{..}$  define unbiased estimators for the model parameters  $\mu$  and  $\alpha_i$ , respectively.

**Proof** It follows from the definition of  $\bar{Y}_{..}$  that

$$\mathbb{E} \bar{Y}_{..} = \frac{1}{\sum_{i=1}^k n_i} \sum_{i=1}^k \sum_{j=1}^{n_i} \mathbb{E} Y_{ij} = \frac{1}{\sum_{i=1}^k n_i} \sum_{i=1}^k n_i \theta_i = \mu + \frac{1}{\sum_{i=1}^k n_i} \sum_{i=1}^k n_i \alpha_i = \mu,$$

where the last equality follows from the reparametrization condition (12). The second part of the statement in (14) can be proved analogously.  $\square$

### 3.1.3 Two-Factor Analysis of Variance

- We now modify the model of one-factor analysis of variance introduced in Section 3.1.1 and assume that the response variables  $Y_1, \dots, Y_n$  depend on *two* predictor variables.
- Thus, we partition the random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  into  $k_1 \cdot k_2$  subsamples  $(Y_{i_1 i_2 j}, j = 1, \dots, n_{i_1 i_2})$ , where  $n_{i_1 i_2} > 1$  for all  $i_1 = 1, \dots, k_1$  and  $i_2 = 1, \dots, k_2$  and

$$\sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} n_{i_1 i_2} = n.$$

- We assume that the sampling variables belonging to the same class have the same expectation  $\theta_{i_1 i_2}$  in each case.
- In other words: We assume that

$$Y_{i_1 i_2 j} = \theta_{i_1 i_2} + \varepsilon_{i_1 i_2 j}, \quad \forall i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2, j = 1, \dots, n_{i_1 i_2}, \quad (15)$$

where  $\theta_{i_1 i_2} \in \mathbb{R}$  are (unknown) parameters and the error terms  $\varepsilon_{i_1 i_2 j} : \Omega \rightarrow \mathbb{R}$  are uncorrelated with

$$\mathbb{E} \varepsilon_{i_1 i_2 j} = 0, \quad \text{Var} \varepsilon_{i_1 i_2 j} = \sigma^2, \quad \forall i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2, j = 1, \dots, n_{i_1 i_2}. \quad (16)$$

#### Remark

- The representation (15) of the sampling variables  $Y_{i_1 i_2 j}$  leads to the same form of linear model as it was considered in Case 1 of Section 3.1.2.
- The numbers  $i_1 = 1, \dots, k_1$  and  $i_2 = 1, \dots, k_2$  of the classes  $(Y_{i_1 i_2 j}, j = 1, \dots, n_{i_1 i_2})$  are again interpreted as *levels* of the corresponding predictor variable.
- Here, the design matrix  $\mathbf{X}$  has the dimension  $n \times (k_1 \cdot k_2)$  and full column rank  $k_1 \cdot k_2$ .

Moreover, we consider a similar reparametrization of the expectations  $\theta_{i_1 i_2}$  as in Section 3.1.2.

- In doing so, we only consider the so-called *balanced case*, i.e.,
  - we additionally assume that all  $k_1 \cdot k_2$  subsamples  $(Y_{i_1 i_2 j}, j = 1, \dots, n_{i_1 i_2})$  have the same sample size.
  - Hence, let  $n_{i_1 i_2} = r$  for all  $i_1 = 1, \dots, k_1$  and  $i_2 = 1, \dots, k_2$ , where  $r = n/(k_1 k_2)$ .
- Let  $\mu \in \mathbb{R}$  and for all  $i_1 \in \{1, \dots, k_1\}$  and  $i_2 \in \{1, \dots, k_2\}$  let  $\alpha_{i_1}^{(1)} \in \mathbb{R}$ ,  $\alpha_{i_2}^{(2)} \in \mathbb{R}$  and  $\alpha_{i_1 i_2} \in \mathbb{R}$  be real numbers, such that

$$\theta_{i_1 i_2} = \mu + \alpha_{i_1}^{(1)} + \alpha_{i_2}^{(2)} + \alpha_{i_1 i_2}, \quad \forall i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2 \quad (17)$$

and

$$\sum_{i_1=1}^{k_1} \alpha_{i_1}^{(1)} = \sum_{i_2=1}^{k_2} \alpha_{i_2}^{(2)} = \sum_{i_1=1}^{k_1} \alpha_{i_1 i_2} = \sum_{i_2=1}^{k_2} \alpha_{i_1 i_2} = 0. \quad (18)$$

- Then, the random sample  $\mathbf{Y}$  can be written in the form  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , where
  - the design matrix  $\mathbf{X}$  is given by a matrix of dimension  $n \times (1 + k_1 + k_2 + k_1 k_2)$ , whose entries only consist of zeros and ones and which does not have full rank.
  - Therefore, the parameter vector  $\boldsymbol{\beta}$  has the following form:

$$\boldsymbol{\beta} = (\mu, \alpha_1^{(1)}, \dots, \alpha_{k_1}^{(1)}, \alpha_1^{(2)}, \dots, \alpha_{k_2}^{(2)}, \alpha_{11}, \dots, \alpha_{k_1 k_2})^\top.$$

**Remark**

- The additional linear conditions (18) for the components of the parameter vector  $\beta$  ensure, in a similar way as in the model of one-factor analysis of variance considered in Section 3.1.2, that the representation (17) – (18) of the expectations  $\theta_{11}, \dots, \theta_{k_1 k_2}$  is unique.
- Here,
  - $\mu$  can be perceived as *general mean* of the expectations  $\mathbb{E} Y_{i_1 i_2 j}$  of the sampling variables  $Y_{i_1 i_2 j}$ ,
  - $\alpha_{i_1}^{(1)}$  is called *main effect* of the  $i_1$ -th level of the first predictor variable,
  - $\alpha_{i_2}^{(2)}$  is called *main effect* of the  $i_2$ -th level of the second predictor variable and
  - $\alpha_{i_1 i_2}$  is called *interaction* between the levels  $i_1$  and  $i_2$  of the level combination  $(i_1, i_2)$ .

For the construction of estimators for the model parameters  $\mu, \alpha_{i_1}^{(1)}, \alpha_{i_2}^{(2)}$  and  $\alpha_{i_1 i_2}$ , we use the following notation: Let

$$Y_{i_1 \cdot \cdot} = \sum_{i_2=1}^{k_2} \sum_{j=1}^r Y_{i_1 i_2 j}, \quad Y_{\cdot i_2 \cdot} = \sum_{i_1=1}^{k_1} \sum_{j=1}^r Y_{i_1 i_2 j}, \quad Y_{i_1 i_2 \cdot} = \sum_{j=1}^r Y_{i_1 i_2 j} \quad (19)$$

and

$$\bar{Y}_{i_1 \cdot \cdot} = \frac{1}{rk_2} Y_{i_1 \cdot \cdot}, \quad \bar{Y}_{\cdot i_2 \cdot} = \frac{1}{rk_1} Y_{\cdot i_2 \cdot}, \quad \bar{Y}_{i_1 i_2 \cdot} = \frac{1}{r} Y_{i_1 i_2 \cdot}, \quad \bar{Y}_{\dots} = \frac{1}{rk_1 k_2} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r Y_{i_1 i_2 j} \quad (20)$$

**Theorem 3.3** *It holds that*

$$\mathbb{E} \bar{Y}_{\dots} = \mu, \quad \mathbb{E} (\bar{Y}_{i_1 \cdot \cdot} - \bar{Y}_{\dots}) = \alpha_{i_1}^{(1)}, \quad \mathbb{E} (\bar{Y}_{\cdot i_2 \cdot} - \bar{Y}_{\dots}) = \alpha_{i_2}^{(2)}, \quad \mathbb{E} (\bar{Y}_{\dots} + \bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot \cdot} - \bar{Y}_{\cdot i_2 \cdot}) = \alpha_{i_1 i_2} \quad (21)$$

for arbitrary  $i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2$ , i.e.,  $\bar{Y}_{\dots}, \bar{Y}_{i_1 \cdot \cdot} - \bar{Y}_{\dots}, \bar{Y}_{\cdot i_2 \cdot} - \bar{Y}_{\dots}$  and  $\bar{Y}_{\dots} + \bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot \cdot} - \bar{Y}_{\cdot i_2 \cdot}$  define unbiased estimators for the model parameters  $\mu, \alpha_{i_1}^{(1)}, \alpha_{i_2}^{(2)}$  and  $\alpha_{i_1 i_2}$ , respectively.

**Proof** It follows from the definition of  $\bar{Y}_{\dots}$  in (20) that

$$\begin{aligned} \mathbb{E} \bar{Y}_{\dots} &= \frac{1}{rk_1 k_2} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r \mathbb{E} Y_{i_1 i_2 j} = \frac{1}{k_1 k_2} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \theta_{i_1 i_2} \\ &= \mu + \frac{1}{k_1 k_2} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\alpha_{i_1}^{(1)} + \alpha_{i_2}^{(2)} + \alpha_{i_1 i_2}) = \mu, \end{aligned}$$

where the last equality follows from the reparametrization conditions (18). The remaining three parts of the statement in (21) can be proved analogously.  $\square$

**Remark**

- The conditions (18), i.e., the assumption that the parameter vector  $\beta$  belongs to a linear subspace of  $\mathbb{R}^{1+k_1+k_2+k_1 k_2}$ , play a fundamental role in the proof of Theorem 3.3.
- Here, the conclusions of Theorem 3.3 can be interpreted as unbiasedness of the considered estimators with respect to this restricted parameter space.
- However, if we allow that  $\beta$  is an *arbitrary* vector of dimension  $1 + k_1 + k_2 + k_1 k_2$ , there is *no* LS-estimator for  $\beta$ , which is unbiased at the same time, cf. the discussion at the end of Section 3.2.1.

The following result contains a *decomposition of sums of squared differences*, cf. also Theorems 2.9 and 3.1.

**Theorem 3.4** *It holds that*

$$\begin{aligned} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y} \dots)^2 &= rk_2 \sum_{i_1=1}^{k_1} (\bar{Y}_{i_1 \dots} - \bar{Y} \dots)^2 + rk_1 \sum_{i_2=1}^{k_2} (\bar{Y}_{\dots i_2} - \bar{Y} \dots)^2 + \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot})^2 \\ &+ r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \dots} - \bar{Y}_{\dots i_2} + \bar{Y} \dots)^2. \end{aligned} \quad (22)$$

**Proof** Using the notation introduced in (19) and (20), we get that

$$\begin{aligned} &\sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y} \dots)^2 \\ &= \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r \left( (\bar{Y}_{i_1 \dots} - \bar{Y} \dots) + (\bar{Y}_{\dots i_2} - \bar{Y} \dots) + (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot}) + (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \dots} - \bar{Y}_{\dots i_2} + \bar{Y} \dots) \right)^2 \\ &= \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (\bar{Y}_{i_1 \dots} - \bar{Y} \dots)^2 + \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (\bar{Y}_{\dots i_2} - \bar{Y} \dots)^2 + \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot})^2 \\ &+ \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \dots} - \bar{Y}_{\dots i_2} + \bar{Y} \dots)^2 + R, \end{aligned}$$

where it can be shown in a similar way as in the proof of Theorem 3.1 that the sum  $R$  of the mixed products is equal to zero.  $\square$

### Remark

- The sum of squares on the left-hand side of (22) can be perceived as a measure for the (total) *variability of the sampling variables*  $\{Y_{i_1 i_2 j}, i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2, j = 1, \dots, r\}$ .
- The first and second sum of squares on the right-hand side of (22) are measures for the variability *between* the levels of the first and second predictor variable, respectively, while the third sum of squares on the right-hand side of (22) is a measure for the variability *within* the pairs of levels  $(i_1, i_2)$  of the two predictor variables, the so-called *residual variance*.
- The fourth sum of squares on the right-hand side of (22) is a measure for the interactions *between* the components of the pairs of levels  $(i_1, i_2)$  of the two predictor variables.
- By similar considerations as in the proof of Theorem 2.5 it can be shown that a suitably normalized version of the residual variance is an unbiased estimator for the variance  $\sigma^2$  of the error terms.
- In particular, it holds that  $\mathbb{E} S^2 = \sigma^2$ , where

$$S^2 = \frac{1}{k_1 k_2 (r-1)} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot})^2.$$

## 3.2 Estimation of Model Parameters

Now we return to the analysis of the linear model with arbitrary design matrix  $\mathbf{X}$  given in (1) – (3). In this section, we assume that

- $\text{rk}(\mathbf{X}) = r < m$ , i.e.,  $\mathbf{X}$  has *not* full column rank and that
- $\boldsymbol{\beta} \in \mathbb{R}^m$  is an *arbitrary*  $m$ -dimensional vector, i.e., at first we do *not* consider any *additional conditions* of type (12) or (18).

### 3.2.1 LS-Estimator for $\beta$

We first recall the following formula for the rank of quadratic matrices.

**Lemma 3.2** *Let  $\mathbf{A}$  be an arbitrary  $n \times n$  matrix. Then it holds that*

$$\text{rk}(\mathbf{A}) = n - \dim \ker(\mathbf{A}), \quad (23)$$

where  $\ker(\mathbf{A}) = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{0}\}$  and  $\dim \ker(\mathbf{A})$  denotes the dimension of  $\ker(\mathbf{A}) \subset \mathbb{R}^n$ .

Moreover, the following property of the rank of products of matrices is useful, which immediately follows from Lemma 3.2.

**Lemma 3.3** *Let  $m, n, r \in \mathbb{N}$  be arbitrary natural numbers and let  $\mathbf{A}, \mathbf{B}$  be arbitrary  $m \times n$  and  $n \times r$  matrices. Then it holds that*

$$\text{rk}(\mathbf{AB}) \leq \min\{\text{rk}(\mathbf{A}), \text{rk}(\mathbf{B})\}. \quad (24)$$

#### Remark

- As we now assume that the design matrix  $\mathbf{X}$  does not have full rank, the  $m \times m$  matrix  $\mathbf{X}^\top \mathbf{X}$  is *not* invertible because Lemma 3.3 implies that  $\text{rk}(\mathbf{X}^\top \mathbf{X}) \leq \text{rk}(\mathbf{X}) < m$ .
- Therefore, the normal equation (2.9), i.e.,

$$\mathbf{X}^\top \mathbf{X} \beta = \mathbf{X}^\top \mathbf{Y}, \quad (25)$$

does not have a uniquely determined solution.

- In order to specify the solution set of (25), we need the notion of the generalized inverse of a matrix.

**Definition** An  $m \times n$  matrix  $\mathbf{A}^-$  is called *generalized inverse* of the  $n \times m$  matrix  $\mathbf{A}$  if

$$\mathbf{A} \mathbf{A}^- \mathbf{A} = \mathbf{A}. \quad (26)$$

In order to show that there always is a solution  $\mathbf{A}^-$  of (26), we use the following general representation formula, which we state without proof at this point.

**Lemma 3.4** *Let  $\mathbf{A}$  be an  $n \times m$  matrix with  $n \geq m$  and  $\text{rk}(\mathbf{A}) = r \leq m$ . Then there are invertible  $n \times n$  and  $m \times m$  matrices  $\mathbf{P}$  and  $\mathbf{Q}$ , such that*

$$\mathbf{PAQ} = \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \quad \text{and} \quad \mathbf{A} = \mathbf{P}^{-1} \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{Q}^{-1}. \quad (27)$$

By use of Lemma 3.4 one can show how solutions  $\mathbf{A}^-$  of (26) can be found.

- Let  $\mathbf{P}$  and  $\mathbf{Q}$  be matrices with the properties considered in Lemma 3.4 and let  $\mathbf{B}$  be an arbitrary  $m \times n$  matrix with

$$\mathbf{B} = \mathbf{Q} \begin{pmatrix} \mathbf{I}_r & \mathbf{R} \\ \mathbf{S} & \mathbf{T} \end{pmatrix} \mathbf{P}, \quad (28)$$

where  $\mathbf{R}$ ,  $\mathbf{S}$  and  $\mathbf{T}$  are arbitrary matrices with dimensions  $r \times (n - r)$ ,  $(m - r) \times r$  and  $(m - r) \times (n - r)$ , respectively.

- Then (27) and (28) imply that

$$\begin{aligned}
\mathbf{ABA} &= \mathbf{P}^{-1} \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{Q}^{-1} \mathbf{Q} \begin{pmatrix} \mathbf{I}_r & \mathbf{R} \\ \mathbf{S} & \mathbf{T} \end{pmatrix} \mathbf{P} \mathbf{P}^{-1} \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{Q}^{-1} \\
&= \mathbf{P}^{-1} \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{I}_r & \mathbf{R} \\ \mathbf{S} & \mathbf{T} \end{pmatrix} \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{Q}^{-1} \\
&= \mathbf{A},
\end{aligned}$$

i.e., the matrix  $\mathbf{B}$  given in (28) is a generalized inverse of  $\mathbf{A}$ .

- Let  $k \in \{r, \dots, m\}$  be an arbitrary natural number. Let

$$\mathbf{R} = \mathbf{0}, \quad \mathbf{S} = \mathbf{0} \quad \text{and} \quad \mathbf{T} = \begin{pmatrix} \mathbf{I}_{k-r} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad (29)$$

then it is  $\text{rk}(\mathbf{B}) = k$ .

Altogether, we obtain the following result.

**Lemma 3.5** *Let  $\mathbf{A}$  be an  $n \times m$  matrix with  $n \geq m$  and  $\text{rk}(\mathbf{A}) = r \leq m$ . Let furthermore  $\mathbf{B}$  be the  $m \times n$  matrix given in (28) – (29), for each  $k \in \{r, \dots, m\}$ . Then it holds that  $\text{rk}(\mathbf{B}) = k$  and  $\mathbf{A}^- = \mathbf{B}$  is a solution of (26).*

Moreover, the following properties of the generalized inverse are useful.

**Lemma 3.6**

- Let  $\mathbf{A}$  be an arbitrary  $n \times m$  matrix with  $n \geq m$  and let  $(\mathbf{A}^\top \mathbf{A})^-$  be a generalized inverse of the symmetric  $m \times m$  matrix  $\mathbf{A}^\top \mathbf{A}$ .
- Then the transposed matrix  $((\mathbf{A}^\top \mathbf{A})^-)^\top$  is a generalized inverse of  $\mathbf{A}^\top \mathbf{A}$  as well.
- Besides, it holds that

$$\mathbf{A}^\top \mathbf{A} (\mathbf{A}^\top \mathbf{A})^- \mathbf{A}^\top = \mathbf{A}^\top. \quad (30)$$

**Proof**

- By definition of the generalized inverse, we have  $\mathbf{A}^\top \mathbf{A} (\mathbf{A}^\top \mathbf{A})^- \mathbf{A}^\top \mathbf{A} = \mathbf{A}^\top \mathbf{A}$ .
- From this equation and from the symmetry of the matrix  $\mathbf{A}^\top \mathbf{A}$  it follows that

$$\mathbf{A}^\top \mathbf{A} = (\mathbf{A}^\top \mathbf{A})^\top = \left( \mathbf{A}^\top \mathbf{A} (\mathbf{A}^\top \mathbf{A})^- \mathbf{A}^\top \mathbf{A} \right)^\top = \mathbf{A}^\top \mathbf{A} \left( (\mathbf{A}^\top \mathbf{A})^- \right)^\top \mathbf{A}^\top \mathbf{A},$$

i.e., the transposed matrix  $((\mathbf{A}^\top \mathbf{A})^-)^\top$  is a generalized inverse of  $\mathbf{A}^\top \mathbf{A}$  as well.

- In order to prove (30), the second part of the statement, we consider the matrix

$$\mathbf{B} = \mathbf{A}^\top \mathbf{A} (\mathbf{A}^\top \mathbf{A})^- \mathbf{A}^\top - \mathbf{A}^\top.$$

- Then it holds that

$$\begin{aligned}
\mathbf{B}\mathbf{B}^\top &= \left(\mathbf{A}^\top\mathbf{A}(\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top - \mathbf{A}^\top\right)\left(\mathbf{A}^\top\mathbf{A}(\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top - \mathbf{A}^\top\right)^\top \\
&= \mathbf{A}^\top\mathbf{A}(\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top\mathbf{A}\left((\mathbf{A}^\top\mathbf{A})^{-1}\right)^\top\mathbf{A}^\top\mathbf{A} \\
&\quad - \mathbf{A}^\top\mathbf{A}(\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top\mathbf{A} - \mathbf{A}^\top\mathbf{A}\left((\mathbf{A}^\top\mathbf{A})^{-1}\right)^\top\mathbf{A}^\top\mathbf{A} + \mathbf{A}^\top\mathbf{A} \\
&= \mathbf{A}^\top\mathbf{A} - \mathbf{A}^\top\mathbf{A} - \mathbf{A}^\top\mathbf{A} + \mathbf{A}^\top\mathbf{A} = \mathbf{0}.
\end{aligned}$$

- Therefore, we get that  $\mathbf{B} = \mathbf{0}$ . □

By use of the generalized inverse  $(\mathbf{X}^\top\mathbf{X})^-$  of  $\mathbf{X}^\top\mathbf{X}$  and its properties (considered in Lemma 3.6), the solution set of the normal equation (25) can be specified.

**Theorem 3.5** *The general solution  $\boldsymbol{\beta}$  of the normal equation  $\mathbf{X}^\top\mathbf{X}\boldsymbol{\beta} = \mathbf{X}^\top\mathbf{Y}$  has the form*

$$\boldsymbol{\beta} = (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{Y} + (\mathbf{I}_m - (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{X})\mathbf{z}, \quad (31)$$

where  $(\mathbf{X}^\top\mathbf{X})^-$  is an arbitrary solution of

$$\mathbf{X}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{X} = \mathbf{X}^\top\mathbf{X} \quad (32)$$

and  $\mathbf{z} \in \mathbb{R}^m$  is an arbitrary  $m$ -dimensional vector.

### Proof

- By plugging (31) into the left-hand side of the normal equation (25), one sees
  - that for each  $\mathbf{z} \in \mathbb{R}^m$  equation (31) gives a solution of (25),
  - as it holds that

$$\begin{aligned}
\mathbf{X}^\top\mathbf{X}\boldsymbol{\beta} &= \mathbf{X}^\top\mathbf{X}\left((\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{Y} + (\mathbf{I}_m - (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{X})\mathbf{z}\right) \\
&= \mathbf{X}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{Y} = \mathbf{X}^\top\mathbf{Y},
\end{aligned}$$

where the last equality follows from Lemma 3.6.

- Let now  $\tilde{\boldsymbol{\beta}}$  be an arbitrary solution and  $\boldsymbol{\beta}$  be a solution of the form (31) of equation (25).
  - Then subtraction on each side of (25) yields

$$\mathbf{X}^\top\mathbf{X}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}) = \mathbf{0}. \quad (33)$$

- Hence, for a  $\mathbf{z} \in \mathbb{R}^m$  it holds that

$$\begin{aligned}
\tilde{\boldsymbol{\beta}} &= \boldsymbol{\beta} - (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) \stackrel{(31)}{=} (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{Y} + (\mathbf{I}_m - (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{X})\mathbf{z} - (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) \\
&= (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{Y} + (\mathbf{I}_m - (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{X})(\mathbf{z} - (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})) + (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{X}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}) \\
&\stackrel{(33)}{=} (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{Y} + (\mathbf{I}_m - (\mathbf{X}^\top\mathbf{X})^- \mathbf{X}^\top\mathbf{X})(\mathbf{z} - (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})).
\end{aligned}$$

- This means that  $\tilde{\boldsymbol{\beta}}$  is a solution of the form (31) of equation (25) as well. □

**Example** (one-factor analysis of variance)

- *Recall:* In the reparametrized model of one-factor analysis of variance (cf. case 2 of the example considered in Section 3.1.2) the design matrix is given by the  $n \times (k+1)$  matrix

$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 1 & 0 & 0 & \dots & 0 & 0 \\ 1 & 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 0 & 1 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & 0 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 0 & 0 & 0 & \dots & 0 & 1 \end{pmatrix}, \quad (34)$$

and the parameter vector  $\boldsymbol{\beta}$  is given by  $\boldsymbol{\beta} = (\mu, \alpha_1, \dots, \alpha_k)^\top$ .

- One can easily see that in this case

$$\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} n & n_1 & n_2 & n_3 & \dots & n_{k-1} & n_k \\ n_1 & n_1 & 0 & 0 & \dots & 0 & 0 \\ n_2 & 0 & n_2 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ n_k & 0 & 0 & 0 & \dots & 0 & n_k \end{pmatrix} \quad (35)$$

and that a generalized inverse of  $\mathbf{X}^\top \mathbf{X}$  is given by

$$(\mathbf{X}^\top \mathbf{X})^- = \begin{pmatrix} \frac{1}{n} & 0 & 0 & 0 & \dots & 0 & 0 \\ -\frac{1}{n} & \frac{1}{n_1} & 0 & 0 & \dots & 0 & 0 \\ -\frac{1}{n} & 0 & \frac{1}{n_2} & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ -\frac{1}{n} & 0 & 0 & 0 & \dots & 0 & \frac{1}{n_k} \end{pmatrix}. \quad (36)$$

- Therefore, the normal equation (25), i.e.,  $\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{Y}$ , has the following form:

$$n\mu + \sum_{i=1}^k n_i \alpha_i = Y_{..}, \quad n_i \mu + n_i \alpha_i = Y_{i.}, \quad \forall i = 1, \dots, k.$$

- If we only consider solutions of this system of equations which are in the *restricted parameter space*  $\Theta \subset \mathbb{R}^{k+1}$ , where

$$\Theta = \left\{ \boldsymbol{\beta} = (\mu, \alpha_1, \dots, \alpha_k) : \sum_{i=1}^k n_i \alpha_i = 0 \right\}, \quad (37)$$



we obtain the (uniquely determined) solution  $\widehat{\boldsymbol{\beta}} = (\widehat{\mu}, \widehat{\alpha}_1, \dots, \widehat{\alpha}_k)$  with

$$\widehat{\mu} = \bar{Y}_{..}, \quad \widehat{\alpha}_i = \bar{Y}_{i.} - \bar{Y}_{..}, \quad \forall i = 1, \dots, k. \quad (38)$$

- One can easily see that the solution  $\widehat{\boldsymbol{\beta}}$  of the normal equation (25) which is given in (38)
  - has the form  $\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y}$ , where the generalized inverse  $(\mathbf{X}^\top \mathbf{X})^-$  is given by (36), and
  - is an unbiased estimator for  $\boldsymbol{\beta} = (\mu, \alpha_1, \dots, \alpha_k)$  with respect to the restricted parameter space  $\Theta$ , which has the form given in (37).
- Without the additional condition considered in (37), there is no LS-estimator which is unbiased at the same time, cf. Theorem 3.8.

Now we consider the linear model with general design matrix  $\mathbf{X}$  again, which is given in (1) – (3). In particular, we consider the solutions of the normal equation (25) discussed in Theorem 3.5 and show that the mean squared error  $e(\boldsymbol{\beta})$  given in (2.8) is minimized for  $\mathbf{z} = \mathbf{o}$ .

**Theorem 3.6** *Let  $(\mathbf{X}^\top \mathbf{X})^-$  be a generalized inverse of  $\mathbf{X}^\top \mathbf{X}$ . Then the sample function*

$$\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y} \quad (39)$$

*minimizes the mean squared error  $e(\boldsymbol{\beta})$ , i.e.,  $\bar{\boldsymbol{\beta}}$  is an LS-estimator for  $\boldsymbol{\beta}$ .*

**Proof**

- For each  $m$ -dimensional vector  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)^\top$  it holds that

$$\begin{aligned} ne(\boldsymbol{\beta}) &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}} + \mathbf{X}(\bar{\boldsymbol{\beta}} - \boldsymbol{\beta}))^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}} + \mathbf{X}(\bar{\boldsymbol{\beta}} - \boldsymbol{\beta})) \\ &= (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}) + (\bar{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\bar{\boldsymbol{\beta}} - \boldsymbol{\beta}) \geq (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}) = ne(\bar{\boldsymbol{\beta}}) \end{aligned}$$

- because

$$(\bar{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\bar{\boldsymbol{\beta}} - \boldsymbol{\beta}) = (\mathbf{X}(\bar{\boldsymbol{\beta}} - \boldsymbol{\beta}))^\top (\mathbf{X}(\bar{\boldsymbol{\beta}} - \boldsymbol{\beta})) \geq 0$$

and

$$(\bar{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}) = (\bar{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top (\mathbf{X}^\top - \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top) \mathbf{Y} = 0,$$

where the last equality follows from Lemma 3.6. □

### 3.2.2 Expectation Vector and Covariance Matrix of the LS-Estimator $\bar{\boldsymbol{\beta}}$

The model assumptions (3) for  $\varepsilon_1, \dots, \varepsilon_n$  and the general calculation rules for the expectation and covariance of real-valued random variables imply that the expectation vector and the covariance matrix of the LS-estimator  $\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y}$  have the following form.

**Theorem 3.7** *It holds that*

$$\mathbb{E} \bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} \quad (40)$$

and

$$\text{Cov} \bar{\boldsymbol{\beta}} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X} ((\mathbf{X}^\top \mathbf{X})^-)^\top. \quad (41)$$

**Proof**

- It follows from  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  and  $\mathbb{E}\boldsymbol{\varepsilon} = \mathbf{o}$  that

$$\mathbb{E}\bar{\boldsymbol{\beta}} = \mathbb{E}\left((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}\right) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}\mathbf{Y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}.$$

- Moreover, it holds that

$$\begin{aligned} \text{Cov}(\bar{\beta}_i, \bar{\beta}_j) &= \text{Cov}\left(\sum_{\ell=1}^n ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)_{i\ell} Y_\ell, \sum_{r=1}^n ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)_{jr} Y_r\right) \\ &= \sum_{\ell=1}^n \sum_{r=1}^n ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)_{i\ell} ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)_{jr} \text{Cov}(Y_\ell, Y_r) \\ &= \sigma^2 \sum_{\ell=1}^n ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)_{i\ell} ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)_{j\ell} \\ &= \sigma^2 \sum_{\ell=1}^n ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)_{i\ell} (\mathbf{X}((\mathbf{X}^\top \mathbf{X})^{-1})^\top)_{\ell j} \\ &= \sigma^2 ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} ((\mathbf{X}^\top \mathbf{X})^{-1})^\top)_{ij} \end{aligned}$$

□

for arbitrary  $i, j \in \{1, \dots, m\}$ .

Together with Lemma 3.3, Theorems 3.5 and 3.7 imply that there is *no* LS-estimator for  $\boldsymbol{\beta}$  which is additionally unbiased. In particular, the LS-estimator  $\bar{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$  given in (39) is biased.

**Theorem 3.8** *If  $\text{rk}(\mathbf{X}) < m$ , there is no unbiased LS-estimator for  $\boldsymbol{\beta}$ .*

**Proof**

- Due to  $\text{rk}(\mathbf{X}) < m$ , it follows from Lemma 3.3 that  $\text{rk}(\mathbf{X}^\top \mathbf{X}) < m$  and

$$\text{rk}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}) < m.$$

- Hence, there is a  $\boldsymbol{\beta} \neq \mathbf{o}$  with  $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} = \mathbf{o}$ , i.e., the equation

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta} \tag{42}$$

does *not* hold for each  $\boldsymbol{\beta} \in \mathbb{R}^m$ .

- So because of (40), the LS-estimator  $\bar{\boldsymbol{\beta}}$  für  $\boldsymbol{\beta}$  given in (39) is biased.

- As (42) does not hold for each  $\boldsymbol{\beta} \in \mathbb{R}^m$ , one additionally obtains that for each arbitrary but fixed  $\mathbf{z} \in \mathbb{R}^m$  the equation

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}(\boldsymbol{\beta} - \mathbf{z}) = \boldsymbol{\beta} - \mathbf{z}$$

or equivalently

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} + (\mathbf{I}_m - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X})\mathbf{z} = \boldsymbol{\beta}$$

does not hold for each  $\boldsymbol{\beta} \in \mathbb{R}^m$ .

- Due to Theorem 3.5, this means that there is no LS-estimator for  $\boldsymbol{\beta}$  which is additionally unbiased. □

### 3.2.3 Estimable Functions

- In Section 3.2.2 we have shown that if the design matrix  $\mathbf{X}$  does not have full rank, there is no unbiased LS-estimator for  $\boldsymbol{\beta}$  in the linear model without additional conditions.
- Hence, instead of the vector  $\boldsymbol{\beta}$ , one considers a class of (real-valued) linear functions  $\mathbf{a}^\top \boldsymbol{\beta}$  of the parameter vector  $\boldsymbol{\beta}$ , for which unbiased LS-estimators can be constructed.
- In other words: Instead of the (vectorial) linear transformation  $\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$  of the random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  one considers a class of (real-valued) linear functions  $\mathbf{c}^\top \mathbf{Y}$  of  $\mathbf{Y}$ , which can be perceived as estimators for  $\mathbf{a}^\top \boldsymbol{\beta}$ .
- This leads to the following conception.

#### Definition

- Let  $\mathbf{a} = (a_1, \dots, a_m)^\top \in \mathbb{R}^m$  be an arbitrary  $m$ -dimensional vector.
- The linear function  $\mathbf{a}^\top \boldsymbol{\beta}$  of the parameter vector  $\boldsymbol{\beta}$  is called *estimable without bias* or an *estimable function* if there is an  $n$ -dimensional vector  $\mathbf{c} = (c_1, \dots, c_n)^\top$  such that

$$\mathbb{E}(\mathbf{c}^\top \mathbf{Y}) = \mathbf{a}^\top \boldsymbol{\beta}, \quad \forall \boldsymbol{\beta} \in \mathbb{R}^m. \quad (43)$$

#### Example (one-factor analysis of variance)

- For the reparametrized model of one-factor analysis of variance with parameter vector  $\boldsymbol{\beta} = (\mu, \alpha_1, \dots, \alpha_k)^\top \in \mathbb{R}^{k+1}$  one can show that for example  $\alpha_1 - \alpha_2$  is an estimable function as defined by (43).
- This is true because for

$$\mathbf{a}^\top = (0, 1, -1, 0, \dots, 0) \quad \text{and} \quad \mathbf{c}^\top = \underbrace{(0, \dots, 0)}_{n_1-1}, 1, -1, 0, \dots, 0$$

it holds that

$$\mathbb{E}(\mathbf{c}^\top \mathbf{Y}) = \mathbb{E}(Y_{1n_1} - Y_{21}) = (\mu + \alpha_1) - (\mu + \alpha_2) = \alpha_1 - \alpha_2 = \mathbf{a}^\top \boldsymbol{\beta}$$

for each  $\boldsymbol{\beta} = (\mu, \alpha_1, \dots, \alpha_k)^\top \in \mathbb{R}^{k+1}$ .

- In a similar way it can be shown that  $\mu + \alpha_i$  and  $\alpha_i - \alpha_{i'}$  are estimable functions of  $\boldsymbol{\beta}$  for  $i = 1, \dots, k$  and  $i, i' = 1, \dots, k$  with  $i \neq i'$ , respectively.

#### Example (two-factor analysis of variance with balanced subsamples)

- For the model of two-factor analysis of variance with balanced subsamples, introduced in Section 3.1.3, the normal equation (25) has the following form:

$$\begin{aligned} rk_1 k_2 \mu + rk_2 \sum_{i_1=1}^{k_1} \alpha_{i_1}^{(1)} + rk_1 \sum_{i_2=1}^{k_2} \alpha_{i_2}^{(2)} + r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \alpha_{i_1 i_2} &= Y_{..} \\ rk_2 \mu + rk_2 \alpha_{i_1}^{(1)} + r \sum_{i_2=1}^{k_2} \alpha_{i_2}^{(2)} + r \sum_{i_2=1}^{k_2} \alpha_{i_1 i_2} &= Y_{i_1 \cdot} \quad \forall i_1 = 1, \dots, k_1 \\ rk_1 \mu + r \sum_{i_1=1}^{k_1} \alpha_{i_1}^{(1)} + rk_1 \alpha_{i_2}^{(2)} + r \sum_{i_1=1}^{k_1} \alpha_{i_1 i_2} &= Y_{\cdot i_2} \quad \forall i_2 = 1, \dots, k_2 \\ r\mu + r\alpha_{i_1}^{(1)} + r\alpha_{i_2}^{(2)} + r\alpha_{i_1 i_2} &= Y_{i_1 i_2} \quad \forall i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2 \end{aligned}$$

- In consideration of the additional condition (18), this system of equations can be solved uniquely. In other words: If only parameter vectors

$$\boldsymbol{\beta} = (\mu, \alpha_1^{(1)}, \dots, \alpha_{k_1}^{(1)}, \alpha_1^{(2)}, \dots, \alpha_{k_2}^{(2)}, \alpha_{11}, \dots, \alpha_{k_1 k_2})^\top$$

from the restricted parameter space

$$\Theta = \left\{ \boldsymbol{\beta} : \sum_{i_1=1}^{k_1} \alpha_{i_1}^{(1)} = \sum_{i_2=1}^{k_2} \alpha_{i_2}^{(2)} = \sum_{i_1=1}^{k_1} \alpha_{i_1 i_2} = \sum_{i_2=1}^{k_2} \alpha_{i_1 i_2} = 0 \right\}$$

are considered, one obtains the unique solution

$$\hat{\boldsymbol{\beta}} = (\hat{\mu}, \hat{\alpha}_1^{(1)}, \dots, \hat{\alpha}_{k_1}^{(1)}, \hat{\alpha}_1^{(2)}, \dots, \hat{\alpha}_{k_2}^{(2)}, \hat{\alpha}_{11}, \dots, \hat{\alpha}_{k_1 k_2})^\top \quad (44)$$

of the normal equation, where

$$\hat{\mu} = \bar{Y} \dots, \quad \hat{\alpha}_{i_1}^{(1)} = \bar{Y}_{i_1 \dots} - \bar{Y} \dots, \quad \hat{\alpha}_{i_2}^{(2)} = \bar{Y}_{\cdot i_2} - \bar{Y} \dots, \quad \hat{\alpha}_{i_1 i_2} = \bar{Y} \dots + \bar{Y}_{i_1 i_2} - \bar{Y}_{i_1 \dots} - \bar{Y}_{\cdot i_2}. \quad (45)$$

for arbitrary  $i_1 = 1, \dots, k_1$ ,  $i_2 = 1, \dots, k_2$ .

- It can be shown that the solution  $\hat{\boldsymbol{\beta}}$  of the normal equation given in (44) – (45) has the form

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y},$$

where  $(\mathbf{X}^\top \mathbf{X})^-$  is a generalized inverse of  $\mathbf{X}^\top \mathbf{X}$  and  $\mathbf{X}$  is the design matrix of the model of two-factor analysis of variance with balanced subsamples.

- *Remark:* The sample function  $\hat{\boldsymbol{\beta}}$  given in (44) – (45) was already discussed in Theorem 3.3, where we have shown that  $\hat{\boldsymbol{\beta}}$  is an unbiased estimator for  $\boldsymbol{\beta}$  with respect to the parameter space  $\Theta$ .
- Furthermore, one can show that the linear functions  $\mu + \alpha_{i_1}^{(1)} + \alpha_{i_2}^{(2)} + \alpha_{i_1 i_2}$  of the parameter vector  $\boldsymbol{\beta}$  are estimable without bias as defined in (43) (without taking into account the additional conditions (18)) for arbitrary  $i_1 = 1, \dots, k_1$ ,  $i_2 = 1, \dots, k_2$ .
- In the model of two-factor analysis of variance without interactions, i.e.,  $\alpha_{i_1 i_2} = 0$  for arbitrary  $i_1 = 1, \dots, k_1$ ,  $i_2 = 1, \dots, k_2$ , also the linear functions  $\alpha_{i_1}^{(1)} - \alpha_{i'_1}^{(1)}$  and  $\alpha_{i_2}^{(2)} - \alpha_{i'_2}^{(2)}$  are estimable without bias for arbitrary  $i_1, i'_1 = 1, \dots, k_1$  with  $i_1 \neq i'_1$  and for arbitrary  $i_2, i'_2 = 1, \dots, k_2$  with  $i_2 \neq i'_2$ , respectively.

The following lemma, which is an extension of Lemma 3.6, is needed to derive two general *criteria for linear functions*  $\mathbf{a}^\top \boldsymbol{\beta}$  of the parameter vector  $\boldsymbol{\beta}$  to be estimable without bias.

**Lemma 3.7** *Let  $(\mathbf{X}^\top \mathbf{X})^-$  be a generalized inverse of  $\mathbf{X}^\top \mathbf{X}$ . Then it holds that*

$$\mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X} = \mathbf{X}. \quad (46)$$

**Proof** In Lemma 3.6 we have shown that

- the transposed matrix  $((\mathbf{X}^\top \mathbf{X})^-)^\top$  is a generalized inverse of  $\mathbf{X}^\top \mathbf{X}$  as well and that  $\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top = \mathbf{X}^\top$ .
- Hence, it is  $\mathbf{X}^\top \mathbf{X}((\mathbf{X}^\top \mathbf{X})^-)^\top \mathbf{X}^\top = \mathbf{X}^\top$ .
- This leads to (46) by swapping columns and rows.  $\square$

**Theorem 3.9** *Let  $\mathbf{a} = (a_1, \dots, a_m)^\top \in \mathbb{R}^m$  be an arbitrary vector. The linear function  $\mathbf{a}^\top \boldsymbol{\beta}$  of the parameter vector  $\boldsymbol{\beta}$  is estimable without bias if and only if one of the following conditions is fulfilled:*

1. There is a  $\mathbf{c} \in \mathbb{R}^n$ , such that

$$\mathbf{a}^\top = \mathbf{c}^\top \mathbf{X}. \quad (47)$$

2. The vector  $\mathbf{a}$  fulfills the following system of equations:

$$\mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = \mathbf{a}^\top. \quad (48)$$

### Proof

- Let  $\mathbf{a}^\top \boldsymbol{\beta}$  be an estimable function of the parameter vector  $\boldsymbol{\beta}$ .
  - Then it follows from (43) that

$$\mathbf{a}^\top \boldsymbol{\beta} = \mathbb{E}(\mathbf{c}^\top \mathbf{Y}) = \mathbf{c}^\top \mathbb{E} \mathbf{Y} = \mathbf{c}^\top \mathbf{X} \boldsymbol{\beta} \quad \text{and hence} \quad (\mathbf{c}^\top \mathbf{X} - \mathbf{a}^\top) \boldsymbol{\beta} = 0$$

for each  $\boldsymbol{\beta} \in \mathbb{R}^m$ .

- Therefore, we get that  $\mathbf{a}^\top = \mathbf{c}^\top \mathbf{X}$ .
- Conversely, let  $\mathbf{c}$  be a vector satisfying condition (47).
  - Then

$$\mathbb{E}(\mathbf{c}^\top \mathbf{Y}) = \mathbf{c}^\top \mathbb{E} \mathbf{Y} = \mathbf{c}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{a}^\top \boldsymbol{\beta}, \quad \forall \boldsymbol{\beta} \in \mathbb{R}^m.$$

- Thus, also the sufficiency of condition (47) is proved.
- In order to show the necessity of condition (48), we use the result of Lemma 3.7.
  - Let  $\mathbf{a}^\top \boldsymbol{\beta}$  be an estimable function of the parameter vector  $\boldsymbol{\beta}$ .
  - Then it follows from (47) and (46) that

$$\mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = \mathbf{c}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = \mathbf{c}^\top \mathbf{X} = \mathbf{a}^\top.$$

- In order to show the sufficiency of condition (48), one only has to observe
  - that (48) implies  $\mathbf{a}^\top = \mathbf{c}^\top \mathbf{X}$  for  $\mathbf{c}^\top = \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top$ .
  - Now the first part of the statement implies that  $\mathbf{a}^\top \boldsymbol{\beta}$  is estimable without bias. □

### Remark

- If the design matrix  $\mathbf{X}$  has full rank, i.e.,  $(\mathbf{X}^\top \mathbf{X})^{-} = (\mathbf{X}^\top \mathbf{X})^{-1}$ , then condition (48) is obviously fulfilled for each  $\mathbf{a} \in \mathbb{R}^m$ .
- In this case, *every* linear function of the parameter vector  $\boldsymbol{\beta}$  is estimable without bias, which has already been shown in Theorem 2.2.

In the case that the design matrix  $\mathbf{X} = (x_{ij})$  does not have full rank, we show

- how the second part of the statement in Theorem 3.9 implies that the following linear functions  $\mathbf{a}^\top \boldsymbol{\beta}$  of  $\boldsymbol{\beta}$  are estimable without bias.
- In doing so, the vector (of weights)  $\mathbf{c}$  of the linear unbiased estimator  $\mathbf{c}^\top \mathbf{Y}$  for  $\mathbf{a}^\top \boldsymbol{\beta}$  can be chosen as in the proof of Theorem 3.9, i.e.,

$$\mathbf{c}^\top = \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top. \quad (49)$$

**Theorem 3.10** *The following linear functions of the parameter vector  $\boldsymbol{\beta}$  are estimable without bias:*

1. the components  $\sum_{j=1}^m x_{1j} \beta_j, \dots, \sum_{j=1}^m x_{nj} \beta_j$  of the expectation vector  $\mathbb{E} \mathbf{Y} = \mathbf{X} \boldsymbol{\beta}$ ,
2. each linear function of estimable functions,

3. the components  $\beta'_1, \dots, \beta'_m$  of the so-called projected parameter vector  $\beta' = (\beta'_1, \dots, \beta'_m)^\top$ , where

$$\beta' = (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} \beta. \quad (50)$$

**Proof**

- Let  $\mathbf{a}_i = (x_{i1}, \dots, x_{im})^\top$  for each  $i \in \{1, \dots, n\}$ . Then it holds that

$$\begin{pmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_n^\top \end{pmatrix} (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = \mathbf{X} = \begin{pmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_n^\top \end{pmatrix},$$

where the last but one equality follows from Lemma 3.7. Now Theorem 3.9 implies that every linear combination  $\sum_{j=1}^m x_{1j} \beta_j, \dots, \sum_{j=1}^m x_{nj} \beta_j$  is an estimable function.

- In order to prove the second part of the statement, we consider a (finite) family  $\mathbf{a}_1^\top \beta, \dots, \mathbf{a}_s^\top \beta$  of  $s$  estimable functions, which we write in the form  $\mathbf{A} \beta$ , where  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_s)^\top$  is an  $s \times m$ -dimensional matrix and  $s \in \mathbb{N}$  is an arbitrary natural number.
- Theorem 3.9 implies that

$$\mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = \mathbf{A}.$$

- Thus, for each  $s$ -dimensional vector  $\mathbf{b} = (b_1, \dots, b_s)^\top \in \mathbb{R}^s$  it holds that

$$\mathbf{b}^\top \mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = \mathbf{b}^\top \mathbf{A}.$$

- Therefore, it follows by use of Theorem 3.9 that the linear function  $\mathbf{b}^\top \mathbf{A} \beta$  of the estimable functions  $\mathbf{A} \beta$  is an estimable function itself.
- In the third part of the statement the family  $\mathbf{A} \beta$  of linear functions of the parameter vector  $\beta$  is considered, where the  $m \times m$  matrix  $\mathbf{A}$  is given by  $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X}$ .

- Hence,

$$\mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} = \mathbf{A},$$

where the last but one equality follows from the definition of the generalized inverse in (26).

- Now Theorem 3.9 implies that the components  $\beta'_1, \dots, \beta'_m$  of the projected parameter vector

$$\beta' = \mathbf{A} \beta = (\mathbf{X}^\top \mathbf{X})^{-} \mathbf{X}^\top \mathbf{X} \beta$$

are estimable functions. □

### 3.2.4 Best Linear Unbiased Estimator; Gauss–Markov Theorem

- In this section we show how BLUE-estimators for estimable functions of the parameter vector  $\beta$  can be constructed.
- *Recall:* A linear unbiased estimator is called BLUE-estimator if there is no linear unbiased estimator with lower variance (BLUE = best linear unbiased estimator).
- In the theory of linear models, the following result is called *Gauss–Markov theorem*.

**Theorem 3.11**

- Let  $\mathbf{a}^\top \boldsymbol{\beta}$  be an estimable function of the parameter vector  $\boldsymbol{\beta}$ , let  $(\mathbf{X}^\top \mathbf{X})^-$  be an arbitrary generalized inverse of the  $m \times m$  matrix  $\mathbf{X}^\top \mathbf{X}$  and let  $\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y}$ .
- Then  $\mathbf{a}^\top \bar{\boldsymbol{\beta}}$  is a BLUE-estimator for  $\mathbf{a}^\top \boldsymbol{\beta}$ , where

$$\text{Var}(\mathbf{a}^\top \bar{\boldsymbol{\beta}}) = \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{a}. \quad (51)$$

**Proof**

- First, we show that  $\mathbf{a}^\top \bar{\boldsymbol{\beta}}$  is a linear unbiased estimator for  $\mathbf{a}^\top \boldsymbol{\beta}$ .

– It is clear that

$$\mathbf{a}^\top \bar{\boldsymbol{\beta}} = \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y}$$

is a linear function of the random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)$ .

– As we assume that  $\mathbf{a}^\top \boldsymbol{\beta}$  is an estimable function of the parameter vector  $\boldsymbol{\beta}$ , Theorem 3.9 implies that there is a  $\mathbf{c} \in \mathbb{R}^n$  such that

$$\mathbf{a}^\top = \mathbf{c}^\top \mathbf{X}. \quad (52)$$

– Hence, it holds that

$$\mathbb{E}(\mathbf{a}^\top \bar{\boldsymbol{\beta}}) = \mathbf{c}^\top \mathbf{X} \mathbb{E} \bar{\boldsymbol{\beta}} = \mathbf{c}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbb{E} \mathbf{Y} = \mathbf{c}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{c}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{a}^\top \boldsymbol{\beta}$$

for each  $\boldsymbol{\beta} \in \mathbb{R}^m$ , where the last but one equality follows from Lemma 3.7.

– Therefore, we have shown that  $\mathbf{a}^\top \bar{\boldsymbol{\beta}}$  is a linear unbiased estimator for  $\mathbf{a}^\top \boldsymbol{\beta}$ .

- It follows from the calculation rules for the variance of random variables (cf. Theorem WR-4.13) that

$$\text{Var}(\mathbf{a}^\top \bar{\boldsymbol{\beta}}) = \text{Var}\left(\sum_{i=1}^m a_i \bar{\beta}_i\right) = \sum_{i=1}^m \sum_{j=1}^m a_i a_j \text{Cov}(\bar{\beta}_i, \bar{\beta}_j).$$

– Moreover, we have shown in Theorem 3.7 that

$$\text{Cov}(\bar{\beta}_i, \bar{\beta}_j) = \sigma^2 \left( (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X} ((\mathbf{X}^\top \mathbf{X})^-)^\top \right)_{ij}.$$

– This leads to

$$\begin{aligned} \text{Var}(\mathbf{a}^\top \bar{\boldsymbol{\beta}}) &= \sigma^2 \sum_{i=1}^m \sum_{j=1}^m a_i a_j \left( (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X} ((\mathbf{X}^\top \mathbf{X})^-)^\top \right)_{ij} \\ &= \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X} ((\mathbf{X}^\top \mathbf{X})^-)^\top \mathbf{a} \\ &= \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X} ((\mathbf{X}^\top \mathbf{X})^-)^\top \mathbf{X}^\top \mathbf{c} \\ &= \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{c}, \end{aligned}$$

where the last but one equality follows from (52) and the last equality follows from Lemma 3.6.

– Applying (52) again yields the variance formula (51).

- It remains to show that the estimator  $\mathbf{a}^\top \bar{\boldsymbol{\beta}}$  has minimum variance in the class of all linear unbiased estimators for  $\mathbf{a}^\top \boldsymbol{\beta}$ .

– Let  $\mathbf{b} \in \mathbb{R}^n$ , such that  $\mathbf{b}^\top \mathbf{Y}$  is a linear unbiased estimator for  $\mathbf{a}^\top \boldsymbol{\beta}$ . Then it holds that

$$\mathbf{a}^\top \boldsymbol{\beta} = \mathbb{E}(\mathbf{b}^\top \mathbf{Y}) = \mathbf{b}^\top \mathbf{X} \boldsymbol{\beta}, \quad \forall \boldsymbol{\beta} \in \mathbb{R}^m$$

and thus

$$\mathbf{b}^\top \mathbf{X} = \mathbf{a}^\top. \quad (53)$$

– For the covariance of  $\mathbf{a}^\top \bar{\boldsymbol{\beta}}$  and  $\mathbf{b}^\top \mathbf{Y}$  it holds that

$$\text{Cov}(\mathbf{a}^\top \bar{\boldsymbol{\beta}}, \mathbf{b}^\top \mathbf{Y}) = \text{Cov}(\mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}, \mathbf{b}^\top \mathbf{Y}) = \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{b} = \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a},$$

where the last equality follows from (53).

– This result and the variance formula (51) imply that

$$\begin{aligned} 0 &\leq \text{Var}(\mathbf{a}^\top \bar{\boldsymbol{\beta}} - \mathbf{b}^\top \mathbf{Y}) = \text{Var}(\mathbf{a}^\top \bar{\boldsymbol{\beta}}) + \text{Var}(\mathbf{b}^\top \mathbf{Y}) - 2 \text{Cov}(\mathbf{a}^\top \bar{\boldsymbol{\beta}}, \mathbf{b}^\top \mathbf{Y}) \\ &= \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} + \text{Var}(\mathbf{b}^\top \mathbf{Y}) - 2 \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} \\ &= \text{Var}(\mathbf{b}^\top \mathbf{Y}) - \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} = \text{Var}(\mathbf{b}^\top \mathbf{Y}) - \text{Var}(\mathbf{a}^\top \bar{\boldsymbol{\beta}}). \end{aligned}$$

□

### Remark

- In the proof of Theorem 3.11 it has never explicitly been used that  $\text{rk}(\mathbf{X}) < m$ .
- In other words: If the design matrix  $\mathbf{X}$  has full rank, i.e.,  $\text{rk}(\mathbf{X}) = m$ , then  $\mathbf{a}^\top \boldsymbol{\beta}$  is estimable without bias for *each*  $m$ -dimensional vector  $\mathbf{a}^\top \in \mathbb{R}^m$  and  $\mathbf{a}^\top \hat{\boldsymbol{\beta}} = \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$  is a BLUE-estimator for  $\mathbf{a}^\top \boldsymbol{\beta}$ .

The following *invariance property* of the generalized inverse  $(\mathbf{X}^\top \mathbf{X})^-$  of  $\mathbf{X}^\top \mathbf{X}$  implies that the BLUE-estimator  $\mathbf{a}^\top \bar{\boldsymbol{\beta}}$ , considered in Theorem 3.11, does *not* depend on the specific choice of  $(\mathbf{X}^\top \mathbf{X})^-$ .

**Lemma 3.8** *Let  $\mathbf{A}$  and  $\mathbf{A}'$  be arbitrary generalized inverses of the matrix  $\mathbf{X}^\top \mathbf{X}$ . Then it holds that*

$$\mathbf{X} \mathbf{A} \mathbf{X}^\top = \mathbf{X} \mathbf{A}' \mathbf{X}^\top. \quad (54)$$

### Proof

- In Lemma 3.7 we have shown that

$$\mathbf{X} \mathbf{A} \mathbf{X}^\top \mathbf{X} = \mathbf{X} = \mathbf{X} \mathbf{A}' \mathbf{X}^\top \mathbf{X} \quad (55)$$

for arbitrary generalized inverses  $\mathbf{A}$  and  $\mathbf{A}'$  of the matrix  $\mathbf{X}^\top \mathbf{X}$ .

- If this chain of equations is multiplied by  $\mathbf{A} \mathbf{X}^\top$  from the right, one obtains

$$\mathbf{X} \mathbf{A} \mathbf{X}^\top \mathbf{X} \mathbf{A} \mathbf{X}^\top = \mathbf{X} \mathbf{A}' \mathbf{X}^\top \mathbf{X} \mathbf{A} \mathbf{X}^\top. \quad (56)$$

- formula (55) implies for the left-hand side of the last equality that  $\mathbf{X} \mathbf{A} \mathbf{X}^\top \mathbf{X} \mathbf{A} \mathbf{X}^\top = \mathbf{X} \mathbf{A} \mathbf{X}^\top$ .
- Furthermore, we have shown in Lemma 3.6 that  $\mathbf{X}^\top \mathbf{X} \mathbf{A} \mathbf{X}^\top = \mathbf{X}^\top$  for each generalized inverse  $\mathbf{A}$  of  $\mathbf{X}^\top \mathbf{X}$ .
- If this is plugged into the right-hand side of (56), one obtains (54). □

By using Lemma 3.8, we can now prove the invariance property of the BLUE-estimator  $\mathbf{a}^\top \bar{\boldsymbol{\beta}}$  considered in Theorem 3.11 which was already mentioned above.

**Theorem 3.12** *Let  $\mathbf{a}^\top \boldsymbol{\beta}$  be an estimable function of the parameter vector  $\boldsymbol{\beta}$ . Then the BLUE-estimator  $\mathbf{a}^\top \bar{\boldsymbol{\beta}} = \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y}$  does not depend on the choice of the generalized inverse  $(\mathbf{X}^\top \mathbf{X})^-$ .*



**Proof**

- Recall: It follows from Theorem 3.9 that for each estimable function  $\mathbf{a}^\top \boldsymbol{\beta}$  of  $\boldsymbol{\beta}$  there is a  $\mathbf{c} \in \mathbb{R}^n$ , such that  $\mathbf{a}^\top = \mathbf{c}^\top \mathbf{X}$ .
- Together with Lemma 3.8 this implies that

$$\mathbf{a}^\top \bar{\boldsymbol{\beta}} = \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y} = \mathbf{c}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y}$$

does not depend on the choice of the generalized inverse  $(\mathbf{X}^\top \mathbf{X})^-$ .  $\square$

**Example** (one-factor analysis of variance)

- For the reparametrized model of one-factor analysis of variance  $\mu + \alpha_i$  and  $\alpha_i - \alpha_{i'}$  are estimable functions of  $\boldsymbol{\beta}$  for  $i = 1, \dots, k$  and  $i, i' = 1, \dots, k$ , respectively, cf. Theorem 3.10.
- Theorem 3.11 implies that  $\hat{\mu} + \hat{\alpha}_i$  and  $\hat{\alpha}_i - \hat{\alpha}_{i'}$  are BLUE-estimators for  $\mu + \alpha_i$  and  $\alpha_i - \alpha_{i'}$ , respectively, where

$$\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y} = (\hat{\mu}, \hat{\alpha}_1, \dots, \hat{\alpha}_k)^\top$$

with  $\hat{\mu} = \bar{Y}_{..}$  and  $\hat{\alpha}_i = \bar{Y}_{i.} - \bar{Y}_{..}$  is the solution of the normal equation (25), which was already considered in Section 3.2.1.

**Example** (two-factor analysis of variance with balanced subsamples)

- For the model of two-factor analysis of variance with balanced subsamples introduced in Section 3.1.3 the linear functions  $\mu + \alpha_{i_1}^{(1)} + \alpha_{i_2}^{(2)} + \alpha_{i_1 i_2}$  of the parameter vector

$$\boldsymbol{\beta} = (\mu, \alpha_1^{(1)}, \dots, \alpha_{k_1}^{(1)}, \alpha_1^{(2)}, \dots, \alpha_{k_2}^{(2)}, \alpha_{11}, \dots, \alpha_{k_1 k_2})$$

are estimable without bias for arbitrary  $i_1 = 1, \dots, k_1$ ,  $i_2 = 1, \dots, k_2$ , cf. Theorem 3.10.

- Theorem 3.11 implies that  $\hat{\mu} + \hat{\alpha}_{i_1}^{(1)} + \hat{\alpha}_{i_2}^{(2)} + \hat{\alpha}_{i_1 i_2}$  is a BLUE-estimator for  $\mu + \alpha_{i_1}^{(1)} + \alpha_{i_2}^{(2)} + \alpha_{i_1 i_2}$ , where

$$\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y} = (\hat{\mu}, \hat{\alpha}_1^{(1)}, \dots, \hat{\alpha}_{k_1}^{(1)}, \hat{\alpha}_1^{(2)}, \dots, \hat{\alpha}_{k_2}^{(2)}, \hat{\alpha}_{11}, \dots, \hat{\alpha}_{k_1 k_2})$$

with

$$\hat{\mu} = \bar{Y}_{...}, \quad \hat{\alpha}_{i_1}^{(1)} = \bar{Y}_{i_1 \dots} - \bar{Y}_{...}, \quad \hat{\alpha}_{i_2}^{(2)} = \bar{Y}_{\dots i_2} - \bar{Y}_{...}, \quad \hat{\alpha}_{i_1 i_2} = \bar{Y}_{\dots} + \bar{Y}_{i_1 i_2 \dots} - \bar{Y}_{i_1 \dots} - \bar{Y}_{\dots i_2}$$

is the solution of the normal equation (25) already considered in Section 3.2.3.

- Furthermore, it follows from Theorem 3.10 that in the model of two-factor analysis of variance without interactions, i.e.,  $\alpha_{i_1 i_2} = 0$  for arbitrary  $i_1 = 1, \dots, k_1$ ,  $i_2 = 1, \dots, k_2$ , also  $\alpha_{i_1}^{(1)} - \alpha_{i_1'}^{(1)}$  and  $\alpha_{i_2}^{(2)} - \alpha_{i_2'}^{(2)}$  are estimable without bias for arbitrary  $i_1, i_1' = 1, \dots, k_1$  with  $i_1 \neq i_1'$  and  $i_2, i_2' = 1, \dots, k_2$  with  $i_2 \neq i_2'$ , respectively.
- Therefore, Theorem 3.11 implies that  $\hat{\alpha}_{i_1}^{(1)} - \hat{\alpha}_{i_1'}^{(1)}$  and  $\hat{\alpha}_{i_2}^{(2)} - \hat{\alpha}_{i_2'}^{(2)}$  are BLUE-estimators for  $\alpha_{i_1}^{(1)} - \alpha_{i_1'}^{(1)}$  and  $\alpha_{i_2}^{(2)} - \alpha_{i_2'}^{(2)}$ , respectively.

### 3.3 Normally Distributed Error Terms

In addition to the model assumptions made at the beginning of Chapter 3, we now assume again that  $n > m$  and that the random error terms  $\varepsilon_1, \dots, \varepsilon_n : \Omega \rightarrow \mathbb{R}$  are independent and (identically) normally distributed, i.e., in particular, that  $\varepsilon_i \sim N(0, \sigma^2)$  for each  $i = 1, \dots, n$ .

### 3.3.1 Maximum–Likelihood Estimation

- In the same way as in the case of a design matrix with full column rank, which has been discussed in Section 2.2.1, Theorem 1.3 implies that the vector  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  of response variables is normally distributed with

$$\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}). \quad (57)$$

- In other words: The distribution of the random vector  $\mathbf{Y}$  is absolutely continuous with density

$$f_{\mathbf{Y}}(\mathbf{y}) = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right) \quad (58)$$

for each  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$ .

- Therefore, the loglikelihood function  $\log L(\mathbf{y}; \boldsymbol{\beta}, \sigma^2) = \log f_{\mathbf{Y}}(\mathbf{y})$  has the form

$$\log L(\mathbf{y}; \boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} |\mathbf{y} - \mathbf{X}\boldsymbol{\beta}|^2. \quad (59)$$

- In order to specify a maximum–likelihood estimator for the parameter vector  $(\boldsymbol{\beta}, \sigma^2)$ , we first consider the mapping

$$\mathbb{R}^m \ni \boldsymbol{\beta} \mapsto \log L(\mathbf{y}; \boldsymbol{\beta}, \sigma^2) \quad (60)$$

for arbitrary but fixed  $\mathbf{y} \in \mathbb{R}^n$  and  $\sigma^2 > 0$  (just as in the proof of Theorem 2.6).

- It follows from (59) and (60) that the following expression  $e(\boldsymbol{\beta})$  has to be minimized, where

$$e(\boldsymbol{\beta}) = \frac{1}{n} |\mathbf{y} - \mathbf{X}\boldsymbol{\beta}|^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

- Hence, Theorem 3.6 implies that the LS–estimator  $\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$  simultaneously is an ML–estimator for  $\boldsymbol{\beta}$  (which does not depend on  $\sigma^2$ ).
- Moreover, one obtains just as in the proof of Theorem 2.6 that an ML–estimator for  $(\boldsymbol{\beta}, \sigma^2)$  is given by  $(\bar{\boldsymbol{\beta}}, \bar{\sigma}^2)$ , where

$$\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \quad \text{and} \quad \bar{\sigma}^2 = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}). \quad (61)$$

#### Remark

- In Section 3.2.2 we have shown that in general  $\bar{\boldsymbol{\beta}}$  is not an unbiased estimator for  $\boldsymbol{\beta}$ .
- Similarly,  $\bar{\sigma}^2$  is not an unbiased estimator for  $\sigma^2$ ; however, an unbiased estimator for  $\sigma^2$  can be derived by a simple modification of  $\bar{\sigma}^2$ .
- In order to show this statement, the following properties of the matrix

$$\mathbf{G} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \quad (62)$$

are useful, which can be perceived as a generalization of the corresponding matrix properties derived in Lemmas 2.1 and 2.2 for the case of a design matrix  $\mathbf{X}$  with full column rank.

**Lemma 3.9** *Let  $\text{rk}(\mathbf{X}) = r \leq m$ . Then for the matrix  $\mathbf{G}$  given in (62) it holds that*

- 1)  $\mathbf{G}$  is idempotent and symmetric,
- 2)  $\mathbf{G}\mathbf{X} = \mathbf{0}$  and 3)  $\text{tr}(\mathbf{G}) = \text{rk}(\mathbf{G}) = n - r$ .

**Proof**

- Lemma 3.6 implies that

$$\begin{aligned} \mathbf{G}^2 &= (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)(\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\ &= \mathbf{I} - 2\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \underbrace{\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{=\mathbf{X}^\top} \\ &= \mathbf{I} - 2\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{G}. \end{aligned}$$

- In Lemma 3.6 we have shown that  $((\mathbf{X}^\top \mathbf{X})^{-1})^\top$  is a generalized inverse of  $\mathbf{X}^\top \mathbf{X}$ . Hence, it follows from Lemma 3.8 that

$$\mathbf{G}^\top = (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top = \mathbf{I} - \mathbf{X}((\mathbf{X}^\top \mathbf{X})^{-1})^\top \mathbf{X}^\top = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{G}.$$

- Thus, the first part of the statement is proved. In order to prove the second part of the statement, it suffices to observe that

$$\mathbf{G}\mathbf{X} = (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)\mathbf{X} = \mathbf{X} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0},$$

where the last but one equality follows from Lemma 3.7.

- The third part of the statement can be proved as follows:

- Lemmas 3.3 and 3.7 imply that

$$r = \text{rk}(\mathbf{X}) = \text{rk}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}) \leq \text{rk}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \leq \text{rk}(\mathbf{X}) = r,$$

i.e.,

$$\text{rk}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = r. \quad (63)$$

- It follows from Lemma 3.6 that the matrix  $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  is idempotent because it holds that

$$(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \underbrace{\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{=\mathbf{X}^\top} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top.$$

- Since additionally  $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  is symmetric, Lemma 1.3 together with (63) implies that

$$\begin{aligned} \text{tr}(\mathbf{G}) &= \text{tr}(\mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\ &= \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\ &= n - \text{rk}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = n - r. \end{aligned} \quad \square$$

By use of Lemma 3.9 one now obtains a formula for the expectation of the ML-estimator  $\bar{\sigma}^2$  considered in (61).

**Theorem 3.13** *It holds that*

$$\mathbb{E} \bar{\sigma}^2 = \frac{n-r}{n} \sigma^2. \quad (64)$$

**Proof** Due to the properties of the matrix  $\mathbf{G} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  derived in Lemma 3.9, it holds for the ML-estimator  $\bar{\sigma}^2 = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})/n$  considered in (61) that

$$\begin{aligned} \mathbb{E} \bar{\sigma}^2 &= \frac{1}{n} \mathbb{E} \left( (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}) \right) \\ &= \frac{1}{n} \mathbb{E} \left( \mathbf{Y}^\top (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{Y} \right) \\ &= \frac{1}{n} \mathbb{E} \left( \mathbf{Y}^\top \mathbf{G}^\top \mathbf{G} \mathbf{Y} \right) = \frac{1}{n} \mathbb{E} (|\mathbf{G}\mathbf{Y}|^2) = \frac{1}{n} \mathbb{E} (|\mathbf{G}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})|^2) \\ &= \frac{1}{n} \mathbb{E} (|\mathbf{G}\boldsymbol{\varepsilon}|^2) = \frac{1}{n} \mathbb{E} \left( \boldsymbol{\varepsilon}^\top \mathbf{G}^\top \mathbf{G} \boldsymbol{\varepsilon} \right) = \frac{1}{n} \mathbb{E} \left( \boldsymbol{\varepsilon}^\top \mathbf{G} \boldsymbol{\varepsilon} \right) \\ &= \frac{\sigma^2}{n} \text{tr}(\mathbf{G}) = \frac{n-r}{n} \sigma^2. \end{aligned}$$

□

**Remark** By using the notation

$$S^2 = \frac{1}{n-r} (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}) \quad \text{and} \quad S^2 = \frac{1}{n-r} \boldsymbol{\varepsilon}^\top \mathbf{G}\boldsymbol{\varepsilon}, \quad (65)$$

Theorem 3.13 implies that  $\mathbb{E}S^2 = \sigma^2$ , i.e.,  $S^2$  is an unbiased estimator for  $\sigma^2$ .

For being able to specify the distributions of the estimators  $\bar{\boldsymbol{\beta}}$  and  $S^2$ , we need the notion of the degenerate multivariate normal distribution, cf. Section 1.2.5.

**Theorem 3.14** *Let  $\text{rk}(\mathbf{X}) = r \leq m$ . Then it holds that*

$$\bar{\boldsymbol{\beta}} \sim \text{N}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} ((\mathbf{X}^\top \mathbf{X})^{-1})^\top) \quad (66)$$

and

$$\frac{(n-r)S^2}{\sigma^2} \sim \chi_{n-r}^2, \quad (67)$$

where the random variables  $\bar{\boldsymbol{\beta}}$  and  $S^2$  are independent.

**Proof**

- For the estimator  $\bar{\boldsymbol{\beta}}$  given in (61) it holds that

$$\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \boldsymbol{\mu} + \mathbf{B}\boldsymbol{\varepsilon},$$

where

$$\boldsymbol{\mu} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}, \quad \mathbf{B} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top, \quad \boldsymbol{\varepsilon} \sim \text{N}(\mathbf{o}, \sigma^2 \mathbf{I}_n).$$

- Now the definition of the (degenerate) multivariate normal distribution implies that  $\bar{\boldsymbol{\beta}} \sim \text{N}(\boldsymbol{\mu}, \mathbf{K})$ , where

$$\mathbf{K} = \sigma^2 \mathbf{B}\mathbf{B}^\top = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} ((\mathbf{X}^\top \mathbf{X})^{-1})^\top.$$

- Thus, (66) is proved. In order to prove (67), we use the identity derived in the proof of Theorem 3.13, i.e.,

$$S^2 = \frac{1}{n-r} \boldsymbol{\varepsilon}^\top \mathbf{G}\boldsymbol{\varepsilon}, \quad \text{where } \mathbf{G} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top. \quad (68)$$

- As  $\boldsymbol{\varepsilon} \sim \text{N}(\mathbf{o}, \sigma^2 \mathbf{I}_n)$  and as we have shown in Lemma 3.9 that

- the matrix  $\mathbf{G}$  is idempotent and symmetric
- with  $\text{rk}(\mathbf{G}) = n-r$ ,

it follows from Theorem 1.9 that the quadratic form  $(n-r)S^2/\sigma^2$  has a (central)  $\chi^2$ -distribution with  $n-r$  degrees of freedom, i.e.,  $(n-r)S^2/\sigma^2 \sim \chi_{n-r}^2$ .

- Since every idempotent and symmetric matrix simultaneously is positive semidefinite and since

$$\mathbf{B}\mathbf{G} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top - \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{=\mathbf{X}^\top} = \mathbf{0}$$

due to Lemma 3.6, it follows from Theorem 1.10 that the random variables  $\mathbf{B}\boldsymbol{\varepsilon}$  and  $\boldsymbol{\varepsilon}^\top \mathbf{G}\boldsymbol{\varepsilon}$  are independent. Thus, also the random variables  $\bar{\boldsymbol{\beta}}$  and  $S^2$  are independent. □

### 3.3.2 Testing Linear Hypotheses

In this section we discuss a generalized version of the test for linear forms of  $\boldsymbol{\beta}$  which we considered in Section 2.2.3 for the case of a design matrix  $\mathbf{X}$  with full column rank, cf. Theorem 2.11. However, we now assume that  $\text{rk}(\mathbf{X}) = r < m$ .

- Let  $s \in \{1, \dots, m\}$ , let  $\mathbf{H}$  be an  $s \times m$  matrix with full rank  $\text{rk}(\mathbf{H}) = s$  and let  $\mathbf{d} \in \mathbb{R}^s$ .
- The hypothesis to be tested is

$$H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d} \quad \text{versus} \quad H_1 : \mathbf{H}\boldsymbol{\beta} \neq \mathbf{d}, \quad (69)$$

where we assume that the entries of the matrix  $\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_s)^\top$  and the components of the vector  $\mathbf{d} = (d_1, \dots, d_s)^\top$  are known.

- In order to verify the null hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  considered in (69), we construct a test statistic whose distribution does *not* depend on the unknown parameter vector  $(\boldsymbol{\beta}, \sigma^2)$  (in a similar way as in Theorem 2.11).
- For this purpose, we introduce the following term for assuming that the components of the vector  $\mathbf{H}\boldsymbol{\beta}$  are estimable without bias.

**Definition** The hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  is called *testable* if all components  $\mathbf{h}_1^\top \boldsymbol{\beta}, \dots, \mathbf{h}_s^\top \boldsymbol{\beta}$  of the vector  $\mathbf{H}\boldsymbol{\beta}$  are estimable functions of the parameter vector  $\boldsymbol{\beta}$ .

**Remark** Theorem 3.9 implies that the hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  is testable if and only if

- there is an  $s \times n$  matrix  $\mathbf{C}$ , such that

$$\mathbf{H} = \mathbf{C}\mathbf{X}, \quad (70)$$

or

- the matrix  $\mathbf{H}$  fulfills the following equation:

$$\mathbf{H}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{X} = \mathbf{H}. \quad (71)$$

In order to construct a test statistic for the verification of the null hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  considered in (69), the following lemma is useful.

#### Lemma 3.10

- Let  $s \leq m$ , let  $\mathbf{H}$  be an  $s \times m$  matrix with full rank  $\text{rk}(\mathbf{H}) = s$  which fulfills (70) or (71) and let  $(\mathbf{X}^\top \mathbf{X})^-$  be an arbitrary generalized inverse of  $\mathbf{X}^\top \mathbf{X}$ .
- Then the  $s \times s$  matrix  $\mathbf{H}(\mathbf{X}^\top \mathbf{X})^- \mathbf{H}^\top$  is positive definite (and thus invertible).

#### Proof

- One can show that the symmetric  $m \times m$  matrix  $\mathbf{X}^\top \mathbf{X}$  with  $\text{rk}(\mathbf{X}^\top \mathbf{X}) = r < m$  can be represented in the form:

$$\mathbf{X}^\top \mathbf{X} = \mathbf{P}^{-1} \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{P}^{-1},$$

where the  $m \times m$  matrix  $\mathbf{P}$  is invertible and symmetric; cf. also Lemma 3.4.

- In the proof of Lemma 3.5 we have shown that in this case

$$\mathbf{P} \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{m-r} \end{pmatrix} \mathbf{P} = \mathbf{P}\mathbf{P}$$

defines a generalized inverse of  $\mathbf{X}^\top \mathbf{X}$ , which obviously is positive definite.

- Now it follows from Lemma 1.8 that also the matrix  $\mathbf{HPPH}^\top$  is positive definite.
- This implies that  $\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top$  is positive definite for any arbitrary generalized inverse  $(\mathbf{X}^\top \mathbf{X})^{-}$  of  $\mathbf{X}^\top \mathbf{X}$  as it follows from (70) that

$$\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top = \mathbf{CX}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top \mathbf{C}^\top = \mathbf{CXPPX}^\top \mathbf{C}^\top = \mathbf{HPPH}^\top,$$

where the second equality holds due to Lemma 3.8.  $\square$

### Remark

- Lemma 3.10 implies that the test statistic  $T_{\mathbf{H}}$  with

$$T_{\mathbf{H}} = \frac{(\overline{\mathbf{H}\boldsymbol{\beta}} - \mathbf{d})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top)^{-1} (\overline{\mathbf{H}\boldsymbol{\beta}} - \mathbf{d})}{sS^2} \quad (72)$$

is well-defined, where  $\overline{\boldsymbol{\beta}}$  and  $S^2$  are the estimators for  $\boldsymbol{\beta}$  and  $\sigma^2$  which are given in (61) and (65), respectively.

- This test statistic is a generalization of the corresponding test statistic  $T_{\mathbf{H}}$  considered in Section 2.2.3 for a design matrix  $\mathbf{X}$  with full rank. The distribution of the test statistic  $T_{\mathbf{H}}$  given in (72) can be specified as follows.

**Theorem 3.15** *Let the hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  be testable. Then, assuming that  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  is true, it holds that  $T_{\mathbf{H}} \sim F_{s,n-r}$ , i.e., the test statistic  $T_{\mathbf{H}}$  given in (72) has an F-distribution with  $(s, n-r)$  degrees of freedom.*

**Proof** Assuming that  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  is true, the following statements hold.

- The definition of  $\overline{\boldsymbol{\beta}}$  in (61) implies that

$$\overline{\mathbf{H}\boldsymbol{\beta}} - \mathbf{d} = \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top \mathbf{Y} - \mathbf{d} = \mathbf{H}((\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top)(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) - \mathbf{d} = \boldsymbol{\mu} + \mathbf{B}\boldsymbol{\varepsilon},$$

where

$$\boldsymbol{\mu} = \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} - \mathbf{d} = \mathbf{C}\underbrace{\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top \mathbf{X}}_{=\mathbf{X}}\boldsymbol{\beta} - \mathbf{d} = \mathbf{H}\boldsymbol{\beta} - \mathbf{d} = \mathbf{o},$$

$\mathbf{B} = \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top$  and  $\boldsymbol{\varepsilon} \sim N(\mathbf{o}, \sigma^2 \mathbf{I}_n)$ .

- Hence, or the numerator  $Z = (\overline{\mathbf{H}\boldsymbol{\beta}} - \mathbf{d})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top)^{-1} (\overline{\mathbf{H}\boldsymbol{\beta}} - \mathbf{d})$  of the test statistic  $T_{\mathbf{H}}$  given in (72) it holds that

$$\begin{aligned} Z &= \boldsymbol{\varepsilon}^\top \mathbf{B}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top)^{-1} \mathbf{B}\boldsymbol{\varepsilon} \\ &= \boldsymbol{\varepsilon}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top)^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top)^{-1} \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top \boldsymbol{\varepsilon} \\ &= \boldsymbol{\varepsilon}^\top \mathbf{A}\boldsymbol{\varepsilon}, \end{aligned}$$

- where the matrix  $\mathbf{A} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top)^{-1} \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top$  is idempotent as due to (71) it holds that

$$\begin{aligned} \mathbf{A}^2 &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top)^{-1} \\ &\quad \times \underbrace{\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top}_{=\mathbf{H}} (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top)^{-1} \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top \\ &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{H}^\top)^{-1} \mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-}\mathbf{X}^\top = \mathbf{A}. \end{aligned}$$

- Since  $\mathbf{A}$  is also symmetric with  $\text{rk}(\mathbf{A}) = s$ , Theorem 1.9 implies that the quadratic form  $Z/\sigma^2$  has a  $\chi^2$ -distribution with  $s$  degrees of freedom.
- Furthermore, it has been shown in Theorem 3.14 that  $(n-r)S^2/\sigma^2 \sim \chi_{n-r}^2$  and that the random variables  $\bar{\boldsymbol{\beta}}$  and  $S^2$  are independent.
- Therefore, the random variables  $Z$  and  $S^2$  are also independent and it holds that

$$T_{\mathbf{H}} = \frac{Z/s\sigma^2}{S^2/\sigma^2} \sim F_{s,n-r}. \quad \square$$

### Remark

- The choice of the test statistic  $T_{\mathbf{H}}$  in (72) can be motivated in the following way: Theorem 1.9 implies (in a similar way as in the proof of Theorem 3.15) that the quadratic form  $Z/\sigma^2$  with

$$Z = (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d})$$

in general (i.e., without assuming that  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  is true) has a noncentral  $\chi^2$ -distribution  $\chi_{s,\lambda}^2$  with

$$\lambda = \frac{(\mathbf{H}\boldsymbol{\beta} - \mathbf{d})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\boldsymbol{\beta} - \mathbf{d})}{\sigma^2}.$$

- This implies that

$$\mathbb{E} \left( \frac{Z}{\sigma^2} \right) = \frac{d}{dt} \mathbb{E} \exp \left( \frac{tZ}{\sigma^2} \right) \Big|_{t=0} = s + \lambda,$$

where the last equality follows from the formula for the moment generating function of the  $\chi_{s,\lambda}^2$ -distribution which has been derived in Theorem 1.8.

- In other words: It holds that

$$\mathbb{E} \left( \frac{Z}{s} \right) = \sigma^2 + \frac{(\mathbf{H}\boldsymbol{\beta} - \mathbf{d})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\boldsymbol{\beta} - \mathbf{d})}{s} \quad (73)$$

and Theorem 3.13 implies that  $\mathbb{E}(S^2) = \sigma^2$ .

- Hence, assuming that the null hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  is true, the expectations of numerator and denominator of the test statistic  $T_{\mathbf{H}}$  are equal.
- On the other hand, Lemmas 1.8 and 3.10 imply that the inverse matrix  $(\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1}$  is positive definite and thus,

$$(\mathbf{H}\boldsymbol{\beta} - \mathbf{d})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\boldsymbol{\beta} - \mathbf{d}) > 0$$

if the hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  is false. In this case, it follows from (73) that

$$\mathbb{E} \left( \frac{Z}{s} \right) > \sigma^2 = \mathbb{E}(S^2). \quad (74)$$

- In general, we have  $\mathbb{E} T_{\mathbf{H}} = \mathbb{E}(Z/s) \mathbb{E}(1/S^2)$  (due to the independence of  $Z$  and  $S^2$ ) and the Jensen inequality implies that  $\mathbb{E}(1/S^2) > 1/\mathbb{E}(S^2)$ .
- So (74) implies that

$$\mathbb{E} T_{\mathbf{H}} > \frac{\mathbb{E} \left( \frac{Z}{s} \right)}{\mathbb{E}(S^2)} > 1$$

in the case that  $H_0$  is false.

- Therefore, it is reasonable to reject the null hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}$  if the test statistic  $T_{\mathbf{H}}$  takes values which are significantly larger than 1.

- Thus, due to the distributional property of the test statistic  $T_{\mathbf{H}}$  which has been derived in Theorem 3.15, the null hypothesis  $H_0$  is rejected if  $T_{\mathbf{H}} > F_{s, n-r, 1-\alpha}$ .

In some cases it is more convenient to consider an *alternative representation* of the test statistic  $T_{\mathbf{H}}$  given in (72). For this purpose, we define the following *sums of squared errors*  $SSE$  and  $SSE_H$  by

$$SSE = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}), \quad \text{where} \quad \bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}, \quad (75)$$

and

$$SSE_H = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H)^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H), \quad \text{where} \quad \bar{\boldsymbol{\beta}}_H = \bar{\boldsymbol{\beta}} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d}). \quad (76)$$

**Theorem 3.16** For the test statistic  $T_{\mathbf{H}}$  given in (72) it holds that

$$T_{\mathbf{H}} = \frac{(SSE_H - SSE)/s}{SSE/(n-r)}. \quad (77)$$

**Proof**

- It holds that

$$\begin{aligned} SSE_H &= (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H)^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H) \\ &= (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}} + \mathbf{X}(\bar{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}}_H))^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}} + \mathbf{X}(\bar{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}}_H)) \\ &= (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}) + (\bar{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}}_H)^\top \mathbf{X}^\top \mathbf{X} (\bar{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}}_H) \end{aligned}$$

because parts 1 and 2 of the statement in Lemma 3.9 imply that

$$\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}) = \mathbf{X}^\top (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{Y} = \underbrace{\mathbf{X}^\top \mathbf{G}}_{=\mathbf{0}} \mathbf{Y} = \mathbf{0}.$$

- Furthermore, it follows from (70), i.e.,  $\mathbf{H} = \mathbf{C}\mathbf{X}$ , and from Lemmas 3.6 and 3.7 that

$$\begin{aligned} SSE_H &= (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}) \\ &\quad + (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} \underbrace{\mathbf{H}((\mathbf{X}^\top \mathbf{X})^{-1})^\top \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top}_{=\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top} \\ &\quad \times (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d}) \\ &= SSE + (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d}). \quad \square \end{aligned}$$

**Remark**

- From the definition of  $\bar{\boldsymbol{\beta}}_H$  in (76) it follows that

$$\mathbf{H}\bar{\boldsymbol{\beta}}_H = \mathbf{H} \left( \bar{\boldsymbol{\beta}} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d}) \right) = \mathbf{d},$$

i.e., the random vector  $\bar{\boldsymbol{\beta}}_H$  given in (76) only takes values in the *restricted parameter space*  $\Theta_H = \{\boldsymbol{\beta} \in \mathbb{R}^m : \mathbf{H}\boldsymbol{\beta} = \mathbf{d}\}$ .

- Furthermore, it can easily be shown that  $\bar{\boldsymbol{\beta}}_H$  minimizes the mean squared error  $e(\boldsymbol{\beta})$  for all  $\boldsymbol{\beta} \in \Theta_H$ , where

$$e(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n (Y_i - (\beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_m x_{im}))^2.$$



### 3.3.3 Confidence Regions

- For the construction of confidence regions we proceed in a similar way as in Section 2.2.4, where we considered the case that the design matrix  $\mathbf{X}$  has full rank, i.e.,  $\text{rk}(\mathbf{X}) = m$ . However, in doing so we now assume that  $\text{rk}(\mathbf{X}) = r < m$  as we did in Section 3.3.2.
- Let  $s \in \{1, \dots, m\}$  and let  $\mathbf{H}$  be an  $s \times m$  matrix with full rank  $\text{rk}(\mathbf{H}) = s$  whose entries are known, where  $\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_s)^\top$ .
- Then Theorem 3.15 immediately leads to the following *confidence region for the vector  $\mathbf{H}\boldsymbol{\beta}$*  with confidence level  $1 - \alpha \in (0, 1)$ .

**Theorem 3.17** *Let all components  $\mathbf{h}_1^\top \boldsymbol{\beta}, \dots, \mathbf{h}_s^\top \boldsymbol{\beta}$  of the vector  $\mathbf{H}\boldsymbol{\beta}$  be estimable functions of  $\boldsymbol{\beta}$ . Then the (random) ellipsoid*

$$E = \left\{ \mathbf{d} \in \mathbb{R}^s : \frac{(\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{d})}{sS^2} \leq F_{s, n-r, 1-\alpha} \right\} \quad (78)$$

*is a confidence region for  $\mathbf{H}\boldsymbol{\beta}$  with confidence level  $1 - \alpha \in (0, 1)$ , where  $\bar{\boldsymbol{\beta}}$  and  $S^2$  are the estimators for  $\boldsymbol{\beta}$  and  $\sigma^2$  given in (61) and (65), respectively.*

In particular, Theorem 3.17 implies the following result.

**Corollary 3.1** *For each  $i \in \{1, \dots, s\}$  the (random) interval*

$$(\underline{\theta}, \bar{\theta}) = \left( \mathbf{h}_i^\top \bar{\boldsymbol{\beta}} - t_{n-r, 1-\alpha/2} S \sqrt{\mathbf{h}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{h}_i}, \mathbf{h}_i^\top \bar{\boldsymbol{\beta}} + t_{n-r, 1-\alpha/2} S \sqrt{\mathbf{h}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{h}_i} \right) \quad (79)$$

*is a confidence interval for  $\mathbf{h}_i^\top \boldsymbol{\beta}$  with confidence level  $1 - \alpha \in (0, 1)$ .*

#### Example

- We consider the following linear model, cf. N. Ravishanker und D.K. Dey (2002) *A First Course in Linear Model Theory*, Chapman & Hall/CRC, S. 235:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \end{pmatrix},$$

where  $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \varepsilon_3)^\top \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ .

- By means of Corollary 3.1, a confidence interval for  $\beta_1 + \beta_2/3 + 2\beta_3/3$  with confidence level  $1 - \alpha = 0.95$  shall be specified.
- As  $\text{rk}(\mathbf{X}) = 2 < m = 3$ , we first need to check whether the function

$$\mathbf{h}^\top \boldsymbol{\beta} = \beta_1 + \beta_2/3 + 2\beta_3/3 \quad (80)$$

of  $\boldsymbol{\beta}^\top = (\beta_1, \beta_2, \beta_3)$  with  $\mathbf{h}^\top = (1, 1/3, 2/3)$  is estimable without bias.

- Due to criterion 1 in Theorem 3.9, this is the case if and only if there is a  $\mathbf{c}^\top = (c_1, c_2, c_3) \in \mathbb{R}^3$ , such that  $\mathbf{h}^\top = \mathbf{c}^\top \mathbf{X}$ , i.e., if

$$\begin{aligned} 1 &= c_1 + c_2 + c_3 \\ 1/3 &= c_1 + c_3 \\ 2/3 &= c_2. \end{aligned}$$

– Since this system of equations obviously is solvable,  $\mathbf{h}^\top \boldsymbol{\beta}$  is estimable without bias.

- Moreover, it holds that

$$\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} 3 & 2 & 1 \\ 2 & 2 & 0 \\ 1 & 0 & 1 \end{pmatrix}$$

and a generalized inverse of  $\mathbf{X}^\top \mathbf{X}$  is given by

$$(\mathbf{X}^\top \mathbf{X})^- = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 3/2 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

- This yields that  $\mathbf{h}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{h} = 1/2$ . Thus,

$$(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top = \begin{pmatrix} 0 & 1 & 0 \\ 1/2 & -1 & 1/2 \\ 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y} = \begin{pmatrix} Y_2 \\ (Y_1 - 2Y_2 + Y_3)/2 \\ 0 \end{pmatrix}.$$

- Hence, one obtains that  $(\mathbf{X}\bar{\boldsymbol{\beta}})^\top = ((Y_1 + Y_3)/2, Y_2, (Y_1 + Y_3)/2)$  and

$$\mathbf{h}^\top \bar{\boldsymbol{\beta}} = (Y_1 + 4Y_2 + Y_3)/6 \quad \text{and} \quad S^2 = (Y_1 - Y_3)^2/2.$$

- Therefore, a confidence interval  $(\underline{\theta}, \bar{\theta})$  for  $\mathbf{h}^\top \boldsymbol{\beta}$  with confidence level  $1 - \alpha = 0.95$  is obtained, which has the form

$$(\underline{\theta}, \bar{\theta}) = \left( (Y_1 + 4Y_2 + Y_3)/6 - Z, (Y_1 + 4Y_2 + Y_3)/6 + Z \right),$$

where  $Z = t_{1,0.975}|Y_1 - Y_3|/2$ .

By generalizing Theorem 2.12, we now derive a so-called *Scheffé confidence band*, i.e., simultaneous confidence intervals for a whole class of estimable functions of the parameter vector  $\boldsymbol{\beta}$ .

- Let  $s \in \{1, \dots, m\}$ , let  $\mathbf{H}$  once more be an  $s \times m$  matrix with full rank, i.e.,  $\text{rk}(\mathbf{H}) = s$ , where  $\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_s)^\top$ , and let all components  $\mathbf{h}_1^\top \boldsymbol{\beta}, \dots, \mathbf{h}_s^\top \boldsymbol{\beta}$  of the vector  $\mathbf{H}\boldsymbol{\beta}$  be estimable functions of  $\boldsymbol{\beta}$ .
- As  $\mathbf{H}$  has full (row) rank, the vectors  $\mathbf{h}_1, \dots, \mathbf{h}_s$  are linearly independent and form the basis of an  $s$ -dimensional linear subspace in  $\mathbb{R}^m$ , which we denote by  $\mathcal{L} = \mathcal{L}(\mathbf{h}_1, \dots, \mathbf{h}_s)$ .
- Due to Theorem 3.10, the function  $\mathbf{h}^\top \boldsymbol{\beta}$  of  $\boldsymbol{\beta}$  is estimable without bias for each  $\mathbf{h} \in \mathcal{L}$ .
- We are looking for a number  $a_\gamma > 0$ , such that

$$\mathbf{h}^\top \bar{\boldsymbol{\beta}} - a_\gamma Z_{\mathbf{h}} \leq \mathbf{h}^\top \boldsymbol{\beta} \leq \mathbf{h}^\top \bar{\boldsymbol{\beta}} + a_\gamma Z_{\mathbf{h}} \quad (81)$$

holds for each  $\mathbf{h} \in \mathcal{L}$  simultaneously with the (given) probability  $\gamma \in (0, 1)$ , where  $Z_{\mathbf{h}} = S\sqrt{\mathbf{h}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{h}}$  and  $\bar{\boldsymbol{\beta}}, S^2$  are the estimators for  $\boldsymbol{\beta}$  and  $\sigma^2$  given in (61) and (65), respectively.

**Theorem 3.18** *Let  $a_\gamma = \sqrt{s F_{s, n-r, \gamma}}$ . Then it holds that*

$$\mathbb{P}_{\boldsymbol{\beta}} \left( \max_{\mathbf{h} \in \mathcal{L}} \frac{(\mathbf{h}^\top \bar{\boldsymbol{\beta}} - \mathbf{h}^\top \boldsymbol{\beta})^2}{S^2 \mathbf{h}^\top (\mathbf{X}^\top \mathbf{X})^- \mathbf{h}} \leq a_\gamma^2 \right) = \gamma. \quad (82)$$

**Proof**

- In a similar way as in the proof of Theorem 2.12, the Cauchy–Schwarz inequality for scalar products, cf. (65), implies that

$$(\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{H}\boldsymbol{\beta})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{H}\boldsymbol{\beta}) = \max_{\mathbf{x} \neq \mathbf{o}} \frac{(\mathbf{x}^\top (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{H}\boldsymbol{\beta}))^2}{\mathbf{x}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top) \mathbf{x}},$$

where the maximum extends over all vectors  $\mathbf{x} \in \mathbb{R}^s$  with  $\mathbf{x} \neq \mathbf{o}$ .

- By means of Theorem 3.15 we now get that

$$\begin{aligned} \gamma &= \mathbb{P}_{\boldsymbol{\beta}} \left( (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{H}\boldsymbol{\beta})^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top)^{-1} (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{H}\boldsymbol{\beta}) \leq sS^2 F_{s,n-r,\gamma} \right) \\ &= \mathbb{P}_{\boldsymbol{\beta}} \left( \max_{\mathbf{x} \neq \mathbf{o}} \frac{(\mathbf{x}^\top (\mathbf{H}\bar{\boldsymbol{\beta}} - \mathbf{H}\boldsymbol{\beta}))^2}{\mathbf{x}^\top (\mathbf{H}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{H}^\top) \mathbf{x}} \leq sS^2 F_{s,n-r,\gamma} \right) \\ &= \mathbb{P}_{\boldsymbol{\beta}} \left( \max_{\mathbf{x} \neq \mathbf{o}} \frac{((\mathbf{H}^\top \mathbf{x})^\top \bar{\boldsymbol{\beta}} - (\mathbf{H}^\top \mathbf{x})^\top \boldsymbol{\beta})^2}{(\mathbf{H}^\top \mathbf{x})^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{H}^\top \mathbf{x})} \leq sS^2 F_{s,n-r,\gamma} \right) \\ &= \mathbb{P}_{\boldsymbol{\beta}} \left( \max_{\mathbf{h} \in \mathcal{L}} \frac{(\mathbf{h}^\top \bar{\boldsymbol{\beta}} - \mathbf{h}^\top \boldsymbol{\beta})^2}{\mathbf{h}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{h}} \leq sS^2 F_{s,n-r,\gamma} \right). \end{aligned} \quad \square$$

**3.4 Examples****3.4.1 F-Test for the ANOVA Null Hypothesis**

- We consider the reparametrized model of *one-factor analysis of variance*, i.e., the design matrix  $\mathbf{X}$  is the  $n \times (k+1)$  matrix given in (13) with  $\text{rk}(\mathbf{X}) = k < m = k+1$ , where

$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 1 & 0 & 0 & \dots & 0 & 0 \\ 1 & 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 0 & 1 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & 0 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & 0 & 0 & 0 & \dots & 0 & 1 \end{pmatrix} \quad (83)$$

and the parameter vector  $\boldsymbol{\beta}$  has the form  $\boldsymbol{\beta} = (\mu, \alpha_1, \dots, \alpha_k)^\top$ .

- It shall be tested whether the levels of the predictor variable are significant, i.e., the hypothesis to be tested is the ANOVA null hypothesis  $H_0 : \alpha_1 = \dots = \alpha_k$  (against the alternative  $H_1 : \alpha_i \neq \alpha_j$  for some pair  $i, j \in \{1, \dots, k\}$  with  $i \neq j$ ). For this purpose, we use the general testing approach introduced in Theorems 3.15 and 3.16.

- An equivalent formulation of the null hypothesis  $H_0 : \alpha_1 = \dots = \alpha_k$  is given by

$$H_0 : \alpha_1 - \alpha_2 = 0, \dots, \alpha_1 - \alpha_k = 0 \quad \text{or} \quad H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}, \quad (84)$$

where  $\mathbf{H}$  is a  $(k-1) \times (k+1)$  matrix with

$$\mathbf{H} = \begin{pmatrix} 0 & 1 & -1 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & -1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 1 & 0 & 0 & \dots & -1 & 0 \\ 0 & 1 & 0 & 0 & \dots & 0 & -1 \end{pmatrix}. \quad (85)$$

- Obviously,  $\mathbf{H}$  is a matrix with full row rank, i.e.,  $\text{rk}(\mathbf{H}) = k-1$ . Furthermore, Theorem 3.10 implies that all components  $\alpha_1 - \alpha_2, \dots, \alpha_1 - \alpha_k$  of the vector  $\mathbf{H}\boldsymbol{\beta}$  are estimable functions of  $\boldsymbol{\beta}$ .
- In other words, the matrix  $\mathbf{H}$  fulfills the requirements of Theorems 3.15 and 3.16. Thus, the test statistic

$$T_{\mathbf{H}} = \frac{(SSE_H - SSE)/(k-1)}{SSE/(n-k)}$$

considered in Theorem 3.16 may be used for the verification of the hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}$ , where the sums of squares  $SSE$  and  $SSE_H$  defined in (75) and (76), respectively, can be determined as follows.

- *Recall:* In Section 3.2.1 we have shown that a generalized inverse of  $\mathbf{X}^\top \mathbf{X}$  is given by (36), i.e.,

$$(\mathbf{X}^\top \mathbf{X})^- = \begin{pmatrix} \frac{1}{n} & 0 & 0 & 0 & \dots & 0 & 0 \\ -\frac{1}{n} & \frac{1}{n_1} & 0 & 0 & \dots & 0 & 0 \\ -\frac{1}{n} & 0 & \frac{1}{n_2} & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ -\frac{1}{n} & 0 & 0 & 0 & \dots & 0 & \frac{1}{n_k} \end{pmatrix}. \quad (86)$$

- Together with (83) this implies that

$$\bar{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y} = (\bar{Y}_{..}, \bar{Y}_{1.} - \bar{Y}_{..}, \dots, \bar{Y}_{k.} - \bar{Y}_{..})^\top$$

and

$$\mathbf{X}\bar{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{Y} = \left( \underbrace{\bar{Y}_{1.}, \dots, \bar{Y}_{1.}}_{n_1}, \dots, \underbrace{\bar{Y}_{k.}, \dots, \bar{Y}_{k.}}_{n_k} \right)^\top.$$

- Hence, one gets for the sum of squares  $SSE = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})$  that

$$SSE = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2. \quad (87)$$

**Remark**

- Due to the special form (83) of the design matrix  $\mathbf{X}$ , formula (87) can also be derived directly from the fact that  $\bar{\boldsymbol{\beta}}$  is an LS-estimator. Indeed, it holds that

$$\begin{aligned} SSE &= (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}) = \min_{\boldsymbol{\beta} \in \mathbb{R}^{k+1}} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \sum_{i=1}^k \left\{ \min_{x \in \mathbb{R}} \sum_{j=1}^{n_i} (Y_{ij} - x)^2 \right\} = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\cdot})^2. \end{aligned}$$

- Moreover, it follows from the remark at the end of Section 3.3.2 that

$$SSE_H = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H)^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H) = \min_{\boldsymbol{\beta} \in \Theta_H} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2,$$

where  $\Theta_H = \{\boldsymbol{\beta} \in \mathbb{R}^{k+1} : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}\}$  and  $\{\mathbf{X}\boldsymbol{\beta} : \boldsymbol{\beta} \in \Theta_H\} \subset \mathbb{R}^n$  is the set of those  $n$ -dimensional vectors whose components are all equal.

- Then one gets that

$$SSE_H = \min_{x \in \mathbb{R}} \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - x)^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 \quad (88)$$

because the mean  $\bar{Y}_{..}$  minimizes the sum of squares  $\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - x)^2$ .

- Together with (87) this implies that

$$SSE_H - SSE = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 - \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\cdot})^2 = \sum_{i=1}^k n_i (\bar{Y}_{i\cdot} - \bar{Y}_{..})^2,$$

where the last equality follows from the decomposition

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\cdot})^2 + \sum_{i=1}^k n_i (\bar{Y}_{i\cdot} - \bar{Y}_{..})^2,$$

cf. formula (9) in Theorem 3.1.

- Therefore, for the test statistic  $T_{\mathbf{H}}$  considered in Theorem 3.16 it holds that

$$T_{\mathbf{H}} = \frac{(SSE_H - SSE)/(k-1)}{SSE/(n-k)} = \frac{(n-k) \sum_{i=1}^k n_i (\bar{Y}_{i\cdot} - \bar{Y}_{..})^2}{(k-1) \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\cdot})^2} \sim F_{k-1, n-k}. \quad (89)$$

**3.4.2 F-Tests for the Two-Factor Analysis of Variance**

Now we construct F-tests for the model of *two-factor analysis of variance* with balanced subsets, which has been introduced in Section 3.1.3, i.e.,

- the parameter vector  $\boldsymbol{\beta}$  has the form

$$\boldsymbol{\beta} = (\mu, \alpha_1^{(1)}, \dots, \alpha_{k_1}^{(1)}, \alpha_1^{(2)}, \dots, \alpha_{k_2}^{(2)}, \alpha_{11}, \dots, \alpha_{k_1 k_2})^\top,$$

- the design matrix  $\mathbf{X}$  has the dimension  $n \times m$ , where  $n = rk_1 k_2$  and  $m = 1 + k_1 + k_2 + k_1 k_2$ ,

- the entries of  $\mathbf{X}$  only consist of zeros and ones and it holds that  $\text{rk}(\mathbf{X}) = k_1 k_2 < m$ .

### Significance of the Predictor Variables

We first construct a test to investigate the question whether the levels of the first predictor variable are significant. For this purpose, we verify the hypothesis that the effects

$$\alpha_{i_1}^{(1)*} = \alpha_{i_1}^{(1)} + \frac{1}{k_2} \sum_{i_2=1}^{k_2} \alpha_{i_1 i_2}$$

of the first predictor variable plus their interactions, averaged over all levels of the second predictor variable, are equal. In other words: The hypothesis to be tested is

$$H_0 : \alpha_1^{(1)*} - \alpha_{i_1}^{(1)*} = 0 \quad \forall i_1 \in \{1, \dots, k_1\} \quad \text{vs.} \quad H_1 : \alpha_1^{(1)*} - \alpha_{i_1}^{(1)*} \neq 0 \quad \text{for some } i_1 \in \{1, \dots, k_1\}, \quad (90)$$

where it is actually sufficient to consider the pair of hypotheses

$$H_0 : \alpha_1^{(1)*} - \alpha_{i_1}^{(1)*} = 0 \quad \forall i_1 \in \{2, \dots, k_1\} \quad \text{vs.} \quad H_1 : \alpha_1^{(1)*} - \alpha_{i_1}^{(1)*} \neq 0 \quad \text{for some } i_1 \in \{2, \dots, k_1\}.$$

- It can easily be shown that the null hypothesis in (90) has the form  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}$ ,  
– where

$$\mathbf{H} = \begin{pmatrix} 0 & 1 & -1 & 0 & \dots & 0 & 0 & \dots & 0 & \frac{1}{k_2} & \dots & \frac{1}{k_2} & \frac{-1}{k_2} & \dots & \frac{-1}{k_2} & 0 & \dots & 0 \\ 0 & 1 & 0 & -1 & & 0 & 0 & \dots & 0 & \frac{1}{k_2} & \dots & \frac{1}{k_2} & 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & 0 & & & -1 & 0 & \dots & 0 & \frac{1}{k_2} & \dots & \frac{1}{k_2} & 0 & \dots & 0 & \frac{-1}{k_2} & \dots & \frac{-1}{k_2} \end{pmatrix}$$

is a  $(k_1 - 1) \times m$  matrix with full row rank  $\text{rk}(\mathbf{H}) = k_1 - 1$  and with blocks of rows of the lengths 1, 1,  $k_1 - 1$ ,  $k_2$ ,  $k_2$  and  $(k_1 - 1)k_2$ , respectively,

- and where all components of the vector  $\mathbf{H}\boldsymbol{\beta}$  are estimable functions of  $\boldsymbol{\beta}$  because of

$$\alpha_1^{(1)*} - \alpha_{i_1}^{(1)*} = \frac{1}{k_2} \sum_{i_2=1}^{k_2} (\theta_{1i_2} - \theta_{i_1 i_2}).$$

- In order to verify the hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}$  we can now again use the test statistic  $T_{\mathbf{H}}$  considered in Theorem 3.16, where

$$T_{\mathbf{H}} = \frac{(SSE_H - SSE)/(k_1 - 1)}{SSE/(k_1 k_2 (r - 1))} \sim F_{k_1 - 1, k_1 k_2 (r - 1)}$$

with

$$SSE = \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2})^2 \quad (91)$$

and

$$SSE_H - SSE = r k_2 \sum_{i_1=1}^{k_1} (\bar{Y}_{i_1 \dots} - \bar{Y}_{\dots})^2. \quad (92)$$

- The formulas (91) and (92) for the sums of squares  $SSE$  and  $SSE_H$  defined in (75) and (76), respectively, can be derived in a similar way as in Section 3.4.1.
- Indeed, the same minimization technique that has been used for the direct derivation of (87) yields (91). Moreover, the sum of squares  $SSE_H$  can be determined as follows.

- Just as before,  $\Theta_H = \{\boldsymbol{\beta} \in \mathbb{R}^{1+k_1+k_2+k_1k_2} : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}\}$  is the restricted parameter space.
- Due to the special form of the matrices  $\mathbf{X}$  and  $\mathbf{H}$ , the minimization in

$$SSE_H = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H)^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H) = \min_{\boldsymbol{\beta} \in \Theta_H} |\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}|^2$$

with respect to the set  $\{\mathbf{X}\boldsymbol{\beta} : \boldsymbol{\beta} \in \Theta_H\} \subset \mathbb{R}^{rk_1k_2}$  can (in a similar way as in formula (88)) be replaced by a minimization with respect to the set  $\mathbb{R} \times \mathbb{R}_H^{k_2} \times \mathbb{R}_H^{k_1k_2} \subset \mathbb{R}^{1+k_2+k_1k_2}$  of those vectors  $\mathbf{x} = (x, x_1, \dots, x_{k_2}, x_{11}, \dots, x_{k_1k_2}) \in \mathbb{R}^{1+k_2+k_1k_2}$  which fulfill the following conditions:

$$\sum_{i_2=1}^{k_2} x_{i_2} = \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} x_{i_1i_2} = 0.$$

- In more detail, it holds that

$$\begin{aligned} SSE_H &= \min_{\mathbf{x} \in \mathbb{R} \times \mathbb{R}_H^{k_2} \times \mathbb{R}_H^{k_1k_2}} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1i_2j} - (x + x_{i_2} + x_{i_1i_2}))^2 \\ &= \min_{\mathbf{x} \in \mathbb{R} \times \mathbb{R}_H^{k_2} \times \mathbb{R}_H^{k_1k_2}} \left\{ \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1i_2j} - \bar{Y}_{i_1i_2\cdot})^2 + k_1k_2r(\bar{Y}\dots - x)^2 + k_2r \sum_{i_1=1}^{k_1} (\bar{Y}_{i_1\cdot\cdot} - \bar{Y}\dots)^2 \right. \\ &\quad \left. + k_1r \sum_{i_2=1}^{k_2} (\bar{Y}\cdot i_2\cdot - \bar{Y}\dots - x_{i_2})^2 + r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1i_2\cdot} - \bar{Y}_{i_1\cdot\cdot} - \bar{Y}\cdot i_2\cdot + \bar{Y}\dots - x_{i_1i_2})^2 \right\}, \end{aligned}$$

i.e.,

$$SSE_H = \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1i_2j} - \bar{Y}_{i_1i_2\cdot})^2 + k_2r \sum_{i_1=1}^{k_1} (\bar{Y}_{i_1\cdot\cdot} - \bar{Y}\dots)^2.$$

- Together with (91) this implies (92).

### Remark

- In the same way, one can construct a test to investigate the question whether the levels of the second predictor variable are significant. For this purpose, we verify the hypothesis that the effects

$$\alpha_{i_2}^{(2)*} = \alpha_{i_2}^{(2)} + \frac{1}{k_1} \sum_{i_1=1}^{k_1} \alpha_{i_1i_2}$$

of the second predictor variable plus their interactions, averaged over all levels of the first predictor variable, are equal.

- Thus, the hypothesis to be tested is

$$H_0 : \alpha_1^{(2)*} - \alpha_2^{(2)*} = 0, \dots, \alpha_1^{(2)*} - \alpha_{k_2}^{(2)*} = 0 \quad \text{vs.} \quad H_1 : \alpha_1^{(2)*} - \alpha_{i_2}^{(2)*} \neq 0 \quad \text{for some } i_2 \in \{1, \dots, k_2\}.$$

- In this case, one obtains the test statistic

$$T_{\mathbf{H}} = \frac{k_1k_2(r-1)rk_1 \sum_{i_2=1}^{k_2} (\bar{Y}\cdot i_2\cdot - \bar{Y}\dots)^2}{(k_2-1) \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1i_2j} - \bar{Y}_{i_1i_2\cdot})^2} \sim F_{k_2-1, k_1k_2(r-1)}.$$

### Interactions between both Predictor Variables

Now we construct a test to check whether there are significant interactions between the two predictor variables. For this purpose, the hypothesis

$$H_0 : \alpha_{i_1 i_2}^* - \alpha_{i_1 i_2}^* = 0 \quad \forall (i_1, i_2) \in \{1, \dots, k_1\} \times \{1, \dots, k_2\} \quad (93)$$

is tested, where

$$\alpha_{i_1 i_2}^* = \alpha_{i_1 i_2} - \bar{\alpha}_{i_1 \cdot} - \bar{\alpha}_{\cdot i_2} + \bar{\alpha}_{\cdot \cdot}$$

and

$$\bar{\alpha}_{i_1 \cdot} = \frac{1}{k_2} \sum_{i_2=1}^{k_2} \alpha_{i_1 i_2}, \quad \bar{\alpha}_{\cdot i_2} = \frac{1}{k_1} \sum_{i_1=1}^{k_1} \alpha_{i_1 i_2}, \quad \bar{\alpha}_{\cdot \cdot} = \frac{1}{k_1 k_2} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \alpha_{i_1 i_2}.$$

- In a similar way as before, one can show that the hypothesis considered in (93) can be written in the form  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}$ , where

- $\mathbf{H}$  is a  $(k_1 k_2 - 1) \times m$  matrix with full row rank, i.e.,  $\text{rk}(\mathbf{H}) = k_1 k_2 - 1$  and
- all components of the vector  $\mathbf{H}\boldsymbol{\beta}$  are estimable functions of  $\boldsymbol{\beta}$  as it holds that

$$\alpha_{i_1 i_2}^* = \theta_{i_1 i_2} - \frac{1}{k_2} \sum_{i_2=1}^{k_2} \theta_{i_1 i_2} - \frac{1}{k_1} \sum_{i_1=1}^{k_1} \theta_{i_1 i_2} + \frac{1}{k_1 k_2} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \theta_{i_1 i_2}.$$

- For the verification of the hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}$  we can thus consider the test statistic

$$T_{\mathbf{H}} = \frac{(SSE_H - SSE)/(k_1 k_2 - 1)}{SSE/(k_1 k_2 (r - 1))}$$

considered in Theorem 3.16, where  $SSE$  is given by (91) as before while the sum of squares  $SSE_H$  results from the following considerations.

- Due to the special form of the matrices  $\mathbf{X}$  and  $\mathbf{H}$ , the set  $\{\mathbf{X}\boldsymbol{\beta} : \boldsymbol{\beta} \in \Theta_H\} \subset \mathbb{R}^{r k_1 k_2}$  in the minimization in

$$SSE_H = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H)^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H) = \min_{\boldsymbol{\beta} \in \Theta_H} |\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}|^2$$

can (in a similar way as before) be replaced by the set  $\mathbb{R} \times \mathbb{R}_H^{k_1} \times \mathbb{R}_H^{k_2} \subset \mathbb{R}^{1+k_1+k_2}$  of those vectors

$$\mathbf{x} = (x, x_1^{(1)}, \dots, x_{k_1}^{(1)}, x_1^{(2)}, \dots, x_{k_2}^{(2)}) \in \mathbb{R}^{1+k_1+k_2}$$

which fulfill the following conditions:

$$\sum_{i_1=1}^{k_1} x_{i_1}^{(1)} = \sum_{i_2=1}^{k_2} x_{i_2}^{(2)} = 0.$$

- Indeed, it holds that

$$SSE_H = \min_{\mathbf{x} \in \mathbb{R} \times \mathbb{R}_H^{k_1} \times \mathbb{R}_H^{k_2}} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - (x + x_{i_1}^{(1)} + x_{i_2}^{(2)}))^2$$

and

$$\begin{aligned} SSE_H &= \min_{\mathbf{x} \in \mathbb{R} \times \mathbb{R}_H^{k_1} \times \mathbb{R}_H^{k_2}} \left\{ \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot})^2 \right. \\ &\quad + k_1 k_2 r (\bar{Y}_{\cdot \cdot} - x)^2 + k_2 r \sum_{i_1=1}^{k_1} (\bar{Y}_{i_1 \cdot} - \bar{Y}_{\cdot \cdot} - x_{i_1}^{(1)})^2 + k_1 r \sum_{i_2=1}^{k_2} (\bar{Y}_{\cdot i_2} - \bar{Y}_{\cdot \cdot} - x_{i_2}^{(2)})^2 \\ &\quad \left. + r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot} - \bar{Y}_{\cdot i_2} + \bar{Y}_{\cdot \cdot})^2 \right\} \\ &= \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot})^2 + r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot} - \bar{Y}_{\cdot i_2} + \bar{Y}_{\cdot \cdot})^2. \end{aligned}$$



- Together with (91) this implies that

$$SSE_H - SSE = r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot \cdot} - \bar{Y}_{\cdot i_2 \cdot} + \bar{Y}_{\cdot \cdot \cdot})^2. \quad (94)$$

- Therefore, for the test statistic  $T_H$  considered in Theorem 3.16 it holds that

$$\begin{aligned} T_H &= \frac{(SSE_H - SSE)/(k_1 k_2 - 1)}{SSE/(k_1 k_2 (r - 1))} \\ &= \frac{k_1 k_2 (r - 1) r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot \cdot} - \bar{Y}_{\cdot i_2 \cdot} + \bar{Y}_{\cdot \cdot \cdot})^2}{(k_1 k_2 - 1) \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot})^2} \sim F_{k_1 k_2 - 1, k_1 k_2 (r - 1)}. \end{aligned}$$

### 3.4.3 Two-Factor Analysis of Variance with Hierarchical Classification

- Instead of the model of two-factor analysis of variance with interactions, introduced in Section 3.1.3, one sometimes considers the following model of two-factor analysis of variance with *hierarchical classification* of the pairs of levels  $i_1, i_2$  of the two predictor variables.
- Here we consider the representation

$$\theta_{i_1 i_2} = \mu + \alpha_{i_1}^{(1)} + \alpha_{i_2|i_1}^{(2|1)}, \quad \forall i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2 \quad (95)$$

of the expectations  $\theta_{i_1 i_2} = \mathbb{E} Y_{i_1 i_2 j}$  of the sampling variables  $Y_{i_1 i_2 j}$ .

- In other words: Into each of the  $k_1$  levels of the first, i.e., superior, predictor variable  $k_2$  levels of the second (inferior) predictor variable are embedded.
- This situation can occur, e.g., in clinical trials which are carried out in  $k_1$  countries (superior predictor variable) and  $k_2$  hospitals in each country (inferior predictor variable).
- Then the parameter vector  $\beta$  has the dimension  $m = 1 + k_1 + k_1 k_2$  with

$$\beta = (\mu, \alpha_1^{(1)}, \dots, \alpha_{k_1}^{(1)}, \alpha_{1|1}^{(2|1)}, \dots, \alpha_{k_2|k_1}^{(2|1)})^\top$$

- and
  - $\mu$  is again perceived as *general mean* of the expectations  $\mathbb{E} Y_{i_1 i_2 j}$  of the sampling variables  $Y_{i_1 i_2 j}$ ,
  - $\alpha_{i_1}^{(1)}$  is called the *effect* of the  $i_1$ -th level of the superior predictor variable and
  - $\alpha_{i_2|i_1}^{(2|1)}$  is called the *effect* of the  $i_2$ -th level of the inferior predictor variable in the case that the  $i_1$ -th level of the superior predictor variable is on hand.
- Again, we only consider the balanced case, i.e., we assume that all  $k_1 \cdot k_2$  subsamples ( $Y_{i_1 i_2 j}$ ,  $j = 1, \dots, n_{i_1 i_2}$ ) have the same sample size.
- Hence, it holds that  $n_{i_1 i_2} = r$  for all  $i_1 = 1, \dots, k_1$  and  $i_2 = 1, \dots, k_2$  with  $r = n/(k_1 k_2)$ , the design matrix  $\mathbf{X}$  has the dimension  $n \times m$  with  $n = r k_1 k_2$  and  $m = 1 + k_1 + k_1 k_2$ , and the entries of  $\mathbf{X}$  only consist of zeros and ones;  $\text{rk}(\mathbf{X}) = k_1 k_2 < m$ .

### Significance of the Superior Predictor Variable

- Just in the same way as in Section 3.4.2 one can construct a test to investigate the question whether the levels of the superior predictor variable are significant. For this purpose, we verify the hypothesis that the mean effects  $\alpha_{i_1}^{(1)*}$  are equal, where

$$\alpha_{i_1}^{(1)*} = \alpha_{i_1}^{(1)} + \frac{1}{k_2} \sum_{i_2=1}^{k_2} \alpha_{i_2|i_1}^{(2|1)}.$$

- In other words: The hypothesis to be tested is

$$H_0 : \alpha_{i_1}^{(1)*} - \alpha_{i_1}^{(1)*} = 0 \quad \forall i_1 \in \{1, \dots, k_1\} \quad \text{versus} \quad H_1 : \alpha_{i_1}^{(1)*} - \alpha_{i_1}^{(1)*} \neq 0 \quad \text{for some } i_1 \in \{1, \dots, k_1\}.$$

- One can show that the null hypothesis has the form  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}$ , where  $\mathbf{H}$  is a  $(k_1 - 1) \times m$  matrix with full row rank, i.e.,  $\text{rk}(\mathbf{H}) = k_1 - 1$ , and all components of the vector  $\mathbf{H}\boldsymbol{\beta}$  are estimable functions of  $\boldsymbol{\beta}$ .
- For the verification of the hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}$  we can thus use the test statistic

$$T_{\mathbf{H}} = \frac{(SSE_H - SSE)/(k_1 - 1)}{SSE/(k_1 k_2 (r - 1))} \sim F_{k_1 - 1, k_1 k_2 (r - 1)}$$

considered in Theorem 3.16 with

$$SSE = \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2})^2, \quad SSE_H - SSE = r k_2 \sum_{i_1=1}^{k_1} (\bar{Y}_{i_1 \dots} - \bar{Y}_{\dots})^2, \quad (96)$$

where the formulas in (96) are proved in the same way as in (91) and (92), respectively.

### Significance of the Inferior Predictor Variable

- In order to check whether the levels of the inferior predictor variable are significant, one can proceed in a similar way as in the last test in Section 3.4.2 (test for significance of interactions). For this purpose, the hypothesis

$$H_0 : \alpha_{1|1}^{(2|1)*} - \alpha_{i_2|i_1}^{(2|1)*} = 0 \quad \forall (i_1, i_2) \in \{1, \dots, k_1\} \times \{1, \dots, k_2\} \quad (97)$$

is tested, where

$$\alpha_{i_2|i_1}^{(2|1)*} = \alpha_{i_2|i_1}^{(2|1)} - \bar{\alpha}_{i_1} + \bar{\alpha}, \quad \bar{\alpha}_{i_1} = \frac{1}{k_2} \sum_{i_2=1}^{k_2} \alpha_{i_2|i_1}^{(2|1)}, \quad \bar{\alpha} = \frac{1}{k_1 k_2} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \alpha_{i_2|i_1}^{(2|1)}.$$

- It can be shown that the hypothesis considered in (93) can be written in the form  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}$ , where  $\mathbf{H}$  is a  $k_1(k_2 - 1) \times m$  matrix with full row rank, i.e.,  $\text{rk}(\mathbf{H}) = k_1(k_2 - 1)$ , and all components of the vector  $\mathbf{H}\boldsymbol{\beta}$  are estimable functions of  $\boldsymbol{\beta}$ .
- For the verification of the hypothesis  $H_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{o}$  one can thus use the test statistic

$$T_{\mathbf{H}} = \frac{(SSE_H - SSE)/(k_1(k_2 - 1))}{SSE/(k_1 k_2 (r - 1))}$$

considered in Theorem 3.16, where the sums of squares  $SSE$  and  $SSE_H$  defined in (75) and (76), respectively, can be determined as follows.

- Just as before, it holds that

$$SSE = \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2})^2, \quad (98)$$

and the set  $\{\mathbf{X}\boldsymbol{\beta} : \boldsymbol{\beta} \in \Theta_H\}$  in the minimization in

$$SSE_H = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H)^\top (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\beta}}_H) = \min_{\boldsymbol{\beta} \in \Theta_H} |\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}|^2$$

can be replaced by the set  $\mathbb{R} \times \mathbb{R}_H^{k_1} \subset \mathbb{R}^{1+k_1}$  of those vectors  $\mathbf{x} = (x, x_1^{(1)}, \dots, x_{k_1}^{(1)}) \in \mathbb{R}^{1+k_1}$  for which  $\sum_{i_1=1}^{k_1} x_{i_1}^{(1)} = 0$  holds, leading to

$$\begin{aligned} SSE_H &= \min_{\mathbf{x} \in \mathbb{R} \times \mathbb{R}_H^{k_1}} \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - (x + x_{i_1}^{(1)}))^2 \\ &= \min_{\mathbf{x} \in \mathbb{R} \times \mathbb{R}_H^{k_1}} \left\{ \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot})^2 + k_1 k_2 r (\bar{Y} \dots - x)^2 \right. \\ &\quad \left. + k_2 r \sum_{i_1=1}^{k_1} (\bar{Y}_{i_1 \cdot \cdot} - \bar{Y} \dots - x_{i_1}^{(1)})^2 + r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot \cdot})^2 \right\} \\ &= \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot})^2 + r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot \cdot})^2. \end{aligned}$$

- Together with (98) this implies that

$$SSE_H - SSE = r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot \cdot})^2.$$

- Therefore, it holds that

$$T_{\mathbf{H}} = \frac{k_1 k_2 (r-1) r \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} (\bar{Y}_{i_1 i_2 \cdot} - \bar{Y}_{i_1 \cdot \cdot})^2}{k_1 (k_2 - 1) \sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} \sum_{j=1}^r (Y_{i_1 i_2 j} - \bar{Y}_{i_1 i_2 \cdot})^2} \sim F_{k_1(k_2-1), k_1 k_2 (r-1)}.$$

## 4 Generalized Linear Models

- In Chapters 2 and 3 we always assumed for the linear model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 
  - that  $\mathbb{E}\boldsymbol{\varepsilon} = \mathbf{o}$ , i.e.,  $(\mathbb{E}Y_1, \dots, \mathbb{E}Y_n)^\top = \mathbf{X}\boldsymbol{\beta}$ ,
  - where furthermore  $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$  holds if  $\boldsymbol{\varepsilon} \sim N(\mathbf{o}, \sigma^2\mathbf{I})$ .
- Now, we generalize this model and allow that the expectations  $\mathbb{E}Y_1, \dots, \mathbb{E}Y_n$  of the sampling variables  $Y_1, \dots, Y_n$ 
  - can be expressed by the components of the vector  $\mathbf{X}\boldsymbol{\beta}$  by using an arbitrary monotone function  $g : G \rightarrow \mathbb{R}$ , the so-called *link function*, such that
 
$$(g(\mathbb{E}Y_1), \dots, g(\mathbb{E}Y_n))^\top = \mathbf{X}\boldsymbol{\beta}, \quad (1)$$
  - where the domain  $G \subset \mathbb{R}$  of  $g$  will be specified more precisely.
- Moreover, the (independent) sampling variables  $Y_1, \dots, Y_n$  do not have to be normally distributed since we just assume that the distributions of  $Y_1, \dots, Y_n$  belong to an *exponential family*.
- In this chapter we will always assume (as in Chapter 2) that the design matrix  $\mathbf{X}$  has full column rank, i.e.,  $\text{rk}(\mathbf{X}) = m$ .

In the same way as in the linear models, which have been investigated in Chapters 2 and 3, the goal is to estimate the parameter vector  $\boldsymbol{\beta}$  from the observation of the random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ , where we assume that the link function  $g : G \rightarrow \mathbb{R}$  is known.

### 4.1 Definition and Basic Properties

#### 4.1.1 Exponential Family

We assume that the sample variables  $Y_1, \dots, Y_n$  are independent (but in general *not* identically distributed),

- where their distributions belong to a one-parametric *exponential family*, i.e., their densities or probability mass functions, respectively, have the following form: For each  $i \in \{1, \dots, n\}$  it holds that
  - in the absolutely continuous case

$$f(y; \theta_i) = \exp\left(\frac{1}{\tau^2} (y\theta_i + a(y, \tau) - b(\theta_i))\right), \quad \forall y \in \mathbb{R}, \quad (2)$$

- in the discrete case

$$\mathbb{P}_{\theta_i}(Y_i = y) = \exp\left(\frac{1}{\tau^2} (y\theta_i + a(y, \tau) - b(\theta_i))\right), \quad \forall y \in C, \quad (3)$$

where  $a : \mathbb{R} \times (0, \infty) \rightarrow \mathbb{R}$  and  $b : \Theta \rightarrow \mathbb{R}$  are certain functions and  $C \subset \mathbb{R}$  is the smallest countable subset of  $\mathbb{R}$ , for which it holds that  $\mathbb{P}_{\theta_i}(Y_i \in C) = 1$ .

- $\tau^2 > 0$  is a so-called *nuisance parameter*, which does not depend on the index  $i$ , where it is often assumed that  $\tau^2$  is known.
- Then

$$\Theta = \left\{ \theta \in \mathbb{R} : \int_{-\infty}^{\infty} \exp\left(\frac{y\theta + a(y, \tau)}{\tau^2}\right) dy < \infty \right\} \quad (4)$$

or

$$\Theta = \left\{ \theta \in \mathbb{R} : \sum_{y \in C} \exp\left(\frac{y\theta + a(y, \tau)}{\tau^2}\right) < \infty \right\} \quad (5)$$

is the *natural parameter space*, where we always assume that the integrability condition in (4) or (5), respectively, is fulfilled for at least two different  $\theta_1, \theta_2 \in \mathbb{R}$ .

**Remark** In the absolutely continuous case the nuisance parameter  $\tau^2$  can be seen as an additional variance parameter, while  $\tau^2$  usually is set to 1 in the discrete case.

**Lemma 4.1** *The parameter space  $\Theta \subset \mathbb{R}$ , given in (4) and (5), respectively, is an interval in  $\mathbb{R}$ .*

**Proof**

- We only consider the absolutely continuous case since the proof of the discrete case proceeds analogously.
- One can easily see that for arbitrary  $x_1, x_2 \in \mathbb{R}$  and  $\alpha \in (0, 1)$  we have

$$(e^{x_1})^\alpha (e^{x_2})^{1-\alpha} \leq \max_{i=1,2} (e^{x_i})^\alpha (e^{x_i})^{1-\alpha} = \max_{i=1,2} e^{x_i} \leq e^{x_1} + e^{x_2}.$$

- This and the notation  $\theta = \alpha\theta_1 + (1-\alpha)\theta_2$  imply that for arbitrary  $\theta_1, \theta_2 \in \Theta$  and  $\alpha \in (0, 1)$

$$\begin{aligned} \int_{-\infty}^{\infty} \exp\left(\frac{y\theta + a(y, \tau)}{\tau^2}\right) dy &= \int_{-\infty}^{\infty} \left(\exp\left(\frac{y\theta_1 + a(y, \tau)}{\tau^2}\right)\right)^\alpha \left(\exp\left(\frac{y\theta_2 + a(y, \tau)}{\tau^2}\right)\right)^{1-\alpha} dy \\ &\leq \int_{-\infty}^{\infty} \left(\exp\left(\frac{y\theta_1 + a(y, \tau)}{\tau^2}\right) + \exp\left(\frac{y\theta_2 + a(y, \tau)}{\tau^2}\right)\right) dy < \infty. \end{aligned}$$

- Therefore, it also holds that  $\theta \in \Theta$ . □

Because of Lemma 4.1 we will always assume in this chapter that  $\Theta \subset \mathbb{R}$  is an *open* interval such that the integrability condition in (4) and (5), respectively, is fulfilled for each  $\theta \in \Theta$ .

**Lemma 4.2**

- *Let the distribution of the random variable  $Y : \Omega \rightarrow \mathbb{R}$  be given by (2) or (3) for an arbitrary  $\theta \in \Theta$  such that*

$$\mathbb{E}(Y^2) < \infty \quad \forall \theta \in \Theta \tag{6}$$

*holds and the function  $b : \Theta \rightarrow \mathbb{R}$  is twice continuously differentiable.*

- *Then it holds that*

$$\mathbb{E}Y = b^{(1)}(\theta) \quad \text{and} \quad \text{Var} Y = \tau^2 b^{(2)}(\theta). \tag{7}$$

**Proof**

- Once more, we only treat the absolutely continuous case since the proof of the discrete case proceeds analogously. It holds that

$$\begin{aligned} \mathbb{E}Y &= \int_{-\infty}^{\infty} y \exp\left(\frac{1}{\tau^2} (y\theta + a(y, \tau) - b(\theta))\right) dy = e^{-b(\theta)/\tau^2} \int_{-\infty}^{\infty} y \exp\left(\frac{1}{\tau^2} (y\theta + a(y, \tau))\right) dy \\ &= e^{-b(\theta)/\tau^2} \tau^2 \int_{-\infty}^{\infty} \frac{d}{d\theta} \exp\left(\frac{1}{\tau^2} (y\theta + a(y, \tau))\right) dy \\ &= e^{-b(\theta)/\tau^2} \tau^2 \frac{d}{d\theta} \int_{-\infty}^{\infty} \exp\left(\frac{1}{\tau^2} (y\theta + a(y, \tau))\right) dy \\ &= e^{-b(\theta)/\tau^2} \tau^2 \frac{d}{d\theta} e^{b(\theta)/\tau^2} \underbrace{\int_{-\infty}^{\infty} \exp\left(\frac{1}{\tau^2} (y\theta + a(y, \tau) - b(\theta))\right) dy}_{=1} \\ &= b^{(1)}(\theta). \end{aligned}$$

- In a similar way we get that  $\mathbb{E}(Y^2) = \tau^2 b^{(2)}(\theta) + (\mathbb{E}Y)^2$ . □

### 4.1.2 Link of the Parameters; Natural Link Function

- Now, we assume that the function  $b : \Theta \rightarrow \mathbb{R}$  is twice continuously differentiable with  $b^{(2)}(\theta) > 0$  for each  $\theta \in \Theta$ .
- Furthermore, let  $G = \{b^{(1)}(\theta) : \theta \in \Theta\}$  and the link function  $g : G \rightarrow \mathbb{R}$  be twice continuously differentiable, such that  $g^{(1)}(x) \neq 0$  for each  $x \in G$ . The inverse function of  $g$  is denoted by  $h = g^{-1}$ .
- We consider the generalized linear model (GLM) given in (1), i.e., it holds that

$$(g(\mathbb{E}Y_1), \dots, g(\mathbb{E}Y_n))^\top = \mathbf{X}\boldsymbol{\beta}. \quad (8)$$

- By using the notation  $\mathbf{X} = (x_{ij})$ ,  $\mathbf{x}_i = (x_{i1}, \dots, x_{im})^\top$  and  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n)^\top$ , where  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ , formula (8) implies for the expectations  $\mu_i = \mathbb{E}Y_i$  ( $i = 1, \dots, n$ ) that

$$\mu_i = h(\eta_i) = h(\mathbf{x}_i^\top \boldsymbol{\beta}) \quad \text{and thus} \quad \boldsymbol{\mu} = (h(\eta_1), \dots, h(\eta_n))^\top, \quad (9)$$

where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top$ .

- Because of (7) and (8) the parameters  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)^\top$  and  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)^\top$  are related as follows: It holds that

$$(g(b^{(1)}(\theta_1)), \dots, g(b^{(1)}(\theta_n)))^\top = \mathbf{X}\boldsymbol{\beta}. \quad (10)$$

- Together with (9) this implies that

$$(b^{(1)}(\theta_1), \dots, b^{(1)}(\theta_n)) = (h(\mathbf{x}_1^\top \boldsymbol{\beta}), \dots, h(\mathbf{x}_n^\top \boldsymbol{\beta}))$$

and, equivalently,

$$(\theta_1, \dots, \theta_n) = (\psi(h(\mathbf{x}_1^\top \boldsymbol{\beta})), \dots, \psi(h(\mathbf{x}_n^\top \boldsymbol{\beta}))), \quad (11)$$

where  $\psi = (b^{(1)})^{-1}$  is the inverse function of  $b^{(1)}$ .

- Furthermore, it is possible to express the variance  $\sigma_i^2 = \text{Var} Y_i$  of the sample variables  $Y_i$  as a function  $\sigma_i^2(\boldsymbol{\beta})$  of  $\boldsymbol{\beta}$  for each  $i = 1, \dots, n$  as it follows from Lemma 4.2 and (11) that

$$\sigma_i^2(\boldsymbol{\beta}) = \tau^2 b^{(2)}(\psi(h(\mathbf{x}_i^\top \boldsymbol{\beta}))) \quad \forall i = 1, \dots, n. \quad (12)$$

**Remark** The link function  $g : G \rightarrow \mathbb{R}$  is called *natural* if  $g = \psi$ . In this case it holds that  $\theta_i = \psi(\mu_i)$  and, therefore,  $\theta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$  for each  $i = 1, \dots, n$ , i.e.,

$$(\theta_1, \dots, \theta_n)^\top = \mathbf{X}\boldsymbol{\beta}. \quad (13)$$

## 4.2 Examples

### 4.2.1 Linear Model with Normally Error Terms

- For the linear model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (14)$$

with normally distributed error terms  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top \sim \text{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ , considered in Section 2.2, it holds that

$$Y_i \sim \text{N}(\mu_i, \sigma^2) \quad \text{with} \quad \mu_i = \mathbf{x}_i^\top \boldsymbol{\beta}, \quad \forall i = 1, \dots, n, \quad (15)$$

where we assume that  $\sigma^2$  is known.

- Then the distribution of  $Y_i$  belongs to the one-parametric exponential family considered in Section 4.1.1 since the density  $f(y; \theta_i)$  of  $Y_i$  can be written in the following form, where  $\theta_i = \mu_i$  for each  $i = 1, \dots, n$ :

– It holds that

$$f(y; \theta_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y - \mu_i)^2\right) = \exp\left(\frac{1}{\tau^2} (y\theta_i + a(y, \tau) - b(\theta_i))\right) \quad \forall y \in \mathbb{R},$$

– where

$$\tau^2 = \sigma^2, \quad a(y) = -\frac{y^2}{2} \quad \text{and} \quad b(\theta_i) = \frac{\theta_i^2}{2} + \sigma^2 \log \sqrt{2\pi\sigma^2}. \quad (16)$$

- Because of (15) the link function  $g : \mathbb{R} \rightarrow \mathbb{R}$  fulfills  $g(x) = x$  for each  $x \in \mathbb{R}$ .
  - Moreover, (16) yields  $x = b^{(1)}(x)$  for each  $x \in \mathbb{R}$ .
  - Therefore, it holds that  $g(x) = x = \psi(x)$  for each  $x \in \mathbb{R}$ , i.e., the natural link function is given by  $g(x) = x$ .

#### 4.2.2 Binary Categorical Regression

- In this section we consider the case that the sample variables  $Y_1, \dots, Y_n$  are Bernoulli-distributed, i.e., they can only take the values 0 and 1 with positive probability.
  - Here, we use the notation

$$\pi_i = \mathbb{P}(Y_i = 1) \quad \left( = \mu_i = \mathbb{E}Y_i \right) \quad \forall i = 1, \dots, n,$$

where it is assumed that  $0 < \pi_i < 1$  for each  $i = 1, \dots, n$ .

- In this case the probabilities  $\pi_1, \dots, \pi_n$  are linked to the parameter vector  $\boldsymbol{\beta}$  by using a link function  $g : (0, 1) \rightarrow \mathbb{R}$ , i.e.,

$$(g(\pi_1), \dots, g(\pi_n))^\top = \mathbf{X}\boldsymbol{\beta}. \quad (17)$$

- For each  $i = 1, \dots, n$  the  $\text{Bin}(1, \pi_i)$ -distribution belongs to the exponential family introduced in Section 4.1.1, where  $\theta_i = \log(\pi_i/(1 - \pi_i))$ .
  - Because for  $y = 0, 1$  it holds that

$$\mathbb{P}_{\theta_i}(Y_i = y) = \pi_i^y (1 - \pi_i)^{1-y} = \exp\left(y \log\left(\frac{\pi_i}{1 - \pi_i}\right) + \log(1 - \pi_i)\right) = \exp\left(\frac{1}{\tau^2} (y\theta_i + a(y, \tau) - b(\theta_i))\right),$$

– where

$$\tau^2 = 1, \quad a(y) = 0 \quad \text{and} \quad b(\theta_i) = \log(1 + e^{\theta_i}). \quad (18)$$

#### Remark

- From (18) it follows that  $(b^{(1)})^{-1}(x) = \log(x/(1 - x))$  for each  $x \in (0, 1)$ , i.e., the natural link function  $g : (0, 1) \rightarrow \mathbb{R}$  is given by

$$g(x) = \log\left(\frac{x}{1 - x}\right) \quad \forall x \in (0, 1). \quad (19)$$

- The GLM, considered in (17), with the natural link function, given in (19), is then called (binary) *logistic regression model*.
- In this case the dependency of the probabilities  $\pi_i = \pi_i(\boldsymbol{\beta})$  of the linear combinations  $\mathbf{x}_i^\top \boldsymbol{\beta}$  is given by

$$\pi_i = \frac{1}{1 + \exp(-\mathbf{x}_i^\top \boldsymbol{\beta})} \quad \forall i = 1, \dots, n. \quad (20)$$

- Another (nonnatural) link function  $g : (0, 1) \rightarrow \mathbb{R}$ , which is considered in this context, is given by

$$g = \Phi^{-1}, \quad (21)$$

- where  $\Phi : \mathbb{R} \rightarrow (0, 1)$  denotes the distribution function of the  $N(0, 1)$ -distribution.
- Then it holds that  $\pi_i = \Phi(\mathbf{x}_i^\top \boldsymbol{\beta})$  for each  $i = 1, \dots, n$  and the GLM is called the model of the *probit analysis*.

### 4.2.3 Poisson–Distributed Sample Variables with Natural Link Function

- Now, let the sample variables  $Y_1, \dots, Y_n$  be Poisson–distributed, i.e., let  $Y_i \sim \text{Poi}(\lambda_i)$  with  $0 < \lambda_i < \infty$  for each  $i = 1, \dots, n$ .
- The  $\text{Poi}(\lambda_i)$ –distribution also belongs to the exponential family introduced in Section 4.1.1, where  $\theta_i = \log \lambda_i$ 
  - because it holds for each  $y = 0, 1, \dots$  that

$$\mathbb{P}_{\theta_i}(Y_i = y) = \frac{\lambda_i^y e^{-\lambda_i}}{y!} = \exp(y \log \lambda_i - \log(y!) - \lambda_i) = \exp\left(\frac{1}{\tau^2} (y\theta_i + a(y, \tau) - b(\theta_i))\right),$$

– where

$$\tau^2 = 1, \quad a(y) = -\log(y!) \quad \text{and} \quad b(\theta_i) = e^{\theta_i}.$$

- The natural link function  $g : (0, \infty) \rightarrow \mathbb{R}$  is given by

$$g(x) = \log x \quad \forall x > 0. \quad (22)$$

## 4.3 Maximum–Likelihood Estimator for $\beta$

- Since we assumed, that the distributions of the sample variables  $Y_1, \dots, Y_n$  belong to an exponential family, it is possible to estimate the parameter vector  $\beta$  by using the maximum–likelihood method.
- In order to show this, we first discuss some properties of the loglikelihood function  $\log L(\mathbf{Y}, \theta)$  of the random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  and its partial derivatives with respect to the components  $\beta_1, \dots, \beta_m$  of  $\beta$ .

### 4.3.1 Loglikelihood Function and its Partial Derivatives

- From (2) – (3) and from (11) it follows that the loglikelihood function  $\log L(\mathbf{Y}, \theta)$  of the random sample  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  can be written as a function  $\log L(\mathbf{Y}, \beta)$  of  $\beta$ .
  - From (2) – (3) it follows that

$$\log L(\mathbf{Y}, \theta) = \sum_{i=1}^n \frac{1}{\tau^2} (Y_i \theta_i + a(Y_i, \tau) - b(\theta_i)). \quad (23)$$

– This and (11) imply that

$$\log L(\mathbf{Y}, \beta) = \sum_{i=1}^n \frac{1}{\tau^2} \left( Y_i \psi(h(\mathbf{x}_i^\top \beta)) + a(Y_i, \tau) - b(\psi(h(\mathbf{x}_i^\top \beta))) \right). \quad (24)$$

- For generalized linear models with natural link function, (13) and (23) imply that

$$\log L(\mathbf{Y}, \beta) = \sum_{i=1}^n \frac{1}{\tau^2} (Y_i \mathbf{x}_i^\top \beta + a(Y_i, \tau) - b(\mathbf{x}_i^\top \beta)). \quad (25)$$

For the computation of the maximum–likelihood estimators, the knowledge of the so–called score function, i.e., the partial derivative of the loglikelihood function, is useful, as well as the Fisher–information matrix, which is defined as follows.



**Definition** For arbitrary  $i, j = 1, \dots, m$  let

$$U_i(\boldsymbol{\beta}) = \frac{\partial}{\partial \beta_i} \log L(\mathbf{Y}, \boldsymbol{\beta}) \quad \text{and} \quad I_{ij}(\boldsymbol{\beta}) = \mathbb{E}(U_i(\boldsymbol{\beta})U_j(\boldsymbol{\beta})).$$

Then the  $m$ -dimensional random vector  $\mathbf{U}(\boldsymbol{\beta}) = (U_1(\boldsymbol{\beta}), \dots, U_m(\boldsymbol{\beta}))^\top$  is called the *score vector* and the (deterministic)  $m \times m$ -matrix  $\mathbf{I}(\boldsymbol{\beta}) = (I_{ij}(\boldsymbol{\beta}))$  is called the *Fisher-information matrix*.

By using the notation

$$\frac{d\mu_i}{d\eta_i}(\boldsymbol{\beta}) = \left. \frac{dh(s)}{ds} \right|_{s=\eta_i} = \left. \left( \frac{dg(t)}{dt} \right)^{-1} \right|_{t=h(\eta_i)} \quad (26)$$

we get the following result.

**Theorem 4.1** For arbitrary  $j, k = 1, \dots, m$  it holds that

$$U_j(\boldsymbol{\beta}) = \sum_{i=1}^n x_{ij}(Y_i - \mu_i(\boldsymbol{\beta})) \frac{d\mu_i}{d\eta_i}(\boldsymbol{\beta}) \frac{1}{\sigma_i^2(\boldsymbol{\beta})} \quad (27)$$

and

$$I_{jk}(\boldsymbol{\beta}) = \sum_{i=1}^n x_{ij}x_{ik} \left( \frac{d\mu_i}{d\eta_i}(\boldsymbol{\beta}) \right)^2 \frac{1}{\sigma_i^2(\boldsymbol{\beta})} \quad (28)$$

or in matrix notation

$$\mathbf{U}(\boldsymbol{\beta}) = \mathbf{X}^\top \mathbf{V}^{-1}(\boldsymbol{\beta}) \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}}(\boldsymbol{\beta}) (\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})) \quad \text{and} \quad \mathbf{I}(\boldsymbol{\beta}) = \mathbf{X}^\top \mathbf{V}^{-1}(\boldsymbol{\beta}) \left( \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}}(\boldsymbol{\beta}) \right)^2 \mathbf{X}, \quad (29)$$

where

$$\mathbf{V}(\boldsymbol{\beta}) = \text{diag}(\sigma_i^2(\boldsymbol{\beta})) \quad \text{and} \quad \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}}(\boldsymbol{\beta}) = \text{diag}\left(\frac{d\mu_i}{d\eta_i}(\boldsymbol{\beta})\right).$$

**Proof**

- The loglikelihood function, given in (23) or (24), respectively, can be written in the form

$$\log L(\mathbf{Y}, \boldsymbol{\beta}) = \sum_{i=1}^n \frac{1}{\tau^2} \ell^{(i)}(\theta_i),$$

where  $\ell^{(i)}(\theta_i) = Y_i\theta_i + a(Y_i, \tau) - b(\theta_i)$  and  $\theta_i = \psi(h(\mathbf{x}_i^\top \boldsymbol{\beta}))$ .

- Therefore, it holds for each  $j = 1, \dots, m$  that

$$U_j(\boldsymbol{\beta}) = \sum_{i=1}^n \frac{1}{\tau^2} \frac{\partial \ell^{(i)}}{\partial \beta_j}(\theta_i), \quad (30)$$

where multiple use of the chain rule yields

$$\frac{\partial \ell^{(i)}}{\partial \beta_j} = \frac{\partial \ell^{(i)}}{\partial \theta_i} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j}. \quad (31)$$

- On the other hand, it obviously holds that  $\partial \eta_i / \partial \beta_j = x_{ij}$  and from Lemma 4.2 it follows that

$$\frac{\partial \ell^{(i)}}{\partial \theta_i} = Y_i - b^{(1)}(\theta_i) \stackrel{\text{Lemma 4.2}}{=} Y_i - \mu_i$$

or

$$\left( \frac{\partial \theta_i}{\partial \mu_i} \right)^{-1} = \frac{\partial \mu_i}{\partial \theta_i} \stackrel{\text{Lemma 4.2}}{=} b^{(2)}(\theta_i) \stackrel{\text{Lemma 4.2}}{=} \frac{1}{\tau^2} \sigma_i^2.$$

- This and (30) – (31) imply (27).
- In order to show (28) it is enough to notice that for arbitrary  $i, j = 1, \dots, n$

$$\mathbb{E}((Y_i - \mu_i)(Y_j - \mu_j)) = \begin{cases} \sigma_i^2 & \text{für } i = j, \\ 0 & \text{für } i \neq j. \end{cases}$$

because of the independence of the sample variables  $Y_1, \dots, Y_n$ .

- From this and (27) it follows that

$$\begin{aligned} I_{jk}(\boldsymbol{\beta}) &= \mathbb{E}(U_j(\boldsymbol{\beta})U_k(\boldsymbol{\beta})) = \sum_{i=1}^n x_{ij}x_{ik} \left( \frac{d\mu_i}{d\eta_i}(\boldsymbol{\beta}) \right)^2 \frac{1}{\sigma_i^4(\boldsymbol{\beta})} \mathbb{E}((Y_i - \mu_i)^2) \\ &= \sum_{i=1}^n x_{ij}x_{ik} \left( \frac{d\mu_i}{d\eta_i}(\boldsymbol{\beta}) \right)^2 \frac{1}{\sigma_i^2(\boldsymbol{\beta})}. \end{aligned}$$

- Therefore, (28) is proved. □

**Corollary 4.1** *Let  $(g(\mathbb{E}Y_1), \dots, g(\mathbb{E}Y_n))^\top = \mathbf{X}\boldsymbol{\beta}$  be a GLM with natural link function  $g : G \rightarrow \mathbb{R}$ . Then it holds for arbitrary  $j, k = 1, \dots, m$  that*

$$U_j(\boldsymbol{\beta}) = \frac{1}{\tau^2} \sum_{i=1}^n x_{ij}(Y_i - \mu_i(\boldsymbol{\beta})) \quad \text{or} \quad \mathbf{U}(\boldsymbol{\beta}) = \frac{1}{\tau^2} \mathbf{X}^\top (\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})) \quad (32)$$

and

$$I_{jk}(\boldsymbol{\beta}) = \frac{1}{\tau^4} \sum_{i=1}^n x_{ij}x_{ik}\sigma_i^2(\boldsymbol{\beta}) \quad \text{or} \quad \mathbf{I}(\boldsymbol{\beta}) = \frac{1}{\tau^4} \mathbf{X}^\top \mathbf{V}(\boldsymbol{\beta})\mathbf{X}. \quad (33)$$

**Proof** Since  $g : G \rightarrow \mathbb{R}$  is a natural link function, we have  $\theta_i = \eta_i$  for each  $i = 1, \dots, n$ . This and Lemma 4.2 imply that

$$\frac{d\mu_i}{d\eta_i} = b^{(2)}(\theta_i) = \frac{1}{\tau^2} \sigma_i^2.$$

Now, the statement follows from Theorem 4.1. □

### 4.3.2 Hessian Matrix

Besides the (score) vector  $\mathbf{U}(\boldsymbol{\beta})$ , which consists of the first partial derivatives of the loglikelihood function  $\log L(\mathbf{Y}, \boldsymbol{\beta})$ , also the *Hessian matrix*, i.e., the  $m \times m$ -matrix

$$\mathbf{W}(\boldsymbol{\beta}) = (W_{ij}(\boldsymbol{\beta})) = \left( \frac{\partial^2}{\partial \beta_i \partial \beta_j} \log L(\mathbf{Y}, \boldsymbol{\beta}) \right),$$

consisting of the second partial derivatives of the loglikelihood function, is needed.

**Theorem 4.2**

- For each GLM it holds that

$$\mathbf{W}(\boldsymbol{\beta}) = \mathbf{X}^\top \mathbf{R}(\boldsymbol{\beta}) \text{diag}(Y_i - \mu_i(\boldsymbol{\beta})) \mathbf{X} - \mathbf{I}(\boldsymbol{\beta}), \quad (34)$$

where  $\mathbf{I}(\boldsymbol{\beta})$  is the Fisher-information matrix, given in (29), and  $\mathbf{R}(\boldsymbol{\beta}) = \text{diag}(v_i(\boldsymbol{\beta}))$  is an  $(n \times n)$ -diagonal matrix with

$$v_i(\boldsymbol{\beta}) = \frac{1}{\tau^2} \left. \frac{d^2 u(s)}{ds^2} \right|_{s=\mathbf{x}_i^\top \boldsymbol{\beta}} \quad \text{and} \quad u = \psi \circ h.$$

- For a GLM with natural link function it particularly holds that

$$\mathbf{W}(\boldsymbol{\beta}) = -\mathbf{I}(\boldsymbol{\beta}). \quad (35)$$

**Proof**

- From formula (27) in Theorem 4.1 it follows that for arbitrary  $j, k = 1, \dots, m$

$$\begin{aligned} W_{jk}(\boldsymbol{\beta}) &= \frac{\partial}{\partial \beta_k} U_j(\boldsymbol{\beta}) \stackrel{(27)}{=} \frac{\partial}{\partial \beta_k} \sum_{i=1}^n x_{ij} (Y_i - \mu_i(\boldsymbol{\beta})) \frac{d\mu_i(\boldsymbol{\beta})}{d\eta_i} \frac{1}{\sigma_i^2(\boldsymbol{\beta})} \\ &= \sum_{i=1}^n x_{ij} \left( (Y_i - \mu_i(\boldsymbol{\beta})) \frac{\partial}{\partial \beta_k} \left( \frac{d\mu_i(\boldsymbol{\beta})}{d\eta_i} \frac{1}{\sigma_i^2(\boldsymbol{\beta})} \right) - \frac{d\mu_i(\boldsymbol{\beta})}{d\eta_i} \frac{1}{\sigma_i^2(\boldsymbol{\beta})} \frac{\partial \mu_i}{\partial \beta_k}(\boldsymbol{\beta}) \right). \end{aligned}$$

- Here with the notation  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$  we get from Lemma 4.2 that

$$\frac{d\mu_i(\boldsymbol{\beta})}{d\eta_i} \frac{1}{\sigma_i^2(\boldsymbol{\beta})} = b^{(2)}(\psi \circ h(\eta_i)) (\psi \circ h)^{(1)}(\eta_i) \frac{1}{\tau^2 b^{(2)}(\psi \circ h(\eta_i))} = \frac{1}{\tau^2} (\psi \circ h)^{(1)}(\eta_i)$$

and therefore

$$\frac{\partial}{\partial \beta_k} \left( \frac{d\mu_i(\boldsymbol{\beta})}{d\eta_i} \frac{1}{\sigma_i^2(\boldsymbol{\beta})} \right) = \frac{1}{\tau^2} (\psi \circ h)^{(2)}(\eta_i) x_{ik}.$$

- Furthermore, it holds that

$$\frac{\partial \mu_i}{\partial \beta_k} = \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta_k}.$$

- Altogether, it follows that

$$W_{jk}(\boldsymbol{\beta}) = \sum_{i=1}^n x_{ij} x_{ik} (Y_i - \mu_i) v_i - \sum_{i=1}^n x_{ij} x_{ik} \left( \frac{d\mu_i}{d\eta_i} \right)^2 \frac{1}{\sigma_i^2}.$$

- This and the representation formula (29) for the Fisher-information matrix  $\mathbf{I}(\boldsymbol{\beta})$  yield (34).
- Since the superposition  $u = \psi \circ h$  of a GLM with natural link function is the identity function, it holds in this case that  $\mathbf{R}(\boldsymbol{\beta}) = \mathbf{0}$ . Therefore, (35) follows from (34).  $\square$

**Remark** For the examples of GLM considered in Section 4.2, we get the following formulas for  $\mathbf{U}(\boldsymbol{\beta})$  and  $\mathbf{W}(\boldsymbol{\beta})$  from Theorems 4.1 and 4.2 or from Corollary 4.1, respectively.

1. For the linear model  $\mathbb{E} \mathbf{Y} = \mathbf{X} \boldsymbol{\beta}$  with normally distributed sample variables (and with the link function  $g(x) = x$ ), one has that  $(d\boldsymbol{\mu}/d\boldsymbol{\eta})(\boldsymbol{\beta})$  is the identity matrix. Therefore, it holds that

$$\mathbf{U}(\boldsymbol{\beta}) = \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}), \quad \mathbf{W}(\boldsymbol{\beta}) = -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}, \quad (36)$$

cf. Section 2.2.

2. For the logistic regression model (with the natural link function) it holds that

$$\mathbf{U}(\boldsymbol{\beta}) = \mathbf{X}^\top (\mathbf{Y} - \boldsymbol{\pi}), \quad \mathbf{W}(\boldsymbol{\beta}) = -\mathbf{X}^\top \text{diag}(\pi_i(1 - \pi_i))\mathbf{X}, \quad (37)$$

where  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_n)^\top$  and the probabilities  $\pi_i$  can be expressed by  $\boldsymbol{\beta}$  (cf.(20)).

3. For Poisson-distributed sample variables with natural link function it holds that

$$\mathbf{U}(\boldsymbol{\beta}) = \mathbf{X}^\top (\mathbf{Y} - \boldsymbol{\lambda}), \quad \mathbf{W}(\boldsymbol{\beta}) = -\mathbf{X}^\top \text{diag}(\lambda_i)\mathbf{X}, \quad (38)$$

where  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)^\top$  and  $\lambda_i = e^{\mathbf{x}_i^\top \boldsymbol{\beta}}$ .

### 4.3.3 Maximum-Likelihood Equation and Numerical Approach

- In order to determine the maximum-likelihood estimator for  $\boldsymbol{\beta}$ , the *maximum-likelihood equation*

$$\mathbf{U}(\boldsymbol{\beta}) = \mathbf{o} \quad (39)$$

is considered, which in general is nonlinear and therefore often can be solved only by using iterative methods.

- Because of Theorem 4.1 the equation (39) is equivalent to

$$\mathbf{X}^\top \mathbf{V}^{-1}(\boldsymbol{\beta}) \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}}(\boldsymbol{\beta}) (\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})) = \mathbf{o}. \quad (40)$$

#### Remark

- From Corollary 4.1 it follows that, in the case of a natural link function, (40) simplifies to:

$$\mathbf{X}^\top (\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})) = \mathbf{o}. \quad (41)$$

- Since we furthermore assume that  $0 < \sigma_i^2(\boldsymbol{\beta}) < \infty$  for each  $i = 1, \dots, n$  and that the design matrix  $\mathbf{X}$  has full column rank, the matrix  $\mathbf{W}(\boldsymbol{\beta}) = -\tau^{-4}\mathbf{X}^\top \mathbf{V}(\boldsymbol{\beta})\mathbf{X}$  of the second partial derivatives is negative definite.
- Hence, it holds that if (41) has a solution, then the solution is a uniquely determined maximum-likelihood estimator  $\hat{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$ .

Now, we discuss the basic ideas of two numerical iteration methods for solving the maximum-likelihood equation (39). We consider a sequence of random vectors  $\hat{\boldsymbol{\beta}}_0, \hat{\boldsymbol{\beta}}_1, \dots : \Omega \rightarrow \mathbb{R}^m$ , which converge under certain conditions to a random vector  $\hat{\boldsymbol{\beta}}$  such that  $\hat{\boldsymbol{\beta}}$  is a solution of (39).

#### 1. Newton's Method

- Let  $\hat{\boldsymbol{\beta}}_0 : \Omega \rightarrow \mathbb{R}^m$  be a suitably chosen start vector and let the iterations  $\hat{\boldsymbol{\beta}}_1, \dots, \hat{\boldsymbol{\beta}}_k$  be already computed.
- For the computation of the  $(k+1)$ -th iteration  $\hat{\boldsymbol{\beta}}_{k+1}$  from  $\hat{\boldsymbol{\beta}}_k$ , the left hand side  $\mathbf{U}(\boldsymbol{\beta})$  of the maximum-likelihood equation (39) is replaced by
  - the first two terms  $\mathbf{U}(\hat{\boldsymbol{\beta}}_k) + \mathbf{W}(\hat{\boldsymbol{\beta}}_k)(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}_k)$  of the Taylor series expansion of  $\mathbf{U}(\boldsymbol{\beta})$  at  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}_k$ .
  - The  $(k+1)$ -th iteration  $\hat{\boldsymbol{\beta}}_{k+1}$  is also a solution of the equation

$$\mathbf{U}(\hat{\boldsymbol{\beta}}_k) + \mathbf{W}(\hat{\boldsymbol{\beta}}_k)(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}_k) = \mathbf{o}. \quad (42)$$

- If the matrix  $\mathbf{W}(\widehat{\boldsymbol{\beta}}_k)$  is invertible, then it follows from (42) that

$$\widehat{\boldsymbol{\beta}}_{k+1} = \widehat{\boldsymbol{\beta}}_k - \mathbf{W}^{-1}(\widehat{\boldsymbol{\beta}}_k) \mathbf{U}(\widehat{\boldsymbol{\beta}}_k), \quad (43)$$

- For the constructed sequence  $\widehat{\boldsymbol{\beta}}_0, \widehat{\boldsymbol{\beta}}_1, \dots$  to converge to  $\widehat{\boldsymbol{\beta}}$ , this vector  $\widehat{\boldsymbol{\beta}}$  has to be a solution of (39) and the start vector  $\widehat{\boldsymbol{\beta}}_0$  has to be close enough to  $\widehat{\boldsymbol{\beta}}$ .

## 2. Fisher–Scoring

- Now we consider a variation of Newton’s method, the so-called *scoring–method of Fisher*, where the Hessian matrix  $\mathbf{W}(\boldsymbol{\beta})$  in (42) is replaced by the expectation matrix  $\mathbb{E} \mathbf{W}(\boldsymbol{\beta})$ .
  - This has the advantage that the  $(m \times m)$ –matrix  $\mathbb{E} \mathbf{W}(\boldsymbol{\beta})$  is invertible.
  - From Theorems 4.1 and 4.2 it follows that

$$\begin{aligned} \mathbb{E} \mathbf{W}(\boldsymbol{\beta}) &\stackrel{(34)}{=} \mathbb{E} \left( \mathbf{X}^\top \mathbf{R}(\boldsymbol{\beta}) \text{diag}(Y_i - \mu_i(\boldsymbol{\beta})) \mathbf{X} - \mathbf{I}(\boldsymbol{\beta}) \right) \\ &= -\mathbf{I}(\boldsymbol{\beta}) \\ &\stackrel{(29)}{=} -\mathbf{X}^\top \mathbf{V}^{-1}(\boldsymbol{\beta}) \left( \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}}(\boldsymbol{\beta}) \right)^2 \mathbf{X}, \end{aligned}$$

where the second equality follows from the identity  $\mathbb{E} Y_i = \mu_i(\boldsymbol{\beta})$ .

- The last term is an invertible  $(m \times m)$ –matrix because we assumed that the design matrix  $\mathbf{X}$  has full column rank and that  $(d\mu_i/d\eta_i)(\boldsymbol{\beta}) \neq 0$  for each  $i = 1, \dots, n$ .
- Therefore, instead of (43) the following iteration equation is considered:

$$\widehat{\boldsymbol{\beta}}_{k+1} = \widehat{\boldsymbol{\beta}}_k + (\mathbf{X}^\top \mathbf{Z}(\widehat{\boldsymbol{\beta}}_k) \mathbf{X})^{-1} \left( \mathbf{X}^\top \mathbf{Z}(\widehat{\boldsymbol{\beta}}_k) \left( \frac{d\boldsymbol{\eta}}{d\boldsymbol{\mu}}(\widehat{\boldsymbol{\beta}}_k) (\mathbf{Y} - \boldsymbol{\mu}(\widehat{\boldsymbol{\beta}}_k)) \right) \right), \quad (44)$$

where

$$\mathbf{Z}(\boldsymbol{\beta}) = \mathbf{V}^{-1}(\boldsymbol{\beta}) \left( \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}}(\boldsymbol{\beta}) \right)^2 \quad \text{and} \quad \left( \frac{d\boldsymbol{\eta}}{d\boldsymbol{\mu}}(\boldsymbol{\beta}) \right) = \left( \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}}(\boldsymbol{\beta}) \right)^{-1}(\boldsymbol{\beta}).$$

- In the case of a natural link function it follows from Lemma 4.2 that

$$\frac{d\mu_i}{d\eta_i} = b^{(2)}(\theta_i) = \frac{1}{\tau^2} \sigma_i^2 \quad \text{or} \quad \mathbf{Z}(\boldsymbol{\beta}) = \frac{1}{\tau^4} \mathbf{V}(\boldsymbol{\beta}).$$

- Then the iteration equation (44) has the form:

$$\widehat{\boldsymbol{\beta}}_{k+1} = \widehat{\boldsymbol{\beta}}_k + \tau^2 (\mathbf{X}^\top \mathbf{V}(\widehat{\boldsymbol{\beta}}_k) \mathbf{X})^{-1} \left( \mathbf{X}^\top (\mathbf{Y} - \boldsymbol{\mu}(\widehat{\boldsymbol{\beta}}_k)) \right).$$

### Remark

- If the random sample  $\mathbf{Y}$  in (44) is replaced by the so-called *pseudorandom variable*

$$\mathbf{Y}(\boldsymbol{\beta}) = \mathbf{X}\boldsymbol{\beta} + \left( \frac{d\boldsymbol{\eta}}{d\boldsymbol{\mu}}(\boldsymbol{\beta}) \right) (\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})),$$

the iteration equation (44) can be written in the following form:

$$(\mathbf{X}^\top \mathbf{Z}(\widehat{\boldsymbol{\beta}}_k) \mathbf{X}) \widehat{\boldsymbol{\beta}}_{k+1} = \mathbf{X}^\top \mathbf{Z}(\widehat{\boldsymbol{\beta}}_k) \mathbf{Y}(\widehat{\boldsymbol{\beta}}_k).$$

- This equation can be considered as a *weighted normal equation* for  $\widehat{\boldsymbol{\beta}}_{k+1}$  with respect to the pseudorandom sample  $\mathbf{Y}(\widehat{\boldsymbol{\beta}}_k)$ , where the weights, i.e., the entries of the diagonal matrix  $\mathbf{Z}(\widehat{\boldsymbol{\beta}}_k)$  also depend on the  $k$ -th iteration  $\widehat{\boldsymbol{\beta}}_k$ .

## 4.3.4 Asymptotic Normality of ML Estimators; Asymptotic Tests

- The notion of the convergence in distribution of random vectors is defined as follows.
  - Let  $m \in \mathbb{N}$  be an arbitrary natural number and let  $\mathbf{Z}, \mathbf{Z}_1, \mathbf{Z}_2, \dots : \Omega \rightarrow \mathbb{R}^m$  be arbitrary random vectors. We say that  $\{\mathbf{Z}_n\}$  converges *in distribution* to  $\mathbf{Z}$  if

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathbf{Z}_n \leq \mathbf{x}) = \mathbb{P}(\mathbf{Z} \leq \mathbf{x}) \quad (45)$$

for each  $\mathbf{x} \in \mathbb{R}^m$  with  $\mathbb{P}(\mathbf{Z} = \mathbf{x}) = 0$ . Notation:  $\mathbf{Z}_n \xrightarrow{d} \mathbf{Z}$ .

- Now, we discuss asymptotic (distributional) properties of maximum-likelihood estimators  $\widehat{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$  or asymptotic tests if the sample size  $n$  tends to infinity.
  - Here we only consider the case of the natural link function  $g : G \rightarrow \mathbb{R}$
  - and index the random sample  $\mathbf{Y}$ , the loglikelihood function  $\log L(\mathbf{Y}, \boldsymbol{\beta})$ , the score vector  $\mathbf{U}(\boldsymbol{\beta})$ , the Fisher-information matrix  $\mathbf{I}(\boldsymbol{\beta})$  and the ML estimator  $\widehat{\boldsymbol{\beta}}$  each with  $n$ .

## 1. Asymptotic distributional properties

Under certain conditions (cf. Section VII.2.6 in Pruscha (2000)) one can show that: For each  $\boldsymbol{\beta} \in \mathbb{R}^m$  with  $\mathbf{x}_i^\top \boldsymbol{\beta} \in \Theta$  for  $i = 1, 2, \dots$  there exists

- a consistent ML estimator  $\widehat{\boldsymbol{\beta}}_n$  for  $\boldsymbol{\beta}$ , i.e., for each  $\varepsilon > 0$  it holds that

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\boldsymbol{\beta}}(|\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}| \leq \varepsilon, \mathbf{U}_n(\widehat{\boldsymbol{\beta}}_n) = \mathbf{o}) = 1, \quad (46)$$

- a sequence  $\{\boldsymbol{\Gamma}_n\}$  of invertible  $(m \times m)$ -matrices, which can depend on  $\boldsymbol{\beta}$  and for which it holds that  $\lim_{n \rightarrow \infty} \boldsymbol{\Gamma}_n = \mathbf{0}$ ,
- as well as a symmetric and positive definite  $(m \times m)$ -matrix  $\mathbf{K}(\boldsymbol{\beta})$ , such that

$$\lim_{n \rightarrow \infty} \boldsymbol{\Gamma}_n^\top \mathbf{I}_n(\boldsymbol{\beta}) \boldsymbol{\Gamma}_n = \mathbf{K}^{-1}(\boldsymbol{\beta}) \quad (47)$$

and

$$\boldsymbol{\Gamma}_n^{-1}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \xrightarrow{d} \mathbf{N}(\mathbf{o}, \mathbf{K}(\boldsymbol{\beta})) \quad \text{or} \quad 2(\log L_n(\mathbf{Y}_n, \widehat{\boldsymbol{\beta}}_n) - \log L_n(\mathbf{Y}_n, \boldsymbol{\beta})) \xrightarrow{d} \chi_m^2. \quad (48)$$

## 2. Asymptotic tests

- For large  $n$  the test statistic

$$T_n = 2(\log L_n(\mathbf{Y}_n, \widehat{\boldsymbol{\beta}}_n) - \log L_n(\mathbf{Y}_n, \boldsymbol{\beta}_0))$$

can be considered for the construction of an asymptotic test for the pair of hypotheses

$$H_0 : \boldsymbol{\beta} = \boldsymbol{\beta}_0 \quad \text{vs.} \quad H_1 : \boldsymbol{\beta} \neq \boldsymbol{\beta}_0.$$

Because of (48),  $H_0$  is rejected if  $T_n > \chi_{m, 1-\alpha}^2$ .

- The null hypothesis  $H_0 : \boldsymbol{\beta} = \mathbf{o}$  is particularly interesting. If it is rejected, more specific hypotheses can be tested, e.g., for each  $i = 1, \dots, m$  the hypothesis  $H_0 : \beta_i = 0$ .

**Remark**

- If  $\mathbf{I}_n(\boldsymbol{\beta})$  is positive definite for each sufficiently large  $n$  and if

$$\lim_{n \rightarrow \infty} \mathbf{I}_n^{-1}(\boldsymbol{\beta}) = \mathbf{0}, \quad (49)$$

then we can put  $\boldsymbol{\Gamma}_n = \mathbf{I}_n^{-1/2}$  in (47) and (48), which implies that  $\mathbf{K}(\boldsymbol{\beta})$  is the identity matrix.

- In (37) we have already shown that in the logistic regression model it holds that  $\mathbf{I}_n(\boldsymbol{\beta}) = \mathbf{X}^\top \text{diag}(\pi_i(1 - \pi_i))\mathbf{X}$ .
  - Since we assume that  $0 < \pi_i(\boldsymbol{\beta}) < 1$  for each  $i = 1, 2, \dots$  and that the design matrix  $\mathbf{X}$  has full column rank, the matrix  $\mathbf{I}_n(\boldsymbol{\beta})$  is positive definite in this case.
  - If furthermore  $\inf_{i \geq 1} \pi_i(1 - \pi_i) > 0$  and if the entries  $x_{ij}$  of the design matrix  $\mathbf{X}$  are chosen in such a way that  $\lim_{n \rightarrow \infty} (\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{0}$ , then (49) also holds.
- Now, let  $\mathbf{K}(\boldsymbol{\beta})$  be the identity matrix. Because of (47) and (48),  $H_0 : \beta_i = 0$  is rejected if

$$\frac{|\widehat{\boldsymbol{\beta}}_n|_i}{\sqrt{(\mathbf{I}_n^{-1}(\widehat{\boldsymbol{\beta}}_n))_{ii}}} > z_{1-\alpha/2}, \quad (50)$$

where  $z_{1-\alpha/2}$  is the  $(1 - \alpha/2)$ -quantile of the  $N(0, 1)$ -distribution.

**4.4 Weighted LS Estimator for Categorical Regression**

Instead of the maximum-likelihood approach to estimate the parameter vector  $\boldsymbol{\beta}$ , discussed in Section 4.3, we now consider a *weighted LS estimator* for  $\boldsymbol{\beta}$  for the categorical regression model.

**4.4.1 Estimation of the Expectation Vector**

*Recall* (cf. Section 4.2.2): In the binary categorical regression model all sample variables  $Y_1, \dots, Y_n$  are Bernoulli-distributed, i.e., they can only take the values 0 and 1 with positive probability.

- We use the notation

$$\pi_i = \mathbb{P}(Y_i = 1) \quad \left( = \mu_i = \mathbb{E}Y_i \right) \quad \forall i = 1, \dots, n,$$

where it is assumed that  $0 < \pi_i < 1$  for each  $i = 1, \dots, n$ .

- In this case the probabilities  $\pi_1, \dots, \pi_n$  are linked to the parameter vector  $\boldsymbol{\beta}$  by using a link function  $g : (0, 1) \rightarrow \mathbb{R}$ , i.e.,

$$(g(\pi_1), \dots, g(\pi_n))^\top = \mathbf{X}\boldsymbol{\beta}. \quad (51)$$

- In order to estimate the vectors  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_n)^\top$  or  $\mathbf{g}(\boldsymbol{\pi}) = (g(\pi_1), \dots, g(\pi_n))^\top$ , we assume that we are able to observe  $n_i > 0$  independent and identically distributed “copies”  $Y_{i1}, \dots, Y_{in_i}$  of  $Y_i$  for each  $i = 1, \dots, n$ . The total sample size is then equal to  $\sum_{i=1}^n n_i$ .
- For each  $i = 1, \dots, n$

$$\widehat{\pi}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} \quad (52)$$

is a natural estimator for  $\pi_i$ .

- This leads to the estimators  $\widehat{\boldsymbol{\pi}} = (\widehat{\pi}_1, \dots, \widehat{\pi}_n)^\top$  or  $\mathbf{g}(\widehat{\boldsymbol{\pi}}) = (g(\widehat{\pi}_1), \dots, g(\widehat{\pi}_n))^\top$  for  $\boldsymbol{\pi}$  or  $\mathbf{g}(\boldsymbol{\pi})$ , respectively.

One can easily see that the estimator  $\widehat{\boldsymbol{\pi}}$  is unbiased for  $\boldsymbol{\pi}$  and that its covariance matrix  $\mathbf{K}(\widehat{\boldsymbol{\pi}}) = (\text{Cov}(\widehat{\pi}_i, \widehat{\pi}_j))$  has the following form.

**Lemma 4.3** *It holds that*

$$\mathbb{E} \widehat{\boldsymbol{\pi}} = \boldsymbol{\pi}, \quad \text{Var} \widehat{\pi}_i = \pi_i(1 - \pi_i)/n_i \quad (53)$$

and

$$\mathbf{K}(\widehat{\boldsymbol{\pi}}) = \text{diag}(\text{Var} \widehat{\pi}_i). \quad (54)$$

**Proof** The statement follows directly from the fact that the random variables  $n_1\widehat{\pi}_1, \dots, n_n\widehat{\pi}_n$  are independent and binomially distributed with  $n_i\widehat{\pi}_i \sim \text{B}(n_i, \pi_i)$  for each  $i = 1, \dots, n$ .  $\square$

Furthermore, the following *central limit theorem* implies that the estimator  $\mathbf{g}(\widehat{\boldsymbol{\pi}}) = (g(\widehat{\pi}_1), \dots, g(\widehat{\pi}_n))^\top$  is asymptotically normally distributed.

**Theorem 4.3** *If  $n_i \rightarrow \infty$  for each  $i = 1, \dots, n$ , such that*

$$\frac{\sum_{j=1}^n n_j}{n_i} \rightarrow \lambda_i \in [1, \infty) \quad \forall i = 1, \dots, n, \quad (55)$$

then it holds that

$$\left( \sum_{j=1}^n n_j \right)^{1/2} (\mathbf{g}(\widehat{\boldsymbol{\pi}}) - \mathbf{g}(\boldsymbol{\pi})) \xrightarrow{\text{d}} \text{N}(\mathbf{0}, \mathbf{K}), \quad (56)$$

where

$$\mathbf{K} = \text{diag}(\alpha_i) \quad \text{and} \quad \alpha_i = \lambda_i (g^{(1)}(\pi_i))^2 \pi_i (1 - \pi_i). \quad (57)$$

**Proof**

- Since we assume that the link function  $g : (0, 1) \rightarrow \mathbb{R}$  is twice continuously differentiable, by Taylor series expansion we get that for each  $i = 1, \dots, n$

$$g(\widehat{\pi}_i) - g(\pi_i) = g^{(1)}(\pi_i)(\widehat{\pi}_i - \pi_i) + g^{(2)}(Z_i)(\widehat{\pi}_i - \pi_i)^2 = g^{(1)}(\pi_i)(\widehat{\pi}_i - \pi_i) + R_i,$$

where  $R_i = g^{(2)}(Z_i)(\widehat{\pi}_i - \pi_i)^2$  and  $Z_i : \Omega \rightarrow \mathbb{R}$  is a random variable taking values between  $\widehat{\pi}_i$  and  $\pi_i$ .

- From the central limit theorem for sums of independent and identically distributed random variables (cf. Theorem WR-5.16) it follows that

$$n_i^{1/2}(\widehat{\pi}_i - \pi_i) \xrightarrow{\text{d}} \text{N}(0, \pi_i(1 - \pi_i)) \quad \forall i = 1, \dots, n. \quad (58)$$

- Since  $\widehat{\pi}_i - \pi_i \rightarrow 0$  and therefore also  $Z_i - \pi_i \rightarrow 0$  or  $g^{(2)}(Z_i) \rightarrow g^{(2)}(\pi_i)$  with probability 1, it also holds that  $n_i^{1/2}R_i \xrightarrow{\text{P}} 0$  or

$$\left( \sum_{j=1}^n n_j \right)^{1/2} R_i = \left( \frac{\sum_{j=1}^n n_j}{n_i} \right)^{1/2} n_i^{1/2} R_i \xrightarrow{\text{P}} 0.$$

- Altogether, Slutsky's theorem (cf. Theorems WR-5.9 und WR-5.11) implies

$$\begin{aligned} \left( \sum_{j=1}^n n_j \right)^{1/2} (g(\widehat{\pi}_i) - g(\pi_i)) &= \left( \frac{\sum_{j=1}^n n_j}{n_i} \right)^{1/2} g^{(1)}(\pi_i) n_i^{1/2} (\widehat{\pi}_i - \pi_i) + \left( \sum_{j=1}^n n_j \right)^{1/2} R_i \\ &\xrightarrow{\text{d}} \text{N}(0, \lambda_i (g^{(1)}(\pi_i))^2 \pi_i (1 - \pi_i)). \end{aligned}$$

- As the random variables  $g(\widehat{\pi}_1) - g(\pi_1), \dots, g(\widehat{\pi}_i) - g(\pi_i)$  are independent, the statement is proved.  $\square$



#### 4.4.2 Asymptotic Normality of the LS–Estimator

The following approach for the estimation of the parameter vector  $\boldsymbol{\beta}$  is motivated by the form of the asymptotic covariance matrix  $\mathbf{K}$  in Theorem 4.3.

- In a similar way as is Section 2.1 we now consider the method of least squares for getting an estimator  $\widehat{\boldsymbol{\beta}}$  for the unknown regression coefficients  $\beta_1, \dots, \beta_m$ .
- A random vector  $\widehat{\boldsymbol{\beta}} = (\widehat{\beta}_1, \dots, \widehat{\beta}_m)^\top$  shall be determined, such that the *weighted squared error*

$$e(\boldsymbol{\beta}) = \sum_{i=1}^n \frac{(g(\widehat{\pi}_i) - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\widehat{\sigma}_{ii}^2} \quad (59)$$

gets minimal for  $\boldsymbol{\beta} = \widehat{\boldsymbol{\beta}}$ , where  $\widehat{\sigma}_{ii}^2 = (\sum_{j=1}^n n_j/n_i)(g^{(1)}(\widehat{\pi}_i))^2 \widehat{\pi}_i(1 - \widehat{\pi}_i)$  and it is assumed that the weights  $\widehat{\sigma}_{ii}^2$  are positive.

#### Remark

- The weighted sum  $e(\boldsymbol{\beta})$  of the squared residuals  $(g(\widehat{\pi}_i) - \mathbf{x}_i^\top \boldsymbol{\beta})^2$  in (59) can be written as follows: By using the notation  $\widehat{\mathbf{K}} = \text{diag}(\widehat{\sigma}_{ii}^2)$  it holds that

$$e(\boldsymbol{\beta}) = (\mathbf{g}(\widehat{\boldsymbol{\pi}}) - \mathbf{X}\boldsymbol{\beta})^\top \widehat{\mathbf{K}}^{-1} (\mathbf{g}(\widehat{\boldsymbol{\pi}}) - \mathbf{X}\boldsymbol{\beta}). \quad (60)$$

- In the same way as in the proof of Theorem 2.1 one can show that the weighted squared error  $e(\boldsymbol{\beta})$  is minimal if and only if  $\boldsymbol{\beta}$  is a solution of the following normal equation:

$$\mathbf{X}^\top \widehat{\mathbf{K}}^{-1} \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \widehat{\mathbf{K}}^{-1} \mathbf{g}(\widehat{\boldsymbol{\pi}}). \quad (61)$$

- Since the matrix  $\mathbf{X}^\top \widehat{\mathbf{K}}^{-1} \mathbf{X}$  is invertible, (61) has the uniquely determined solution

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \widehat{\mathbf{K}}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \widehat{\mathbf{K}}^{-1} \mathbf{g}(\widehat{\boldsymbol{\pi}}). \quad (62)$$

Now we show that the weighted LS–estimator  $\widehat{\boldsymbol{\beta}}$  in (62) is asymptotically normally distributed if the (sub–) sample sizes  $n_i$  grow unboundedly for each  $i = 1, \dots, n$ .

Here we need the following vectorial versions of Slutsky’s theorem (cf. Theorems WR–5.9 and WR–5.11) and of the “continuous mapping theorem” (cf. Theorem WR–5.12).

#### Lemma 4.4

- Let  $m \in \mathbb{N}$ , let  $\mathbf{Y}, \mathbf{Y}_n, \mathbf{Z}_n : \Omega \rightarrow \mathbb{R}^m$  be arbitrary random vectors over the very same probability space and let  $c \in \mathbb{R}^m$ .
- If  $\mathbf{Y}_n \xrightarrow{d} \mathbf{Y}$  and  $\mathbf{Z}_n \xrightarrow{d} c$ , then  $\mathbf{Y}_n + \mathbf{Z}_n \xrightarrow{d} \mathbf{Y} + c$  and  $\mathbf{Y}_n^\top \mathbf{Z}_n \xrightarrow{d} c^\top \mathbf{Y}$ .

#### Lemma 4.5

- Let  $m \in \mathbb{N}$ , let  $\mathbf{Z}, \mathbf{Z}_1, \mathbf{Z}_2, \dots : \Omega \rightarrow \mathbb{R}^m$  be arbitrary random vectors and let  $\varphi : \mathbb{R}^m \rightarrow \mathbb{R}$  be a continuous function.
- Then it holds that  $\varphi(\mathbf{Z}_n) \xrightarrow{d} \varphi(\mathbf{Z})$  provided that  $\mathbf{Z}_n \xrightarrow{d} \mathbf{Z}$ .

The *proofs* of Lemmas 4.4 and 4.5 are similar to the proofs of Theorems WR–5.9, WR–5.11 and WR–5.12. Therefore, they are omitted.

**Theorem 4.4** *If  $n_i \rightarrow \infty$  for each  $i = 1, \dots, n$ , such that*

$$\frac{\sum_{j=1}^n n_j}{n_i} \rightarrow \lambda_i \in [1, \infty) \quad \forall i = 1, \dots, n, \quad (63)$$

*then it holds that*

$$\left( \sum_{j=1}^n n_j \right)^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} N(\mathbf{o}, (\mathbf{X}^\top \mathbf{K}^{-1} \mathbf{X})^{-1}), \quad (64)$$

*where  $\mathbf{K} = \text{diag}(\alpha_i)$  is the diagonal matrix considered in Theorem 4.3.*

**Proof**

- It follows from the definition of  $\hat{\boldsymbol{\beta}}$  in (62) that

$$\begin{aligned} \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} &= (\mathbf{X}^\top \hat{\mathbf{K}}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{K}}^{-1} \mathbf{g}(\hat{\boldsymbol{\pi}}) - \boldsymbol{\beta} = (\mathbf{X}^\top \hat{\mathbf{K}}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{K}}^{-1} (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{X}\boldsymbol{\beta}) \\ &= (\mathbf{X}^\top \hat{\mathbf{K}}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{K}}^{-1} (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{g}(\boldsymbol{\pi})), \end{aligned}$$

where in the last equality we used that  $\mathbf{g}(\boldsymbol{\pi}) = \mathbf{X}\boldsymbol{\beta}$ ; cf. (51).

- In Theorem 4.3 we have already shown that

$$\left( \sum_{j=1}^n n_j \right)^{1/2} (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{g}(\boldsymbol{\pi})) \xrightarrow{d} N(\mathbf{o}, \mathbf{K}),$$

where the asymptotic covariance matrix  $\mathbf{K}$  is given in (57).

- Moreover, it holds that  $\hat{\mathbf{K}} \xrightarrow{P} \mathbf{K}$  if  $n_i \rightarrow \infty$  for each  $i = 1, \dots, n$ .
- Altogether with
  - Slutsky's theorem (cf. Lemma 4.4),
  - the “continuous mapping theorems” (cf. Lemma 4.5) as well as
  - Theorem 1.3 about linear transformations of normally distributed random vectors,

it follows that

$$\begin{aligned} \left( \sum_{j=1}^n n_j \right)^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) &= \left( \sum_{j=1}^n n_j \right)^{1/2} (\mathbf{X}^\top \hat{\mathbf{K}}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{K}}^{-1} (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{g}(\boldsymbol{\pi})) \\ &= \left( \sum_{j=1}^n n_j \right)^{1/2} \underbrace{(\mathbf{X}^\top \hat{\mathbf{K}}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{K}}^{-1} ((\mathbf{X}^\top \mathbf{K}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{K}^{-1})^{-1}}_{\xrightarrow{P} \mathbf{I}} \\ &\quad \times (\mathbf{X}^\top \mathbf{K}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{K}^{-1} (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{g}(\boldsymbol{\pi})) \\ &\xrightarrow{d} N(\mathbf{o}, ((\mathbf{X}^\top \mathbf{K}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{K}^{-1}) \mathbf{K} ((\mathbf{X}^\top \mathbf{K}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{K}^{-1})^\top) \\ &= N(\mathbf{o}, (\mathbf{X}^\top \mathbf{K}^{-1} \mathbf{X})^{-1}). \end{aligned} \quad \square$$

**Remark**

- If  $n_1 + \dots + n_n$  is a large number, the test statistic

$$T = \left( \sum_{j=1}^n n_j \right)^{1/2} \hat{\beta}_i / \sqrt{\tilde{k}_{ii}}$$

can be considered for the construction of an asymptotic test for the pair of hypotheses  $H_0 : \beta_i = 0$  vs.  $H_1 : \beta_i \neq 0$ , where  $\tilde{k}_{ii}$  is the  $i$ -th diagonal entry of the matrix  $\tilde{\mathbf{K}} = (\mathbf{X}^\top \hat{\mathbf{K}}^{-1} \mathbf{X})^{-1}$ .

- Because of Theorem 4.4,  $H_0$  is rejected if  $|T| > z_{1-\alpha/2}$ .

### 4.4.3 Evaluation of the Goodness of Fit

An important problem is the choice of a suitable design matrix  $\mathbf{X}$  in order to fit the model  $\mathbf{g}(\boldsymbol{\pi}) = \mathbf{X}\boldsymbol{\beta}$  in the best possible way to given data. It is possible to answer this question by using the following result.

**Theorem 4.5** *Under the conditions of Theorem 4.3 it holds that*

$$\sum_{j=1}^n n_j (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{X}\hat{\boldsymbol{\beta}})^\top \hat{\mathbf{K}}^{-1} (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{X}\hat{\boldsymbol{\beta}}) \xrightarrow{d} \chi_{n-m}^2. \quad (65)$$

#### Proof

- In Theorem 4.3 we have already shown that the random vector  $(\sum_{j=1}^n n_j)^{1/2} (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{X}\boldsymbol{\beta})$  is approximately  $N(\mathbf{o}, \mathbf{K})$  distributed.
- Since  $\hat{\boldsymbol{\beta}} \xrightarrow{P} \boldsymbol{\beta}$  and  $\hat{\mathbf{K}} \xrightarrow{P} \mathbf{K}$ , the statement follows with
  - Slutsky’s Theorem (cf. Lemma 4.4),
  - the “continuous mapping theorems” (cf. Lemma 4.5) as well as
  - Theorem 1.9 about quadratic forms of normally distributed random vectors. □

#### Remark

- Because of Theorem 4.5 the quantity  $\sum_{j=1}^n n_j (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{X}\hat{\boldsymbol{\beta}})^\top \hat{\mathbf{K}}^{-1} (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{X}\hat{\boldsymbol{\beta}})$  can be seen as a measure for the goodness of fit of the model  $\mathbf{g}(\boldsymbol{\pi}) = \mathbf{X}\boldsymbol{\beta}$  to given data.
- The goodness of fit is appraised as sufficiently good if

$$\sum_{j=1}^n n_j (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{X}\hat{\boldsymbol{\beta}})^\top \hat{\mathbf{K}}^{-1} (\mathbf{g}(\hat{\boldsymbol{\pi}}) - \mathbf{X}\hat{\boldsymbol{\beta}}) < \chi_{n-m, 1-\alpha}^2.$$

- On the other hand, the dimensions of  $\mathbf{X}$  should be as small as possible, which means in particular that for each  $i = 1, \dots, m$  the null hypothesis of the test  $H_0 : \beta_i = 0$  vs.  $H_1 : \beta_i \neq 0$  should be clearly rejected.

## 5 Goodness-of-Fit Tests

- In this chapter the sample variables are denoted by  $X_1, \dots, X_n$ , where we will assume from now on that  $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$  is a sequence of independent *and* identically distributed random variables.
- The assumptions, we made so far regarding the distribution  $P$  of the sample variables  $X_1, \dots, X_n$ , either have been *strictly qualitative* (discrete or absolutely continuous distribution) or *parametric*, where in the latter case it has been assumed
  - that  $P$  belongs to a parametric family  $\{P_\theta, \theta \in \Theta\}$  of distributions with  $\Theta \subset \mathbb{R}^m$  for some integer  $m \geq 1$ ,
  - and that merely the parameter vector  $\theta = (\theta_1, \dots, \theta_m)^\top$  or some of its components, respectively, are unknown.
- In the following, we discuss so-called *goodness-of-fit tests*.
  - To begin with, we consider a test to verify the hypothesis  $H_0 : P = P_0$  that the distribution  $P$  of the sample variables is equal to a given (hypothetical) distribution  $P_0$ .
  - Afterwards, we construct tests to check if  $P$  belongs to a given (parametric) *class* of distributions  $\{P_\theta, \theta \in \Theta\}$ .

### 5.1 Kolmogorov–Smirnov Test

#### 5.1.1 Empirical Distribution Function; KS Test Statistic

- There are different tests suggested in literature to verify the hypothesis  $H_0 : P = P_0$  that the distribution  $P$  of the independent and identically distributed random variables  $X_1, \dots, X_n$  is equal to a given distribution  $P_0$ .
- This kind of null hypothesis is considered for the *Kolmogorov–Smirnov test*, which is based on the analysis of the *empirical distribution function*  $\widehat{F}_n : \mathbb{R} \times \mathbb{R}^n \rightarrow [0, 1]$ , introduced in Section I-1.5, where

$$\widehat{F}_n(t; x_1, \dots, x_n) = \frac{\#\{i : 1 \leq i \leq n, x_i \leq t\}}{n}. \quad (1)$$

- The sample function  $T_n : \mathbb{R}^n \rightarrow [0, \infty)$  is considered with

$$T_n(x_1, \dots, x_n) = \sqrt{n} \sup_{t \in \mathbb{R}} |\widehat{F}_n(t; x_1, \dots, x_n) - F_0(t)|. \quad (2)$$

- In Section I-1.5.3 we have already shown that the distribution of the *KS test statistic*  $T_n(X_1, \dots, X_n)$  does not depend on  $P_0$  if it is assumed that the distribution function  $F_0 : \mathbb{R} \rightarrow [0, 1]$ , which corresponds to  $P_0$ , is continuous, cf. Theorem I-1.19.
- Let  $s_{n,1-\alpha}$  be the  $(1 - \alpha)$ -quantile of the distribution of  $T_n(X_1, \dots, X_n)$  under an arbitrary continuous distribution function  $F_0$ . The Kolmogorov–Smirnov test rejects the null hypothesis  $H_0 : P = P_0$  if

$$T_n(x_1, \dots, x_n) > s_{n,1-\alpha}. \quad (3)$$

#### Remark

- The quantiles  $s_{n,1-\alpha}$  can be determined using Monte Carlo simulation, where the distribution function  $F_0$  of the standard-uniform distribution on  $[0, 1]$  can be taken as a basis, cf Corollary I-1.3.
- If it is not assumed that  $F_0$  is continuous, then the decision rule, considered in (3), provides a test whose level can be smaller than  $\alpha$ .
- However, if it is possible to determine the quantile  $s'_{n,1-\alpha}$  of  $T_n(X_1, \dots, X_n)$  under  $F_0$ , e.g., using MC simulation, then  $T_n(x_1, \dots, x_n) > s'_{n,1-\alpha}$  is a test, which taps the full level  $\alpha$  even if  $F_0$  is discontinuous.

### 5.1.2 Asymptotic Distribution

We now analyze the asymptotic distribution of the KS test statistic  $T_n(X_1, \dots, X_n)$ , introduced in (2), as  $n \rightarrow \infty$ . To begin with, we provide some auxiliary tools.

In particular we need the following *continuity theorem* for characteristic functions of random vectors, which is a multidimensional generalization of Theorem WR-5.20 and which we will state without proof.

**Lemma 5.1** *Let  $m \in \mathbb{N}$  and let  $\mathbf{Z}, \mathbf{Z}_1, \mathbf{Z}_2, \dots : \Omega \rightarrow \mathbb{R}^m$  be arbitrary random vectors with characteristic functions  $\varphi_{\mathbf{Z}_n}$  and  $\varphi_{\mathbf{Z}}$ . It holds that  $\mathbf{Z}_n \xrightarrow{d} \mathbf{Z}$  if and only if*

$$\lim_{n \rightarrow \infty} \varphi_{\mathbf{Z}_n}(\mathbf{t}) = \varphi_{\mathbf{Z}}(\mathbf{t}) \quad \forall \mathbf{t} \in \mathbb{R}^m. \quad (4)$$

Moreover, we need a *multivariate central limit theorem* for sums of independent and identically distributed random vectors,

- whose proof can be reduced, using Lemma 5.1, to the corresponding central limit theorem for real-valued random variables (cf. Theorem WR-5.16).
- In literature this approach is sometimes called the *Cramèr–Wold device*.

#### Lemma 5.2

- Let  $m \in \mathbb{N}$  and let  $\mathbf{Z}_1, \mathbf{Z}_2, \dots : \Omega \rightarrow \mathbb{R}^m$  be a sequence of independent and identically distributed random vectors with expectation vector  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)^\top$  and covariance matrix  $\mathbf{K}$ .
- Then it holds that

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\frac{(\mathbf{Z}_1 + \dots + \mathbf{Z}_n) - n\boldsymbol{\mu}}{\sqrt{n}} \leq \mathbf{x}\right) = \Phi_{\mathbf{K}}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathbb{R}^m, \quad (5)$$

where  $\Phi_{\mathbf{K}} : \mathbb{R}^m \rightarrow [0, 1]$  denotes the distribution function of the  $N(\mathbf{o}, \mathbf{K})$ -distribution.

#### Proof

- Let  $\mathbf{Z}_n = (Z_{n1}, \dots, Z_{nm})^\top$ . Because of Lemma 5.1 the convergence in distribution stated in (5) is equivalent to

$$\lim_{n \rightarrow \infty} \varphi_n(\mathbf{t}) = \varphi(\mathbf{t}) \quad \forall \mathbf{t} \in \mathbb{R}^m, \quad (6)$$

– where  $\varphi_n(\mathbf{t})$  is the characteristic function of  $(\mathbf{Z}_1 + \dots + \mathbf{Z}_n - n\boldsymbol{\mu})/\sqrt{n}$  with

$$\varphi_n(\mathbf{t}) = \mathbb{E} \exp\left(i \sum_{j=1}^m t_j \frac{(Z_{1j} + \dots + Z_{nj}) - n\mu_j}{\sqrt{n}}\right)$$

– and  $\varphi(\mathbf{t})$  is the characteristic function of the  $N(\mathbf{o}, \mathbf{K})$ -distribution with

$$\varphi(\mathbf{t}) = \exp\left(-\frac{1}{2} \mathbf{t}^\top \mathbf{K} \mathbf{t}\right). \quad (7)$$

- Furthermore, one can easily see that

$$\varphi_n(\mathbf{t}) = \mathbb{E} \exp\left(i \sum_{k=1}^n \frac{\sum_{j=1}^m t_j (Z_{kj} - \mu_j)}{\sqrt{n}}\right) \quad \forall \mathbf{t} = (t_1, \dots, t_m)^\top \in \mathbb{R}^m \quad (8)$$

and

$$\mathbb{E}\left(\sum_{j=1}^m t_j (Z_{kj} - \mu_j)\right) = 0, \quad \text{Var}\left(\sum_{j=1}^m t_j (Z_{kj} - \mu_j)\right) = \mathbf{t}^\top \mathbf{K} \mathbf{t} \quad \forall k \in \mathbb{N}. \quad (9)$$

- If  $\mathbf{t}^\top \mathbf{K} \mathbf{t} = 0$ , then it follows from (9)
  - that  $\sum_{j=1}^m t_j (Z_{kj} - \mu_j) = 0$  with probability 1 for arbitrary  $k = 1, \dots, n$  and  $n \geq 1$ .
  - This and (7) – (8) imply that  $\varphi_n(\mathbf{t}) = 1 = \varphi(\mathbf{t})$  for each  $n \geq 1$ , i.e., (6).
- Now let  $\mathbf{t}^\top \mathbf{K} \mathbf{t} > 0$ .
  - From (8) it follows that  $\varphi_n(\mathbf{t})$  is equal to the value of the characteristic function of the real-valued random variable  $\sum_{k=1}^n \sum_{j=1}^m t_j (Z_{kj} - \mu_j) / \sqrt{n}$  at 1.
  - Moreover, it follows from (7) that  $\varphi(\mathbf{t})$  is the value of the characteristic function of the one-dimensional normal distribution  $N(0, \mathbf{t}^\top \mathbf{K} \mathbf{t})$  at 1.
  - On the other hand, Theorem WR-5.16, i.e., the (1-dimensional) central limit theorem for sums of independent and identically distributed (real-valued) random variables, implies that for  $n \rightarrow \infty$

$$\sum_{k=1}^n \frac{\sum_{j=1}^m t_j (Z_{kj} - \mu_j)}{\sqrt{n}} \xrightarrow{d} N(0, \mathbf{t}^\top \mathbf{K} \mathbf{t}). \quad (10)$$

- This, (7) – (8) and Theorem WR-5.20, i.e., the continuity theorem for characteristic functions of real-valued random variables, imply the validity of (6).  $\square$

The following limit theorem, already mentioned in Section I-1.5.3, provides an approximation formula for the distribution function of  $T_n(X_1, \dots, X_n)$  for a large sample size  $n$ .

**Theorem 5.1** *Let the distribution function  $F_0 : \mathbb{R} \rightarrow [0, 1]$  be continuous. Assuming that  $H_0 : P = P_0$  is true, it holds that*

$$\lim_{n \rightarrow \infty} \mathbb{P}(T_n(X_1, \dots, X_n) \leq x) = K(x) \quad \forall x \in \mathbb{R},$$

where  $K : \mathbb{R} \rightarrow [0, 1]$  is the distribution function of the so-called Kolmogorov distribution with

$$K(x) = \begin{cases} 1 - 2 \sum_{k=1}^{\infty} (-1)^{k-1} \exp(-2k^2 x^2), & \text{if } x > 0, \\ 0, & \text{if } x \leq 0. \end{cases} \quad (11)$$

### Proof

- We only sketch the idea of the proof since the full proof of Theorem 5.1 (cf., e.g., A. van der Vaart and J. Wellner (1996)) exceeds the scope of these lecture notes
  - as it requires profound tools from the theory of stochastic processes.
  - In particular, the term of convergence in distribution in *function spaces* as well as a so-called *functional* central limit theorem is needed,
  - which can be seen as an (infinite dimensional) generalization of the classical central limit theorems for sums of real-valued random variables (cf. Section WR-5.3) or of finite-dimensional random vectors (cf. Lemma 5.2).
- As the distribution of  $T_n(X_1, \dots, X_n)$  does not depend on  $F_0$  (cf. Theorem I-1.19), we can w.l.o.g. assume that  $F_0$  is the distribution function of the uniform distribution on  $[0, 1]$ , i.e.,  $F_0(t) = t$  for each  $t \in [0, 1]$ .
  - In order to analyze the asymptotic distribution of  $T_n(X_1, \dots, X_n)$  for  $n \rightarrow \infty$ , we use the abbreviating notation

$$B_n(t) = \sqrt{n} (\hat{F}_n(t; X_1, \dots, X_n) - F_0(t)) \quad \forall t \in [0, 1], \quad (12)$$

- where the family of random variables  $\{B_n(t), t \in [0, 1]\}$  is a stochastic process, which is called an *empirical process* in literature.

- For arbitrary  $t_1, \dots, t_m \in [0, 1]$  it then holds that

$$\sqrt{n}(B_n(t_1), \dots, B_n(t_m)) = \sum_{i=1}^n (Y_i(t_1) - t_1, \dots, Y_i(t_m) - t_m),$$

where

$$Y_i(t_j) = \begin{cases} 1, & \text{if } X_i \leq t_j, \\ 0, & \text{if } X_i > t_j. \end{cases}$$

- For each  $n \geq 1$  the random vector  $\sqrt{n}(B_n(t_1), \dots, B_n(t_m))$  can be written as a sum of  $n$  independent and identically distributed random vectors with expectation vector  $\mathbf{o}$ , whose covariance matrix  $\mathbf{K} = (\sigma_{ij}^2)$  is given by  $\sigma_{ij}^2 = \min\{t_i, t_j\} - t_i t_j$ .
- From Lemma 5.2 it now follows that for  $n \rightarrow \infty$

$$(B_n(t_1), \dots, B_n(t_m)) \xrightarrow{d} (B(t_1), \dots, B(t_m)), \quad (13)$$

where  $(B(t_1), \dots, B(t_m))$  is a normally distributed random vector with  $(B(t_1), \dots, B(t_m)) \sim N(\mathbf{o}, \mathbf{K})$ .

- This and the continuous mapping theorem for random vectors (cf. Lemma 4.5) imply that

$$\max_{i=1, \dots, m} \sqrt{n} |\widehat{F}_n(t_i; X_1, \dots, X_n) - F_0(t_i)| \xrightarrow{d} \max_{i=1, \dots, m} |B(t_i)|. \quad (14)$$

- It is easy to see that the distribution  $N(\mathbf{o}, \mathbf{K})$  of the random vector  $(B(t_1), \dots, B(t_m))$ 
  - can be considered as the finite-dimensional distribution of the so-called *Brownian bridge process*  $\{B(t), t \in [0, 1]\}$  with  $B(t) = X(t) - tX(1)$ , where  $\{X(t), t \in [0, 1]\}$  is a (standard-) Wiener-process,
  - i.e.,  $\{X(t), t \in [0, 1]\}$  is a stochastic process with continuous trajectories and independent increments, such that  $X(0) = 0$  and  $X(t_2) - X(t_1) \sim N(0, t_2 - t_1)$  for arbitrary  $t_1, t_2 \in [0, 1]$  with  $t_1 < t_2$ , cf. Section 2.4 of the lecture notes “Elementare Wahrscheinlichkeitsrechnung und Statistik”.
- Using the theory of convergence in distribution in function spaces as well as a corresponding functional central limit theorem, it is possible to show that not only the “finite-dimensional” convergences (13) and (14) hold but also

$$(B_n(t), t \in [0, 1]) \xrightarrow{d} (B(t), t \in [0, 1]) \quad (15)$$

and

$$\max_{t \in [0, 1]} \sqrt{n} |\widehat{F}_n(t; X_1, \dots, X_n) - F_0(t)| \xrightarrow{d} \max_{t \in [0, 1]} |B(t)|. \quad (16)$$

- Furthermore, one can show that the distribution function of the maximum  $\max_{t \in [0, 1]} |B(t)|$  of the Brownian bridge  $\{B(t), t \in [0, 1]\}$  is given by (11).  $\square$

### Remark

- Because of Theorem 5.1 the hypothesis  $H_0 : P = P_0$  is rejected for a sufficiently large sample size (as a rule of thumb it holds that  $n > 40$ , cf. the remark at the end of Section I-1.5.3) if

$$T_n(x_1, \dots, x_n) > \xi_{1-\alpha},$$

- where  $\xi_{1-\alpha}$  denotes the  $(1 - \alpha)$ -quantile of the Kolmogorov-distribution, given in (11), i.e.,  $\xi_{1-\alpha}$  is a solution of the equation  $K(\xi_{1-\alpha}) = 1 - \alpha$ .

### 5.1.3 Pointwise and Uniform Consistency

In this section we consider some properties of the Kolmogorov–Smirnov test.

In order to show the (pointwise) consistency of the KS test, we need the *Glivenko–Cantelli theorem* (cf. Theorem I-1.18), i.e.,

$$\mathbb{P}_{F_0} \left( \lim_{n \rightarrow \infty} \sup_{t \in \mathbb{R}} |\widehat{F}_n(t; X_1, \dots, X_n) - F_0(t)| = 0 \right) = 1. \quad (17)$$

**Theorem 5.2** *Let the distribution function  $F_0 : \mathbb{R} \rightarrow [0, 1]$  be continuous. Then the Kolmogorov–Smirnov test is pointwise consistent for each distribution function  $F$  of the sample variables with  $F \neq F_0$ , i.e., it holds that*

$$\lim_{n \rightarrow \infty} \mathbb{P}_F (T_n(X_1, \dots, X_n) > s_{n,1-\alpha}) = 1. \quad (18)$$

#### Proof

- From (17) it follows that for each  $F \neq F_0$

$$\mathbb{P}_F \left( \lim_{n \rightarrow \infty} \sup_{t \in \mathbb{R}} |\widehat{F}_n(t; X_1, \dots, X_n) - F_0(t)| > 0 \right) = 1.$$

- This implies that  $T_n(X_1, \dots, X_n) \rightarrow \infty$  with probability 1 under  $F \neq F_0$ .
- Since  $s_{n,1-\alpha} \rightarrow \xi_{1-\alpha} < \infty$  for  $n \rightarrow \infty$ , where  $\xi_{1-\alpha}$  is the  $(1-\alpha)$ -quantile of the Kolmogorov distribution, given in (11), it also holds that  $T_n(X_1, \dots, X_n) - (s_{n,1-\alpha} - \xi_{1-\alpha}) \xrightarrow{\text{a.s.}} \infty$  and therefore

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}_F (T_n(X_1, \dots, X_n) > s_{n,1-\alpha}) &= \lim_{n \rightarrow \infty} \mathbb{P}_F (T_n(X_1, \dots, X_n) - (s_{n,1-\alpha} - \xi_{1-\alpha}) > \xi_{1-\alpha}) \\ &= \lim_{n \rightarrow \infty} \mathbb{P}_F (T_n(X_1, \dots, X_n) > \xi_{1-\alpha}) = 1. \quad \square \end{aligned}$$

#### Remark

- As an strengthening of Theorem 5.2 one can show that the KS test is also *uniformly consistent* if the *Kolmogorov distance*

$$d_K(\Delta_n; F_0) = \inf_{F \in \Delta_n} \sup_{t \in \mathbb{R}} |F(t) - F_0(t)| \quad (19)$$

of the family  $\Delta_n$  of alternative distribution functions and the (hypothetical) distribution function  $F_0$  does not converge too fast to 0 as the sample size  $n$  increases.

- In this context, we need the following strengthening of Glivenko–Cantelli’s theorem, which is called *inequality of Dworetsky–Kiefer–Wolfowitz* in literature and which we state without proof.

**Lemma 5.3** *For arbitrary  $c > 0$  and  $n \geq 1$  it holds that*

$$\mathbb{P}_F \left( \sup_{t \in \mathbb{R}} |\widehat{F}_n(t; X_1, \dots, X_n) - F(t)| > c \right) \leq C \exp(-2nc^2), \quad (20)$$

where  $C \leq 2$  is an universal constant, which does not depend on  $F$ .

#### Remark

- From Lemma 5.3 it follows that for each  $\varepsilon > 0$  there is a  $c' > 0$  which does not depend on  $F$  and which fulfills

$$\inf_{n \geq 1} \mathbb{P}_F \left( \sqrt{n} \sup_{t \in \mathbb{R}} |\widehat{F}_n(t; X_1, \dots, X_n) - F(t)| \leq c' \right) \geq 1 - \varepsilon. \quad (21)$$



- In order to see this, it is sufficient to choose the threshold value  $c$  in (20) for  $\varepsilon \in (0, 1)$  in such a way that  $c = c'/\sqrt{n}$ , where

$$c' = \sqrt{-\frac{1}{2} \log\left(\frac{\varepsilon}{C}\right)}.$$

- Since  $c'$  does not depend on  $F$ , it also holds that for each  $\varepsilon > 0$  there is a  $c' > 0$  such that

$$\inf_{n \geq 1} \mathbb{P}_{F_n} \left( \sqrt{n} \sup_{t \in \mathbb{R}} |\widehat{F}_n(t; X_1, \dots, X_n) - F_n(t)| \leq c' \right) \geq 1 - \varepsilon, \quad (22)$$

where  $\{F_n\}$  is an arbitrary sequence of distribution functions.

With these tools it is possible to show the uniform consistence property of the KS test for the case that the Kolmogorov distance  $d_K(\Delta_n; F_0)$  of the family  $\Delta_n$  of the alternative distribution functions and the (hypothetical) distribution function  $F_0$  does not converge too fast to 0 as the sample size  $n$  increases.

**Theorem 5.3** *If there is a sequence  $\{\delta_n\}$  of positive numbers with  $\delta_n \rightarrow \infty$  such that*

$$\sqrt{n} d_K(\Delta_n; F_0) \geq \delta_n \quad \forall n \geq 1, \quad (23)$$

*then it holds that*

$$\lim_{n \rightarrow \infty} \inf_{F \in \Delta_n} \mathbb{P}_F(T_n(X_1, \dots, X_n) > s_{n,1-\alpha}) = 1. \quad (24)$$

**Proof**

- Let  $\{\delta_n\}$  be a sequence of positive numbers with  $\delta_n \rightarrow \infty$ , which fulfills (23), and let  $\{F_n\}$  be an arbitrary sequence of distribution functions such that for  $n \geq 1$

$$F_n \in \Delta_n \quad \text{and therefore} \quad \sqrt{n} d_K(F_n; F_0) \geq \delta_n, \quad (25)$$

where  $d_K(F_n; F_0) = \sup_{t \in \mathbb{R}} |F_n(t) - F_0(t)|$ .

- It is sufficient to show that

$$\lim_{n \rightarrow \infty} \mathbb{P}_{F_n}(T_n(X_1, \dots, X_n) > s_{n,1-\alpha}) = 1. \quad (26)$$

– From the triangle inequality it follows that

$$d_K(F_n, F_0) \leq d_K(F_n, \widehat{F}_n) + d_K(\widehat{F}_n, F_0).$$

– This and (25) imply that

$$T_n(X_1, \dots, X_n) \geq \delta_n - \sqrt{n} d_K(F_n, \widehat{F}_n).$$

– Therefore, it holds that

$$\mathbb{P}_{F_n}(T_n(X_1, \dots, X_n) > s_{n,1-\alpha}) \geq \mathbb{P}_{F_n}(\sqrt{n} d_K(F_n, \widehat{F}_n) < \delta_n - s_{n,1-\alpha}). \quad (27)$$

- Since  $s_{n,1-\alpha} \rightarrow \xi_{1-\alpha} < \infty$  and therefore  $\delta_n - s_{n,1-\alpha} \rightarrow \infty$  for  $n \rightarrow \infty$ , formulas (22) and (27) imply the validity of (26).  $\square$

**Remark**

- In particular, condition (23) is fulfilled if  $d_K(\Delta_n; F_0) \geq \delta$  for each  $n \geq 1$  and  $\delta > 0$  is a constant, which does not depend on  $n$ .

- If  $\sqrt{n} d_K(F_n; F_0) \geq \delta_n$  and  $\delta_n > s_{n,1-\alpha}$ , then (27) implies that for each  $n \geq 1$  holds the following (non-asymptotic) lower threshold for the *power* of the KS test:

$$\mathbb{P}_{F_n}(T_n(X_1, \dots, X_n) > s_{n,1-\alpha}) \geq 1 - 2 \exp\left(-2(\delta_n - s_{n,1-\alpha})^2\right), \quad (28)$$

where we have obtained a lower bound of the right-hand side of the inequality in (27) by applying Lemma 5.3 (for  $F = F_n$ ).

- However, notice that for a given (finite) sample size  $n < \infty$  it is possible that the “rejection probability”  $\mathbb{P}_{F_0}(T_n(X_1, \dots, X_n) > s_{n,1-\alpha})$  is smaller than  $\alpha$ .

On the other hand, the (asymptotic) power of the KS test can become arbitrarily small, i.e., arbitrarily close to  $\alpha$ , if the Kolmogorov distance  $d_K(\Delta_n; F_0)$  of the family  $\Delta_n$  of alternative distribution functions and the (hypothetical) distribution function  $F_0$  converges sufficiently fast to 0 as the sample size  $n$  increases.

#### Theorem 5.4

- Let  $\{F_n\}$  be an arbitrary sequence of continuous distribution functions such that

$$\lim_{n \rightarrow \infty} \sqrt{n} d_K(F_n; F_0) = 0. \quad (29)$$

- Then it holds that

$$\limsup_{n \rightarrow \infty} \mathbb{P}_{F_n}(T_n(X_1, \dots, X_n) > s_{n,1-\alpha}) \leq \alpha. \quad (30)$$

#### Proof

- From the triangle inequality it follows that

$$\mathbb{P}_{F_n}(T_n(X_1, \dots, X_n) > s_{n,1-\alpha}) \leq \mathbb{P}_{F_n}(\sqrt{n} d_K(\hat{F}_n; F_n) + \sqrt{n} d_K(F_n; F_0) > s_{n,1-\alpha}).$$

- From this and (29) the validity of (30) follows because the distribution of  $\sqrt{n} d_K(\hat{F}_n; F_n)$  under  $F_n$  does not depend on  $n$ .  $\square$

## 5.2 $\chi^2$ -Goodness-of-Fit Test

We now discuss an asymptotic goodness-of-fit test, where a test statistic is considered, which is approximately  $\chi^2$ -distributed for a large sample size. However, in this context the hypothesis

$$H_0 : P = P_0 \quad (\text{versus } H_1 : P \neq P_0), \quad (31)$$

analyzed in Section 5.1, is usually not considered since we “coarsen” the model of the random sample  $(X_1, \dots, X_n)$  by use of aggregation.

### 5.2.1 Aggregation; Pearson-Statistic

- For a (sufficiently large) natural number  $r$  we partition the range of the random variables  $X_1, \dots, X_n$  into  $r$  classes  $(a_1, b_1], \dots, (a_r, b_r]$  with

$$-\infty \leq a_1 < b_1 = a_2 < b_2 = \dots = a_r < b_r \leq \infty.$$

- Instead of the sample variables  $X_1, \dots, X_n$  we consider the “class sizes”  $Z_1, \dots, Z_r$ , where

$$Z_j = \#\{i : 1 \leq i \leq n, a_j < X_i \leq b_j\} \quad \forall j = 1, \dots, r. \quad (32)$$

To begin with, we show that the random vector  $(Z_1, \dots, Z_r)$  has a multinomial distribution with parameters  $n \geq 1$  and

$$\mathbf{p} = (p_1, \dots, p_{r-1})^\top \in [0, 1]^{r-1}, \quad \text{where } p_j = \mathbb{P}(a_j < X_1 \leq b_j) \quad \forall j = 1, \dots, r-1.$$

**Lemma 5.4** For arbitrary natural numbers  $k_1, \dots, k_r \geq 0$  with  $k_1 + \dots + k_r = n$  it holds that

$$\mathbb{P}(Z_1 = k_1, \dots, Z_r = k_r) = \frac{n!}{k_1! \cdot \dots \cdot k_r!} p_1^{k_1} \dots p_r^{k_r}, \quad (33)$$

where  $p_r = 1 - (p_1 + \dots + p_{r-1})$ .

**Proof**

- Since the random variables  $X_1, \dots, X_n$  are independent and identically distributed, it holds that

$$\mathbb{P}(X_1 \in (a_{i_1}, b_{i_1}], \dots, X_n \in (a_{i_n}, b_{i_n}]) = \prod_{j=1}^n \mathbb{P}(a_{i_j} < X_1 \leq b_{i_j}) = p_1^{k_1} \dots p_r^{k_r} \quad (34)$$

for each sequence of intervals  $(a_{i_1}, b_{i_1}], \dots, (a_{i_n}, b_{i_n}]$ , which contains  $k_1$ -times the interval  $(a_1, b_1], \dots, k_r$ -times the interval  $(a_r, b_r]$ .

- The statement (33) follows from summation of the probabilities, considered in (34), over all permutations of sequences  $(a_{i_1}, b_{i_1}], \dots, (a_{i_n}, b_{i_n}]$  of this kind.  $\square$

**Remark**

- We denote the multinomial distribution with the parameters  $n \geq 1$  and  $\mathbf{p} = (p_1, \dots, p_{r-1})^\top \in [0, 1]^{r-1}$  by  $M_{r-1}(n, \mathbf{p})$ . It is easy to see that for  $r = 2$  the multinomial distribution  $M_1(n, p_1)$  coincides with the binomial distribution  $\text{Bin}(n, p_1)$ .
- Instead of analyzing the test problem (31), we verify the hypothesis  $H_0 : \mathbf{p} = \mathbf{p}_0$  (against the alternative  $H_1 : \mathbf{p} \neq \mathbf{p}_0$ ) for a given (hypothetical) parameter vector

$$\mathbf{p}_0 = (p_{01}, \dots, p_{0,r-1})^\top \in (0, 1)^{r-1} \quad \text{with} \quad \sum_{i=1}^{r-1} p_{0i} < 1.$$

- We thus partition the family  $\Delta$  of all considered distributions of the sample variables  $X_1, \dots, X_n$  into the subsets

$$\Delta_0 = \{P : \mathbb{P}_P(a_j < X_1 \leq b_j) = p_{0j} \forall j = 1, \dots, r-1\} \quad \text{and} \quad \Delta_1 = \Delta \setminus \Delta_0. \quad (35)$$

- In this context we consider the sample function  $T_n : \mathbb{R}^n \rightarrow [0, \infty)$  with

$$T_n(x_1, \dots, x_n) = \sum_{j=1}^r \frac{(Z_j(x_1, \dots, x_n) - np_{0j})^2}{np_{0j}}, \quad (36)$$

- where  $Z_j(x_1, \dots, x_n)$  denotes the number of the sample values  $x_1, \dots, x_n$ , which belong to  $(a_j, b_j]$ .
- Assuming  $H_0 : \mathbf{p} = \mathbf{p}_0$  is true, it holds that  $\mathbb{E} Z_j(X_1, \dots, X_n) = np_{0j}$  for each  $j \in \{1, \dots, r\}$ .
  - Therefore, it makes sense to reject the hypothesis  $H_0 : \mathbf{p} = \mathbf{p}_0$  if  $T_n(x_1, \dots, x_n)$  is significantly larger than 0.
  - In order to make a decision, we need knowledge of the distribution of the test statistic  $T_n(X_1, \dots, X_n)$ , introduced in (36), which is called the *Pearson-statistic*.

### 5.2.2 Asymptotic Distribution

We show that  $T_n(X_1, \dots, X_n)$  converges in distribution to a  $\chi^2$ -distribution with  $r - 1$  degrees of freedom if  $n \rightarrow \infty$ . This is the base of the  $\chi^2$ -goodness-of-fit test, which has been introduced by Karl Pearson (1857–1936).

**Theorem 5.5** *For each  $P \in \Delta_0$  it holds that*

$$\lim_{n \rightarrow \infty} \mathbb{P}_P(T_n(X_1, \dots, X_n) > \chi_{r-1, 1-\alpha}^2) = \alpha, \quad \forall \alpha \in (0, 1), \quad (37)$$

where  $\chi_{r-1, 1-\alpha}^2$  denotes the  $(1 - \alpha)$ -quantile of the  $\chi^2$ -distribution with  $r - 1$  degrees of freedom.

#### Proof

- In Lemma 5.4 we have shown that the random vector  $\mathbf{Z}_n = (Z_{n1}, \dots, Z_{nr})^\top$ , given in (32), where  $Z_{nj} = Z_j(X_1, \dots, X_n)$ , has a multinomial distribution under  $P \in \Delta_0$  with parameters  $n \in \mathbb{N}$  and

$$\mathbf{p}_0 = (p_{01}, \dots, p_{0,r-1})^\top \in [0, 1]^{r-1}, \quad \text{where } p_{0j} = \mathbb{P}_P(a_j < X_1 \leq b_j) \quad \forall j = 1, \dots, r-1.$$

- This in particular implies that for arbitrary  $i, j \in \{1, \dots, r\}$

$$\mathbb{E} Z_{ni} = np_{0i}, \quad \text{Cov}(Z_{ni}, Z_{nj}) = \begin{cases} -np_{0i}p_{0j}, & \text{if } i \neq j, \\ np_{0i}(1 - p_{0i}), & \text{if } i = j. \end{cases} \quad (38)$$

- Moreover, it follows from (32) that  $Z_j = \sum_{i=1}^n \mathbb{1}_{\{a_j < X_i \leq b_j\}}$ , i.e.,  $\mathbf{Z}_n$  can be written as a sum of  $n$  independent and identically distributed random vectors, where  $\mathbb{1}_{\{a_j < X_i \leq b_j\}}$  is the indicator of the event  $\{a_j < X_i \leq b_j\}$ .
- With the notation

$$\mathbf{Z}'_n = \left( \frac{Z_{n1}}{\sqrt{n}} - \sqrt{np_{01}}, \dots, \frac{Z_{n,r-1}}{\sqrt{n}} - \sqrt{np_{0,r-1}} \right)^\top \quad (39)$$

it therefore follows from Lemma 5.2 that for  $n \rightarrow \infty$  it holds that

$$\mathbf{Z}'_n \xrightarrow{d} \mathbf{Z}' \sim N(\mathbf{o}, \mathbf{K}), \quad (40)$$

- where the  $(r - 1)$ -dimensional random vector  $\mathbf{Z}'$  has a (nondegenerate) multivariate normal distribution,
- whose covariance matrix  $\mathbf{K} = (\sigma_{ij}^2)$  is given by

$$\sigma_{ij}^2 = \begin{cases} -p_{0i}p_{0j}, & \text{if } i \neq j, \\ p_{0i}(1 - p_{0i}), & \text{if } i = j. \end{cases} \quad (41)$$

- It is easy to see that  $\mathbf{K}$  is invertible, where the entries  $a_{ij}$  of the inverse matrix  $\mathbf{A} = \mathbf{K}^{-1}$  are given by

$$a_{ij} = \begin{cases} \frac{1}{p_{0r}}, & \text{if } i \neq j, \\ \frac{1}{p_{0i}} + \frac{1}{p_{0r}}, & \text{if } i = j. \end{cases} \quad (42)$$

- From (40) and the properties of linear transformations of normally distributed random vectors (cf. Theorem 1.3) it now follows with Lemma 4.5 that  $\mathbf{A}^{1/2} \mathbf{Z}'_n \xrightarrow{d} N(\mathbf{o}, \mathbf{I}_{r-1})$ , where  $\mathbf{I}_{r-1}$  is the  $(r - 1) \times (r - 1)$ -dimensional identity matrix.
- The repeated application of Lemma 4.5 yields

$$(\mathbf{A}^{1/2} \mathbf{Z}'_n)^\top (\mathbf{A}^{1/2} \mathbf{Z}'_n) \xrightarrow{d} \chi_{r-1}^2.$$

- Now it is sufficient to note that

$$(\mathbf{A}^{1/2}\mathbf{Z}'_n)^\top (\mathbf{A}^{1/2}\mathbf{Z}'_n) = T_n(X_1, \dots, X_n).$$

- It namely holds that

$$\begin{aligned} (\mathbf{A}^{1/2}\mathbf{Z}'_n)^\top (\mathbf{A}^{1/2}\mathbf{Z}'_n) &= (\mathbf{Z}'_n)^\top \mathbf{A} \mathbf{Z}'_n \\ &= n \sum_{j=1}^{r-1} \frac{1}{p_{0j}} \left( \frac{Z_{nj}}{n} - p_{0j} \right)^2 + \frac{n}{p_{0r}} \sum_{i=1}^{r-1} \sum_{j=1}^{r-1} \left( \frac{Z_{ni}}{n} - p_{0i} \right) \left( \frac{Z_{nj}}{n} - p_{0j} \right), \end{aligned}$$

- where the second summand of the last term can be written in the form

$$\frac{n}{p_{0r}} \sum_{i=1}^{r-1} \sum_{j=1}^{r-1} \left( \frac{Z_{ni}}{n} - p_{0i} \right) \left( \frac{Z_{nj}}{n} - p_{0j} \right) = \frac{n}{p_{0r}} \left( \sum_{j=1}^{r-1} \left( \frac{Z_{nj}}{n} - p_{0j} \right) \right)^2 = \frac{n}{p_{0r}} \left( \frac{Z_{nr}}{n} - p_{0r} \right)^2,$$

- because it obviously holds that  $\sum_{j=1}^{r-1} Z_{nj} = n - Z_{nr}$  and  $\sum_{j=1}^{r-1} p_{0j} = 1 - p_{0r}$ .  $\square$

### Remark

- In order to practically use the  $\chi^2$ -goodness-of-fit test for the verification of the hypothesis  $H_0 : \mathbf{p} = \mathbf{p}_0$ , first the value of the test statistic  $T_n(x_1, \dots, x_n)$ , defined in (36), has to be computed.
- For sufficiently large sample sizes  $n$  the hypothesis  $H_0 : \mathbf{p} = \mathbf{p}_0$  is rejected if

$$T_n(x_1, \dots, x_n) > \chi_{r-1, 1-\alpha}^2,$$

- where  $\chi_{r-1, 1-\alpha}^2$  denotes the  $(1-\alpha)$ -quantile of the  $\chi^2$ -distribution with  $(r-1)$  degrees of freedom.
- A rule of thumb for  $n$  being sufficiently large, is the validity of the inequality  $np_{0,j} \geq a$  for each  $j \in \{1, \dots, r\}$  and for a constant  $a > 0$ .
- In literature there are different opinions about the required size of  $a > 0$ , which range from  $a = 2$  to  $a = 5$ . Some authors even demand  $a = 10$ .
- Other authors think that for a large number of classes (about  $r \geq 10$ ) the approximation is sufficiently good, even for  $a = 1$ .

### 5.2.3 Goodness-of-Fit; Local Alternatives

It is not difficult to show the following (pointwise) consistency of the  $\chi^2$ -goodness-of-fit test.

**Theorem 5.6** *The  $\chi^2$ -goodness-of-fit test is pointwise consistent for each vector  $\mathbf{p} = (p_1, \dots, p_{r-1})^\top \in [0, 1]^{r-1}$  with  $\mathbf{p} \neq \mathbf{p}_0$ , i.e., it holds that*

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\mathbf{p}}(T_n(X_1, \dots, X_n) > \chi_{r-1, 1-\alpha}^2) = 1. \quad (43)$$

### Proof

- From  $\mathbf{p} \neq \mathbf{p}_0$  it follows that

$$p_j \neq p_{0,j} \quad (44)$$

for some  $j \in \{1, \dots, r-1\}$ .

- Furthermore, it follows from the strong law of large numbers (cf. Theorem WR-5.15) that  $Z_{nj}/n \xrightarrow{\text{a.s.}} p_j$  for  $n \rightarrow \infty$  under  $\mathbb{P}_{\mathbf{p}}$ .

- This and (44) imply that under  $\mathbb{P}_{\mathbf{p}}$

$$T_n(X_1, \dots, X_n) \geq n \left( \frac{Z_{nj}}{n} - p_{0,j} \right)^2 \xrightarrow{\text{a.s.}} \infty.$$

- Hence, the validity of (43) is shown.  $\square$

### Remark

- Instead of considering a fixed vector  $\mathbf{p} \neq \mathbf{p}_0$ , there are also local alternatives  $\mathbf{p}_n = (p_{n1}, \dots, p_{n,r-1})^\top$  possible of the form

$$p_{nj} = p_{0j} + \frac{h_j}{\sqrt{n}} \quad \forall j = 1, \dots, r-1, \quad (45)$$

which may depend on the sample size  $n$ , where

$$\sum_{j=1}^r h_j = 0. \quad (46)$$

- Then one can show that for  $n \rightarrow \infty$  the asymptotic power of the  $\chi^2$ -goodness-of-fit test vs. such alternatives may be smaller than 1.

In order to prove the statement, we need the following estimate, which is called the *Berry-Esseen theorem* in literature.

**Lemma 5.5** *Let  $Y_1, Y_2, \dots : \Omega \rightarrow \mathbb{R}$  be a sequence of independent and identically distributed random variables with  $\mathbb{E}(|Y_1|^3) < \infty$ . If  $\mathbb{E}Y_1 = 0$  and  $\text{Var}Y_1 = 1$ , then it holds for each  $n \geq 1$*

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P} \left( \frac{Y_1 + \dots + Y_n}{\sqrt{n}} \leq x \right) - \Phi(x) \right| \leq C \frac{\mathbb{E}(|Y_1|^3)}{\sqrt{n}}, \quad (47)$$

where  $\Phi : \mathbb{R} \rightarrow [0, 1]$  denotes the distribution function of the  $N(0, 1)$ -distribution and  $C < \infty$  is a universal constant, which does not depend on the distribution of the random variables  $Y_1, Y_2, \dots$ .

**Theorem 5.7** *Let  $\{\mathbf{p}_n\}$  be a sequence of vectors, which is given by (45) and (46).*

- Then it holds for each  $x \geq 0$  that

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\mathbf{p}_n} (T_n(X_1, \dots, X_n) \leq x) = F_{r-1, \lambda}(x), \quad (48)$$

where  $F_{r-1, \lambda} : \mathbb{R} \rightarrow [0, 1]$  is the distribution function of the noncentral  $\chi^2$ -distribution with  $r-1$  degrees of freedom, whose noncentrality parameter  $\lambda$  is given by

$$\lambda = \sum_{j=1}^r \frac{h_j^2}{p_{0j}}. \quad (49)$$

- If  $h_j \neq 0$  for some  $j = 1, \dots, r$ , then the power of the  $\chi^2$ -goodness-of-fit test converges, in the case of the local alternatives  $\{\mathbf{p}_n\}$ , to a limit, which is larger than  $\alpha$  and smaller than 1, i.e.,

$$\alpha < \lim_{n \rightarrow \infty} \mathbb{P}_{\mathbf{p}_n} (T_n(X_1, \dots, X_n) > \chi_{r-1, 1-\alpha}^2) < 1. \quad (50)$$

**Proof**

- The proof of the first part is analogous to the proof of Theorem 5.5 since due to (45) and (46) it holds that

$$\mathbf{Z}'_n = \sqrt{n} \left( \frac{Z_{n1}}{n} - p_{n1}, \dots, \frac{Z_{n,r-1}}{n} - p_{n,r-1} \right)^\top + \mathbf{h}, \quad \text{where } \mathbf{h} = (h_1, \dots, h_{r-1})^\top. \quad (51)$$

for the random vector  $\mathbf{Z}'_n$ , introduced in (39).

- Because of (51) one can show in a similar way as in the proof of the multivariate central limit theorem in Lemma 5.2 that

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\mathbf{p}_n}(\mathbf{Z}'_n \leq \mathbf{x}) = F_{N(\mathbf{h}, \mathbf{K})}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathbb{R}^{r-1}. \quad (52)$$

- In this context it is sufficient to notice that one can show, by using Berry–Esseen's theorem in Lemma 5.5 that (analogously to formula (10) in the proof of Lemma 5.2)

– for arbitrary  $\mathbf{t} = (t_1, \dots, t_{r-1})^\top \in \mathbb{R}^{r-1}$  and  $x \in \mathbb{R}$  it holds that

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\mathbf{p}_n} \left( n^{-1/2} \sum_{j=1}^{r-1} t_j (Z_{nj} - np_{nj}) \leq x \right) = F_{N(0, \mathbf{t}^\top \mathbf{K} \mathbf{t})}(x),$$

– where  $\mathbf{K}$  is the covariance matrix, introduced in (41).

- In the same way as in the proof of Theorem 5.5 we now get from (52) that

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\mathbf{p}_n}(\mathbf{A}^{1/2} \mathbf{Z}'_n \leq \mathbf{x}) = F_{N(\mathbf{A}^{1/2} \mathbf{h}, \mathbf{I}_{r-1})}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathbb{R}^{r-1}.$$

– From this and from the definition of the noncentral  $\chi^2$ -distribution in Section 1.3.2 it follows that

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\mathbf{p}_n}((\mathbf{A}^{1/2} \mathbf{Z}'_n)^\top (\mathbf{A}^{1/2} \mathbf{Z}'_n) \leq x) = F_{r-1, \lambda}(x) \quad \forall x \in \mathbb{R},$$

– where  $\mathbf{A}$  is the inverse matrix  $\mathbf{A} = \mathbf{K}^{-1}$  in (42) and the noncentrality parameter  $\lambda$  is given by

$$\lambda = (\mathbf{A}^{1/2} \mathbf{h})^\top (\mathbf{A}^{1/2} \mathbf{h}) = \mathbf{h}^\top \mathbf{A} \mathbf{h} = \sum_{j=1}^r \frac{h_j^2}{p_{0j}}.$$

- Thus, (48) is proved and because of  $\alpha < 1 - F_{r-1, \lambda}(\chi_{r-1, 1-\alpha}^2) < 1$  also (50) is valid.  $\square$

**5.3 Pearson–Fisher Test**

- The null hypothesis  $H_0 : \mathbf{p} = \mathbf{p}_0$ , considered in Section 5.2, is in fact a compound hypothesis since it is equivalent to the hypothesis

$$H_0 : P \in \Delta_0,$$

where  $\Delta_0$  is the subset of distributions of the sample variables, which has been introduced in (35).

- If it shall be verified whether the distribution  $P$  of the independent and identically distributed sample variables  $X_1, \dots, X_n$  belongs to a given (parametric) *class of distributions*  $\{P_\theta, \theta \in \Theta\}$  with  $\Theta \subset \mathbb{R}^m$ , then we can proceed in a similar way as in the case of the  $\chi^2$ -goodness-of-fit test, discussed in Section 5.2.
- The sample function  $T_n : \mathbb{R}^n \rightarrow [0, \infty)$ , which has been considered in the definition of the Pearson–statistic  $T_n(X_1, \dots, X_n)$  in (36), is replaced by a modified sample function  $\hat{T}_n : \mathbb{R}^n \rightarrow [0, \infty)$ .

### 5.3.1 Pearson–Fisher Test Statistic

- In the same way as in Section 5.2.1, we “coarsen” the model, i.e.,
  - we partition the range of the sample variables  $X_1, \dots, X_n$  into  $r$  classes  $(a_1, b_1], \dots, (a_r, b_r]$  with  $-\infty \leq a_1 < b_1 = a_2 < b_2 = \dots = a_r < b_r \leq \infty$ , where  $r$  is a (sufficiently large) natural number.
  - Instead of the sample variables  $X_1, \dots, X_n$ , we once more consider the “class sizes”  $Z_1, \dots, Z_r$ , which have been introduced in (32), where

$$Z_j = \#\{i : 1 \leq i \leq n, a_j < X_i \leq b_j\} \quad \forall j = 1, \dots, r.$$

- According to Lemma 5.4 it holds that  $(Z_1, \dots, Z_r) \sim M_{r-1}(n, \mathbf{p})$ , where we now assume
  - that the parameter  $\mathbf{p} = (p_1, \dots, p_{r-1})^\top \in [0, 1]^{r-1}$  of the multinomial distribution  $M_{r-1}(n, \mathbf{p})$
  - is a (known) function  $\boldsymbol{\theta} \mapsto \mathbf{p}(\boldsymbol{\theta})$  of the (unknown) parameter vector  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_m)^\top \in \Theta \subset \mathbb{R}^m$  with  $m < r - 1$ .
- The hypothesis to be tested is  $H_0 : \mathbf{p} \in \{\mathbf{p}(\boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$ .
  - In order to be able to proceed with the verification of this hypothesis in a similar way as in Section 5.2, one first has to determine an estimator  $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_m)^\top$  for  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_m)^\top$ .
  - This also provides an estimator  $(\hat{p}_1, \dots, \hat{p}_r) = (p_1(\hat{\boldsymbol{\theta}}), \dots, p_r(\hat{\boldsymbol{\theta}}))$  for the probabilities  $(p_1, \dots, p_r) = (p_1(\boldsymbol{\theta}), \dots, p_r(\boldsymbol{\theta}))$ , where

$$p_j(\boldsymbol{\theta}) = \mathbb{P}_{\boldsymbol{\theta}}(a_j < X_1 \leq b_j) \quad \forall j = 1, \dots, r.$$

**Definition** The random variable  $\hat{T}_n(X_1, \dots, X_n)$ , which is given by the sample function  $\hat{T}_n : \mathbb{R}^n \rightarrow [0, \infty)$  with

$$\hat{T}_n(x_1, \dots, x_n) = \sum_{j=1}^r \frac{(Z_j(x_1, \dots, x_n) - n\hat{p}_j(x_1, \dots, x_n))^2}{n\hat{p}_j(x_1, \dots, x_n)} \quad (53)$$

is called *Pearson–Fisher statistic*.

#### Remark

- If the mapping  $\boldsymbol{\theta} \mapsto \mathbf{p}(\boldsymbol{\theta})$  is continuous and  $\hat{\boldsymbol{\theta}}$  is a (weakly) consistent estimator for  $\boldsymbol{\theta}$ ,
  - then it follows from the law of large numbers (cf. Theorem WR–5.15) that for arbitrary  $j \in \{1, \dots, r\}$  and  $\boldsymbol{\theta} \in \Theta$

$$\lim_{n \rightarrow \infty} \mathbb{E}_{\boldsymbol{\theta}} \left| \frac{1}{n} Z_j(X_1, \dots, X_n) - \hat{p}_j(X_1, \dots, X_n) \right| = 0.$$

- Therefore, it is reasonable to reject the null hypothesis  $H_0 : \mathbf{p} \in \{\mathbf{p}(\boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$  if  $\hat{T}_n(x_1, \dots, x_n)$  is significantly larger than 0.
- In order to be able to make this decision,
  - we first of all discuss conditions for the mapping  $\boldsymbol{\theta} \mapsto \mathbf{p}(\boldsymbol{\theta})$  which enable the construction of a sequence of consistent (maximum–likelihood) estimators  $\hat{\boldsymbol{\theta}}_n$  for  $\boldsymbol{\theta}$  that are asymptotically normally distributed.
  - Then we determine the (asymptotic limit) distribution of the test statistic  $\hat{T}_n(X_1, \dots, X_n)$ , introduced in (53), for  $n \rightarrow \infty$ .



### 5.3.2 Multivariate Central Limit Theorem for ML Estimators

In a similar way as in Section I-2.4.2, where the case  $m = 1$  has been considered, it is possible to derive a *multivariate central limit theorem* for consistent sequences of maximum-likelihood estimators for the parameter vector  $\theta$ .

In this context we need the following *regularity conditions*.

- The family  $\{P_\theta, \theta \in \Theta\}$  either consists only of discrete distributions or only of absolutely continuous distributions, where  $\Theta \subset \mathbb{R}^m$  is an open set.
- It holds that

$$P_\theta \neq P_{\theta'} \quad \text{if and only if} \quad \theta \neq \theta'.$$

- The set  $B = \{x \in \mathbb{R} : L(x; \theta) > 0\}$  does not depend on  $\theta \in \Theta$ , where the likelihood function  $L(x; \theta)$  is given by

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

and  $p(x; \theta)$  or  $f(x; \theta)$  is the probability function or density of  $P_\theta$ , respectively.

- Furthermore, let the mapping  $\theta \rightarrow L(x; \theta)$  for each  $x \in B$  be three times continuously differentiable and suppose that for each  $x \in B$  it holds that

$$\frac{\partial^k}{\partial \theta_{i_1} \dots \partial \theta_{i_k}} \int_B L(x; \theta) dx = \int_B \frac{\partial^k}{\partial \theta_{i_1} \dots \partial \theta_{i_k}} L(x; \theta) dx \quad \forall k \in \{1, 2, 3\}, i_1, \dots, i_k \in \{1, \dots, m\}, \theta \in \Theta, \quad (54)$$

where the integrals have to be replaced by the corresponding sums in the discrete case.

- For each  $\theta_0 \in \Theta$ , a constant  $c_{\theta_0} > 0$  and a measurable function  $g_{\theta_0} : B \rightarrow [0, \infty)$  exist such that for each triple  $(i_1, i_2, i_3) \in \{1, \dots, m\}^3$

$$\left| \frac{\partial^3}{\partial \theta_{i_1} \partial \theta_{i_2} \partial \theta_{i_3}} \log L(x; \theta) \right| \leq g_{\theta_0}(x) \quad \forall x \in B, \quad \forall \theta \in \Theta \text{ with } |\theta - \theta_0| < c_{\theta_0} \quad (55)$$

and

$$\mathbb{E}_{\theta_0} g_{\theta_0}(X_1) < \infty. \quad (56)$$

#### Remark

- *Recall* :
  - In general, the maximum-likelihood estimator  $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$  for  $\theta$  is defined as the solution of the following optimization problem (cf. Section I-2.2.2).
  - In this context  $\hat{\theta} : \mathbb{R}^n \rightarrow \Theta \subset \mathbb{R}^m$  is a sample function with

$$L(x_1, \dots, x_n; \theta) \leq L(x_1, \dots, x_n; \hat{\theta}(x_1, \dots, x_n)) \quad \forall (x_1, \dots, x_n) \in \mathbb{R}^n, \theta \in \Theta \quad (57)$$

and

$$L(x_1, \dots, x_n; \theta) = \begin{cases} p(x_1; \theta) \dots p(x_n; \theta) & \text{in the discrete case,} \\ f(x_1; \theta) \dots f(x_n; \theta) & \text{in the absolutely continuous case.} \end{cases}$$

- Under the above mentioned regularity conditions one can show that for arbitrary  $x_1, \dots, x_n \in \mathbb{R}$  the estimate  $\hat{\theta}(x_1, \dots, x_n)$  fulfills the following system of equations:

$$\frac{\partial}{\partial \theta_i} L(x_1, \dots, x_n; \hat{\theta}(x_1, \dots, x_n)) = 0 \quad \forall i = 1, \dots, m. \quad (58)$$

- To formulate the multivariate central limit theorem, we need the notion of the *Fisher-information matrix*, which has already been introduced in Section 4.3.1.
  - For each  $\boldsymbol{\theta} \in \Theta$  the  $m \times m$  matrix  $\mathbf{I}(\boldsymbol{\theta}) = (I_{ij}(\boldsymbol{\theta}))$  is considered with

$$I_{ij}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}} \left( \frac{\partial}{\partial \theta_i} \log L(X_1; \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log L(X_1; \boldsymbol{\theta}) \right), \quad (59)$$

- where it is assumed that the expectation in (59) exists for arbitrary  $i, j \in \{1, \dots, m\}$  (and is a finite real number).

As a generalization of Theorem I–2.11, where the 1-dimensional case has been considered, it is possible to derive the following multivariate central limit theorem for weakly consistent sequences of maximum-likelihood estimators  $\{\hat{\boldsymbol{\theta}}(X_1, \dots, X_n), n \geq 1\}$  for the parameter vector  $\boldsymbol{\theta}$  which fulfill the system of equations (58).

### Theorem 5.8

- Let the Fisher-information matrix  $\mathbf{I}(\boldsymbol{\theta})$  be positive definite (and therefore invertible) for each  $\boldsymbol{\theta} \in \Theta$  and let  $\{\hat{\boldsymbol{\theta}}(X_1, \dots, X_n), n \geq 1\}$  be a weakly consistent sequence of maximum-likelihood estimators for  $\boldsymbol{\theta}$ .
- Then it holds for  $n \rightarrow \infty$  that

$$\sqrt{n}(\hat{\boldsymbol{\theta}}(X_1, \dots, X_n) - \boldsymbol{\theta}) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\theta})). \quad (60)$$

The *proof* of Theorem 5.8 proceeds in a similar way as the proof of Theorem I–2.11. It is therefore omitted, cf. for instance E.L. Lehmann und G. Casella (1998) *The Theory of Point Estimation*, Springer-Verlag, New York.

### 5.3.3 Fisher-Information Matrix and Central Limit Theorem in the Coarsened Model

- We now return to the “coarsened” model, already considered in Section 5.3.1.
    - Here we assume that  $L : \mathbb{R} \times \Theta \rightarrow (0, 1)$  is the likelihood function with
$$L(x; \boldsymbol{\theta}) = p_j(\boldsymbol{\theta}), \quad \text{if } x \in (a_j, b_j], \quad (61)$$
  - where the probabilities  $p_j(\boldsymbol{\theta}) = \mathbb{P}_{\boldsymbol{\theta}}(a_j < X_1 \leq b_j)$  are positive and smaller than 1.
- Furthermore, we assume that the regularity conditions, formulated in Section 5.3.2, are fulfilled for the likelihood function given in (61).

**Lemma 5.6** *Then it holds for the Fisher-information matrix  $\mathbf{I}(\boldsymbol{\theta})$  that*

$$\mathbf{I}(\boldsymbol{\theta}) = \mathbf{C}(\boldsymbol{\theta})^\top \mathbf{C}(\boldsymbol{\theta}), \quad (62)$$

where

$$\mathbf{C}(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial p_1(\boldsymbol{\theta})/\partial \theta_1}{\sqrt{p_1(\boldsymbol{\theta})}} & \frac{\partial p_1(\boldsymbol{\theta})/\partial \theta_2}{\sqrt{p_1(\boldsymbol{\theta})}} & \cdots & \frac{\partial p_1(\boldsymbol{\theta})/\partial \theta_m}{\sqrt{p_1(\boldsymbol{\theta})}} \\ \frac{\partial p_2(\boldsymbol{\theta})/\partial \theta_1}{\sqrt{p_2(\boldsymbol{\theta})}} & \frac{\partial p_2(\boldsymbol{\theta})/\partial \theta_2}{\sqrt{p_2(\boldsymbol{\theta})}} & \cdots & \frac{\partial p_2(\boldsymbol{\theta})/\partial \theta_m}{\sqrt{p_2(\boldsymbol{\theta})}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial p_r(\boldsymbol{\theta})/\partial \theta_1}{\sqrt{p_r(\boldsymbol{\theta})}} & \frac{\partial p_r(\boldsymbol{\theta})/\partial \theta_2}{\sqrt{p_r(\boldsymbol{\theta})}} & \cdots & \frac{\partial p_r(\boldsymbol{\theta})/\partial \theta_m}{\sqrt{p_r(\boldsymbol{\theta})}} \end{pmatrix}. \quad (63)$$

**Proof**

- Because of (61) it holds for each  $x \in \mathbb{R}$  that

$$\log L(x; \boldsymbol{\theta}) = \sum_{j=1}^r \mathbb{1}_{\{a_j < x \leq b_j\}} \log p_j(\boldsymbol{\theta}).$$

- From this it follows for the entries  $I_{ij}(\boldsymbol{\theta})$  of  $\mathbf{I}(\boldsymbol{\theta})$  that

$$\begin{aligned} I_{ij}(\boldsymbol{\theta}) &= \mathbb{E}_{\boldsymbol{\theta}} \left( \frac{\partial}{\partial \theta_i} \log L(X_1; \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log L(X_1; \boldsymbol{\theta}) \right) = \sum_{k=1}^r \left( \frac{\partial}{\partial \theta_i} \log p_k(\boldsymbol{\theta}) \right) \left( \frac{\partial}{\partial \theta_j} \log p_k(\boldsymbol{\theta}) \right) p_k(\boldsymbol{\theta}) \\ &= \sum_{k=1}^r \left( \frac{\partial}{\partial \theta_i} p_k(\boldsymbol{\theta}) \right) (p_k(\boldsymbol{\theta}))^{-1} \left( \frac{\partial}{\partial \theta_j} p_k(\boldsymbol{\theta}) \right) = \left( \mathbf{C}(\boldsymbol{\theta})^\top \mathbf{C}(\boldsymbol{\theta}) \right)_{ij}. \quad \square \end{aligned}$$

Thus, Theorem 5.8 implies the following result.

**Corollary 5.1** *If the matrix  $\mathbf{I}(\boldsymbol{\theta}) = \mathbf{C}(\boldsymbol{\theta})^\top \mathbf{C}(\boldsymbol{\theta})$ , given in (63), is positive definite for each  $\boldsymbol{\theta} \in \Theta$ , then it holds that*

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}(X_1, \dots, X_n) - \boldsymbol{\theta}) \xrightarrow{d} \mathbf{N}(\mathbf{o}, (\mathbf{C}(\boldsymbol{\theta})^\top \mathbf{C}(\boldsymbol{\theta}))^{-1}) \quad (64)$$

for each weakly consistent sequence  $\{\widehat{\boldsymbol{\theta}}(X_1, \dots, X_n), n \geq 1\}$  of maximum-likelihood estimators for  $\boldsymbol{\theta}$  which are obtained from observations of the “coarsened” model.

**Remark**

- From (61) it follows for the likelihood function  $L(x_1, \dots, x_n; \boldsymbol{\theta})$  that

$$L(x_1, \dots, x_n; \boldsymbol{\theta}) = \prod_{j=1}^r p_j(\boldsymbol{\theta})^{Z_j(x_1, \dots, x_n)},$$

or for the loglikelihood function  $\log L(x_1, \dots, x_n; \boldsymbol{\theta})$  that

$$\log L(x_1, \dots, x_n; \boldsymbol{\theta}) = \sum_{j=1}^r Z_j(x_1, \dots, x_n) \log p_j(\boldsymbol{\theta}). \quad (65)$$

- Each maximum-likelihood estimate  $\widehat{\boldsymbol{\theta}} = \widehat{\boldsymbol{\theta}}(Z_1(x_1, \dots, x_n), \dots, Z_r(x_1, \dots, x_n))$  for  $\boldsymbol{\theta}$  which is obtained from the coarsened data  $Z_1(x_1, \dots, x_n), \dots, Z_r(x_1, \dots, x_n)$  satisfies the system of equations

$$\frac{\partial \log L(x_1, \dots, x_n; \boldsymbol{\theta})}{\partial \theta_i} = 0, \quad \forall i = 1, \dots, m. \quad (66)$$

because of the above-mentioned regularity conditions.

- Here it follows from (65) that for arbitrary  $i = 1, \dots, m$  and  $\boldsymbol{\theta} \in \Theta$

$$\frac{\partial \log L(x_1, \dots, x_n; \boldsymbol{\theta})}{\partial \theta_i} = \sum_{j=1}^r \frac{Z_j(x_1, \dots, x_n)}{p_j(\boldsymbol{\theta})} \frac{\partial p_j(\boldsymbol{\theta})}{\partial \theta_i}$$

or

$$\frac{\partial \log L(x_1, \dots, x_n; \boldsymbol{\theta})}{\partial \theta_i} = \sum_{j=1}^r \frac{Z_j(x_1, \dots, x_n) - np_j(\boldsymbol{\theta})}{p_j(\boldsymbol{\theta})} \frac{\partial p_j(\boldsymbol{\theta})}{\partial \theta_i}, \quad (67)$$

where the last equality is due to the fact that

$$\sum_{j=1}^r \frac{\partial p_j(\boldsymbol{\theta})}{\partial \theta_i} = 0, \quad \forall i = 1, \dots, m. \quad (68)$$

### 5.3.4 Asymptotic Distribution of the Pearson–Fisher Statistic

The following theorem is the basis for the  $\chi^2$ -goodness-of-fit test of Pearson–Fisher. Here we always assume that

- the likelihood function of the coarsened model, considered in (61), fulfills the regularity conditions of Section 5.3.2,
- the Fisher–information matrix  $I(\boldsymbol{\theta})$ , given in (62), is positive definite and  $\{\widehat{\boldsymbol{\theta}}_n\} = \{\widehat{\boldsymbol{\theta}}(X_1, \dots, X_n), n \geq 1\}$  is a weakly consistent sequence of ML estimators for  $\boldsymbol{\theta}$  which is obtained by considering the coarsened model.

#### Theorem 5.9

- Let  $\widehat{T}_n(X_1, \dots, X_n)$  be the Pearson–Fisher test statistic given in (53), i.e.,

$$\widehat{T}_n(X_1, \dots, X_n) = \sum_{j=1}^r \frac{(Z_j(X_1, \dots, X_n) - n\widehat{p}_j(X_1, \dots, X_n))^2}{n\widehat{p}_j(X_1, \dots, X_n)}, \quad (69)$$

where  $\widehat{p}_j(X_1, \dots, X_n) = p_j(\widehat{\boldsymbol{\theta}}(X_1, \dots, X_n))$ .

- Then it holds that

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\boldsymbol{\theta}}(\widehat{T}_n(X_1, \dots, X_n) > \chi_{r-1-m, 1-\alpha}^2) = \alpha, \quad \forall \alpha \in (0, 1) \quad (70)$$

for each  $\boldsymbol{\theta} \in \Theta$ , where  $\chi_{r-1-m, 1-\alpha}^2$  denotes the  $(1-\alpha)$ -quantile of the  $\chi^2$ -distribution with  $r-1-m$  degrees of freedom.

It is possible to give a full proof of Theorem 5.9 by reinterpreting the  $\chi^2$ -goodness-of-fit test of Pearson–Fisher as a *likelihood-ratio test*, cf. for instance Section 4.7 in H. Pruscha (2000) *Vorlesungen über mathematische Statistik*, Teubner-Verlag, Stuttgart.

However, as this method of proof is rather complex, we only show a derivation of Theorem 5.9, which is partly *heuristic*.

- Let  $\mathbf{p}(\boldsymbol{\theta}) = (p_1(\boldsymbol{\theta}), \dots, p_r(\boldsymbol{\theta}))^\top$  and  $\widetilde{\mathbf{Z}}_n(\boldsymbol{\theta}) = (\widetilde{Z}_{n1}(\boldsymbol{\theta}), \dots, \widetilde{Z}_{nr}(\boldsymbol{\theta}))^\top$  with

$$\widetilde{Z}_{nj}(\boldsymbol{\theta}) = \frac{Z_j(X_1, \dots, X_n) - np_j(\boldsymbol{\theta})}{\sqrt{np_j(\boldsymbol{\theta})}}, \quad j = 1, \dots, r. \quad (71)$$

- Since  $\mathbb{E}\widetilde{\mathbf{Z}}_n(\boldsymbol{\theta}) = \mathbf{o}$  and since it is possible to write  $\widetilde{\mathbf{Z}}_n(\boldsymbol{\theta})$  as a sum of  $n$  independent and identically distributed random vectors, it follows from the multivariate central limit theorem (in the same way as in the proof of Theorem 5.5) that

$$\widetilde{\mathbf{Z}}_n(\boldsymbol{\theta}) \xrightarrow{d} \widetilde{\mathbf{Z}}(\boldsymbol{\theta}) \sim N(\mathbf{o}, \mathbf{B}(\boldsymbol{\theta})\mathbf{K}(\boldsymbol{\theta})\mathbf{B}(\boldsymbol{\theta})), \quad (72)$$

where

$$\mathbf{B}(\boldsymbol{\theta}) = \begin{pmatrix} 1/\sqrt{p_1(\boldsymbol{\theta})} & 0 & \dots & 0 \\ 0 & 1/\sqrt{p_2(\boldsymbol{\theta})} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1/\sqrt{p_r(\boldsymbol{\theta})} \end{pmatrix}, \quad \mathbf{K}(\boldsymbol{\theta}) = \begin{cases} -p_i(\boldsymbol{\theta})p_j(\boldsymbol{\theta}), & \text{if } i \neq j, \\ p_i(\boldsymbol{\theta})(1-p_j(\boldsymbol{\theta})), & \text{if } i = j. \end{cases}$$

– Thus, for the covariance matrix  $\mathbf{B}(\boldsymbol{\theta})\mathbf{K}(\boldsymbol{\theta})\mathbf{B}(\boldsymbol{\theta})$  in (72) we get that

$$\mathbf{B}(\boldsymbol{\theta})\mathbf{K}(\boldsymbol{\theta})\mathbf{B}(\boldsymbol{\theta}) = \mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta}), \quad \text{where } \mathbf{q}(\boldsymbol{\theta}) = (\sqrt{p_1(\boldsymbol{\theta})}, \dots, \sqrt{p_r(\boldsymbol{\theta})})^\top. \quad (73)$$

– For the matrix  $\mathbf{C}(\boldsymbol{\theta})$ , introduced in (63), it holds because of (68) that  $\mathbf{q}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}) = \mathbf{o}$  and thus

$$\begin{aligned} & \left( (\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta}) \right) \left( \mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta}) \right) \left( (\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta}) \right)^\top \\ &= \left( (\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta}) \right) \left( \mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta}) \right) \mathbf{C}(\boldsymbol{\theta}) (\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1} \\ &= (\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}. \end{aligned}$$

– From (72) and (73) it now follows that

$$(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})\tilde{\mathbf{Z}}(\boldsymbol{\theta}) \sim \mathbf{N}(\mathbf{o}, (\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}). \quad (74)$$

• Furthermore, Corollary 5.1 and the Taylor series expansion yield

$$\sqrt{n} \mathbf{B}(\boldsymbol{\theta})(\mathbf{p}(\hat{\boldsymbol{\theta}}_n) - \mathbf{p}(\boldsymbol{\theta})) = \sqrt{n} \mathbf{C}(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) + o(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \xrightarrow{d} \mathbf{N}(\mathbf{o}, \mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})).$$

– This and (74) imply that

$$\sqrt{n} \mathbf{B}(\boldsymbol{\theta})(\mathbf{p}(\hat{\boldsymbol{\theta}}_n) - \mathbf{p}(\boldsymbol{\theta})) \xrightarrow{d} \mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})\tilde{\mathbf{Z}}(\boldsymbol{\theta}). \quad (75)$$

– Moreover, it follows from (69) and (71) that

$$\begin{aligned} \hat{T}_n(X_1, \dots, X_n) &= \sum_{j=1}^r (\tilde{Z}_{nj}(\hat{\boldsymbol{\theta}}_n))^2 \\ &= \sum_{j=1}^r \left( \tilde{Z}_{nj}(\boldsymbol{\theta}) + \tilde{Z}_{nj}(\boldsymbol{\theta}) \left( \sqrt{\frac{p_j(\boldsymbol{\theta})}{p_j(\hat{\boldsymbol{\theta}}_n)}} - 1 \right) - \frac{\sqrt{n}}{\sqrt{p_j(\hat{\boldsymbol{\theta}}_n)}} (p_j(\hat{\boldsymbol{\theta}}_n) - p_j(\boldsymbol{\theta})) \right)^2 \\ &= \sum_{j=1}^r \left( \tilde{Z}_{nj}(\boldsymbol{\theta}) - \frac{\sqrt{n}}{\sqrt{p_j(\hat{\boldsymbol{\theta}}_n)}} (p_j(\hat{\boldsymbol{\theta}}_n) - p_j(\boldsymbol{\theta})) + o(1) \right)^2, \end{aligned}$$

– where the last equality follows from the convergence

$$\tilde{Z}_{nj}(\boldsymbol{\theta}) \left( \sqrt{\frac{p_j(\boldsymbol{\theta})}{p_j(\hat{\boldsymbol{\theta}}_n)}} - 1 \right) \xrightarrow{P} 0, \quad \forall j = 1, \dots, r,$$

which is due to  $\hat{\boldsymbol{\theta}}_n \xrightarrow{P} \boldsymbol{\theta}$  and the continuous mapping theorem for random vectors (cf. Lemma 4.5).

• In other words: With the notation  $\varphi(z_1, \dots, z_r) = \sum_{j=1}^r z_j^2$  it holds that

$$\hat{T}_n(X_1, \dots, X_n) = \varphi\left(\tilde{\mathbf{Z}}_n(\boldsymbol{\theta}) - \sqrt{n} \mathbf{B}(\boldsymbol{\theta})(\mathbf{p}(\hat{\boldsymbol{\theta}}_n) - \mathbf{p}(\boldsymbol{\theta})) + o(1)\right). \quad (76)$$

– Together with (72) and (75), the asymptotic approximation formula (76) suggests the *conjecture* that for  $n \rightarrow \infty$

$$\hat{T}_n(X_1, \dots, X_n) \xrightarrow{d} \varphi\left(\left(\mathbf{I}_r - \mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})\right)\tilde{\mathbf{Z}}(\boldsymbol{\theta})\right). \quad (77)$$

– However, the convergence in (77) *does not directly* follow from (72), (75) and (76), but needs a *separate proof*, which is omitted here.

- We now show that

$$\varphi\left(\left(\mathbf{I}_r - \mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})\right)\tilde{\mathbf{Z}}(\boldsymbol{\theta})\right) \sim \chi_{r-1-m}^2. \quad (78)$$

- In (72) and (73) we have shown that

$$\tilde{\mathbf{Z}}(\boldsymbol{\theta}) \sim N(\mathbf{o}, \mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta})), \quad \text{where } \mathbf{q}(\boldsymbol{\theta}) = (\sqrt{p_1(\boldsymbol{\theta})}, \dots, \sqrt{p_r(\boldsymbol{\theta})})^\top.$$

- Furthermore, it follows from  $\mathbf{q}^\top(\boldsymbol{\theta})\mathbf{q}(\boldsymbol{\theta}) = 1$  that

$$\begin{aligned} \left(\mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta})\right)^2 &= \mathbf{I}_r - 2\mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta}) + \mathbf{q}(\boldsymbol{\theta})\underbrace{\mathbf{q}^\top(\boldsymbol{\theta})\mathbf{q}(\boldsymbol{\theta})}_{=1}\mathbf{q}^\top(\boldsymbol{\theta}) \\ &= \mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta}), \end{aligned}$$

- i.e., the covariance matrix

$$\mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta}) = \begin{pmatrix} 1 - p_1(\boldsymbol{\theta}) & -\sqrt{p_1(\boldsymbol{\theta})p_2(\boldsymbol{\theta})} & -\sqrt{p_1(\boldsymbol{\theta})p_3(\boldsymbol{\theta})} & \dots & -\sqrt{p_1(\boldsymbol{\theta})p_r(\boldsymbol{\theta})} \\ -\sqrt{p_1(\boldsymbol{\theta})p_2(\boldsymbol{\theta})} & 1 - p_2(\boldsymbol{\theta}) & -\sqrt{p_2(\boldsymbol{\theta})p_3(\boldsymbol{\theta})} & \dots & -\sqrt{p_2(\boldsymbol{\theta})p_r(\boldsymbol{\theta})} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -\sqrt{p_1(\boldsymbol{\theta})p_r(\boldsymbol{\theta})} & -\sqrt{p_2(\boldsymbol{\theta})p_r(\boldsymbol{\theta})} & -\sqrt{p_3(\boldsymbol{\theta})p_r(\boldsymbol{\theta})} & \dots & 1 - p_r(\boldsymbol{\theta}) \end{pmatrix}$$

of the random vector  $\tilde{\mathbf{Z}}(\boldsymbol{\theta})$  is symmetric and idempotent.

- This and the last part of Theorem 1.4 imply the representation formula

$$\tilde{\mathbf{Z}}(\boldsymbol{\theta}) \stackrel{d}{=} \left(\mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta})\right)\boldsymbol{\varepsilon}, \quad (79)$$

where  $\boldsymbol{\varepsilon} \sim N(\mathbf{o}, \mathbf{I}_r)$ .

- Moreover, also the matrix  $\mathbf{I}_r - \mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})$  is symmetric and idempotent and from  $\mathbf{q}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}) = \mathbf{o}$  it follows that

$$\left(\mathbf{I}_r - \mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})\right)\left(\mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta})\right) = \mathbf{I}_r - \mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta}) - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta}).$$

- This implies that the matrix  $\mathbf{R} = \left(\mathbf{I}_r - \mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})\right)\left(\mathbf{I}_r - \mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta})\right)$  is symmetric and idempotent as well.
- From (79) and Theorem 1.9 we now get that

$$\varphi\left(\left(\mathbf{I}_r - \mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})\right)\tilde{\mathbf{Z}}(\boldsymbol{\theta})\right) \stackrel{d}{=} \varphi(\mathbf{R}\boldsymbol{\varepsilon}) = \boldsymbol{\varepsilon}^\top \mathbf{R} \boldsymbol{\varepsilon} \sim \chi_{\text{rk}(\mathbf{R})}^2.$$

- Due to Lemma 1.3 it holds for the rank  $\text{rk}(\mathbf{R})$  of the symmetric and idempotent matrix  $\mathbf{R}$  that

$$\begin{aligned} \text{rk}(\mathbf{R}) &= \text{tr}(\mathbf{R}) \\ &= \text{tr}(\mathbf{I}_r) - \text{tr}\left(\mathbf{C}(\boldsymbol{\theta})(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta}))^{-1}\mathbf{C}^\top(\boldsymbol{\theta})\right) - \text{tr}\left(\mathbf{q}(\boldsymbol{\theta})\mathbf{q}^\top(\boldsymbol{\theta})\right) \\ &= \text{tr}(\mathbf{I}_r) - \text{tr}\left(\left(\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta})\right)^{-1}\mathbf{C}^\top(\boldsymbol{\theta})\mathbf{C}(\boldsymbol{\theta})\right) - \text{tr}\left(\mathbf{q}^\top(\boldsymbol{\theta})\mathbf{q}(\boldsymbol{\theta})\right) \\ &= r - m - 1. \end{aligned}$$

- This proves the validity of (78).

**Remark** For the practical usage of the  $\chi^2$ - goodness-of-fit test of Pearson-Fisher one can proceed in a similar way as described in Section 5.2.2 in order to verify the hypothesis  $H_0 : P \in \{P_\theta, \boldsymbol{\theta} \in \Theta\}$ .

- First of all, a ML estimation  $\widehat{\boldsymbol{\theta}}(x_1, \dots, x_n) = (\widehat{\theta}_1(x_1, \dots, x_n), \dots, \widehat{\theta}_m(x_1, \dots, x_n))^\top$  for  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_m)^\top$  is determined by solving the system of equations (66).
- Then the value of the test statistic  $T_n(x_1, \dots, x_n)$ , defined in (53), is calculated.
- For sufficiently large sample sizes  $n$  the hypothesis  $H_0 : P \in \{P_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta\}$  is rejected if

$$T_n(x_1, \dots, x_n) > \chi_{r-1-m, 1-\alpha}^2,$$

- where  $\chi_{r-1-m, 1-\alpha}^2$  denotes the  $(1 - \alpha)$ -quantile of the  $\chi^2$ -distribution with  $(r - 1 - m)$  degrees of freedom.

## 5.4 Examples

### 5.4.1 Pearson–Fisher Test for Poisson–Distribution

- By observing the (independent and identically distributed) sample variables  $X_1, \dots, X_n$  it shall be tested if the distribution  $P$  of  $X_i$  belongs to the family of Poisson–distributions.
  - So let  $\Theta = (0, \infty)$  with  $\boldsymbol{\theta} = \lambda$  and let  $\{P_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta\} = \{\text{Poi}(\lambda), \lambda > 0\}$  be the family of Poisson–distributions.
  - We consider the following  $r$  classes  $\{0\}, \{1\}, \dots, \{r-2\}$  and  $\{r-1, r, r+1, \dots\}$ , i.e.
 
$$(a_1, b_1] = (-\infty, 0], \quad (a_2, b_2] = (0, 1], \quad \dots \quad (a_{r-1}, b_{r-1}] = (r-3, r-2], \quad (a_r, b_r] = (r-2, \infty].$$
  - The probabilities  $p_j(\lambda) = \mathbb{P}_\lambda(a_j < X_1 \leq b_j)$  then are given by

$$p_j(\lambda) = \frac{\lambda^{j-1}}{(j-1)!} e^{-\lambda} \quad \forall j = 1, \dots, r-1 \quad \text{and} \quad p_r(\lambda) = \sum_{i=r}^{\infty} \frac{\lambda^{i-1}}{(i-1)!} e^{-\lambda}. \quad (80)$$

- According to (66), every maximum–likelihood estimate  $\widehat{\lambda}$  for  $\lambda$  which is obtained from grouped data satisfies the equation

$$\sum_{j=1}^r Z_j(x_1, \dots, x_n) \frac{d}{d\lambda} \frac{p_j(\lambda)}{p_j(\lambda)} = 0. \quad (81)$$

- Here it follows from (80) that

$$\frac{d}{d\lambda} \frac{p_j(\lambda)}{p_j(\lambda)} = \frac{j-1}{\lambda} - 1 \quad \forall j = 1, \dots, r-1 \quad \text{and} \quad \frac{d}{d\lambda} \frac{p_r(\lambda)}{p_r(\lambda)} = \frac{\sum_{i=r}^{\infty} \left(\frac{i-1}{\lambda} - 1\right) \lambda^{i-1}}{\sum_{i=r}^{\infty} \lambda^{i-1}}.$$

- This and (81) imply that the ML estimate  $\widehat{\lambda}$  fulfills the following equation:

$$\sum_{j=1}^{r-1} Z_j(x_1, \dots, x_n) \left(\frac{j-1}{\lambda} - 1\right) + Z_r(x_1, \dots, x_n) \frac{\sum_{i=r}^{\infty} \left(\frac{i-1}{\lambda} - 1\right) \lambda^{i-1}}{\sum_{i=r}^{\infty} \lambda^{i-1}} = 0. \quad (82)$$

- For each  $n$  there is one  $r_0 = r_0(n) \in \mathbb{N}$ , such that  $Z_r(x_1, \dots, x_n) = 0$  for each  $r > r_0$ . This and (82) imply that for  $r \rightarrow \infty$

$$\widehat{\lambda}_n = \widehat{\lambda}(x_1, \dots, x_n) \rightarrow \bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i.$$

- For a sufficiently large number  $r$  of classes  $\{0\}, \{1\}, \dots, \{r-2\}, \{r-1, r, r+1, \dots\}$  the sample mean  $\bar{x}_n$ , which is a ML estimate for  $\lambda$  in the non-aggregated Poisson-model, is a good approximation for the ML estimation  $\widehat{\lambda}_n$  for  $\lambda$  in aggregated Poisson-models.
- The null hypothesis  $H_0 : P \in \{\text{Poi}(\lambda), \lambda > 0\}$  is therefore rejected if

$$\widehat{T}_n(x_1, \dots, x_n) = \sum_{j=1}^r \frac{(Z_j(x_1, \dots, x_n) - n\widehat{p}_j(x_1, \dots, x_n))^2}{n\widehat{p}_j(x_1, \dots, x_n)} > \chi_{r-2, 1-\alpha}^2,$$

where  $\widehat{p}_j(x_1, \dots, x_n) = p_j(\bar{x}_n)$  with the function  $p_j : (0, \infty) \rightarrow [0, 1]$  given in (80) and the estimation  $\bar{x}_n$  for  $\lambda$ .

#### 5.4.2 Pearson–Fisher Test for Normal Distribution

- Now let  $\Theta = \mathbb{R} \times (0, \infty)$  with  $\boldsymbol{\theta} = (\mu, \sigma^2)^\top$  and let  $\{P_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta\} = \{\text{N}(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$  be the family of (one-dimensional) normal distributions.

- The probabilities  $p_j(\boldsymbol{\theta}) = \mathbb{P}_{\boldsymbol{\theta}}(a_j < X_1 \leq b_j)$  are then given by

$$p_j(\boldsymbol{\theta}) = \int_{a_j}^{b_j} f(x; \boldsymbol{\theta}) dx, \quad \text{where } f(x; \boldsymbol{\theta}) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (83)$$

- According to (66), each maximum-likelihood estimate

$$\widehat{\boldsymbol{\theta}}(x_1, \dots, x_n) = (\widehat{\mu}(x_1, \dots, x_n), \widehat{\sigma}^2(x_1, \dots, x_n))^\top$$

for  $\boldsymbol{\theta} = (\mu, \sigma^2)^\top$  which is obtained from the aggregated data satisfies the system of equations

$$\sum_{j=1}^r Z_j(x_1, \dots, x_n) \frac{\int_{a_j}^{b_j} \frac{\partial}{\partial \theta_i} f(x; \boldsymbol{\theta}) dx}{\int_{a_j}^{b_j} f(x; \boldsymbol{\theta}) dx} = 0 \quad \text{for } i = 1, 2. \quad (84)$$

- Here it follows from (83) that

$$\frac{\partial}{\partial \mu} f(x; \boldsymbol{\theta}) = \frac{x-\mu}{\sigma^2} f(x; \boldsymbol{\theta}) \quad \text{or} \quad \frac{\partial}{\partial \sigma^2} f(x; \boldsymbol{\theta}) = f(x; \boldsymbol{\theta}) \left( \frac{(x-\mu)^2}{2\sigma^4} - \frac{1}{2\sigma^2} \right).$$

- This and (84) imply that the ML estimate  $\widehat{\boldsymbol{\theta}}(x_1, \dots, x_n)$  satisfies the following system of equations:

$$\sum_{j=1}^r Z_j(x_1, \dots, x_n) \frac{\int_{a_j}^{b_j} (x-\mu) f(x; \boldsymbol{\theta}) dx}{\int_{a_j}^{b_j} f(x; \boldsymbol{\theta}) dx} = 0, \quad \sum_{j=1}^r Z_j(x_1, \dots, x_n) \frac{\int_{a_j}^{b_j} (x-\mu)^2 f(x; \boldsymbol{\theta}) dx}{\int_{a_j}^{b_j} f(x; \boldsymbol{\theta}) dx} - n\sigma^2 = 0,$$

where the first equality of this system of equations is equivalent to

$$\sum_{j=1}^r Z_j(x_1, \dots, x_n) \frac{\int_{a_j}^{b_j} x f(x; \boldsymbol{\theta}) dx}{\int_{a_j}^{b_j} f(x; \boldsymbol{\theta}) dx} - \underbrace{\mu \sum_{j=1}^r Z_j(x_1, \dots, x_n)}_{=n} = 0.$$



- Thus, the ML estimate  $\widehat{\boldsymbol{\theta}}(x_1, \dots, x_n) = (\widehat{\mu}(x_1, \dots, x_n), \widehat{\sigma}^2(x_1, \dots, x_n))^\top$  fulfills the system of equations

$$\mu = \frac{1}{n} \sum_{j=1}^r Z_j(x_1, \dots, x_n) \frac{\int_{a_j}^{b_j} x f(x; \mu, \sigma^2) dx}{\int_{a_j}^{b_j} f(x; \mu, \sigma^2) dx}, \quad \sigma^2 = \frac{1}{n} \sum_{j=1}^r Z_j(x_1, \dots, x_n) \frac{\int_{a_j}^{b_j} (x - \mu)^2 f(x; \mu, \sigma^2) dx}{\int_{a_j}^{b_j} f(x; \mu, \sigma^2) dx},$$

- which for a sufficiently large number  $r$  of classes  $(a_1, b_1], \dots, (a_r, b_r]$  can be solved approximately in the following way:

$$\widehat{\mu} \approx \frac{1}{n} \sum_{j=1}^r c_j Z_j(x_1, \dots, x_n), \quad \widehat{\sigma}^2 \approx \frac{1}{n} \sum_{j=1}^r (c_j - \widehat{\mu})^2 Z_j(x_1, \dots, x_n), \quad (85)$$

- where  $c_1 = b_1$  is the right endpoint of the first class,  $c_r = b_{r-1}$  is the left endpoint of the  $r$ -th class and  $c_j = (b_{j-1} + b_j)/2$  is the center of the  $j$ -th class for  $j = 2, \dots, r-1$ .
- The null hypothesis  $H_0 : P \in \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$  is rejected if

$$\widehat{T}_n(x_1, \dots, x_n) = \sum_{j=1}^r \frac{(Z_j(x_1, \dots, x_n) - n\widehat{p}_j(x_1, \dots, x_n))^2}{n\widehat{p}_j(x_1, \dots, x_n)} > \chi_{r-3, 1-\alpha}^2,$$

where  $\widehat{p}_j(x_1, \dots, x_n) = p_j(\widehat{\mu}, \widehat{\sigma}^2)$  with the function  $p_j : \mathbb{R} \times (0, \infty) \rightarrow [0, 1]$  given in (83) and  $(\widehat{\mu}, \widehat{\sigma}^2)$  is the estimation for  $(\mu, \sigma^2)$ , given in (85).

### Remark

- The approximate solution (85) of the system of equations (84) shall now be used if the number  $r$  of classes is large enough.
  - This requires a sufficiently large sample size  $n$ .
  - In other words: If the sample size  $n$  is small, then the  $\chi^2$ -goodness-of-fit test is not suited to verify the hypothesis for normality.
- Alternative tests for normal distribution are the following *goodness-of-fit tests of Shapiro–Wilk-type*, which lead to acceptable results even for a small sample size  $n$ .

### 5.4.3 Goodness-of-Fit Tests of Shapiro–Wilk-Type

- In this section we discuss two goodness-of-fit tests of Shapiro–Wilk-type which can also be used to verify the hypothesis  $H_0 : P \in \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$ .
- Here the *order statistics*  $X_{(1)}, \dots, X_{(n)}$  of the (independent and identically distributed) sample variables  $X_1, \dots, X_n$  are considered, which have already been introduced in Section I-1.4.
  - *Recall:* The order statistics are defined by using the sample function  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , where

$$(x_1, \dots, x_n) \rightarrow (x_{(1)}, \dots, x_{(n)}) = \varphi(x_1, \dots, x_n) \quad \text{with } x_{(i)} = \min\{x_j : \#\{k : x_k \leq x_j\} \geq i\} \quad (86)$$

for each  $i \in \{1, \dots, n\}$ .

- The mapping  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , given in (86), is a *permutation* of the components of the vector  $(x_1, \dots, x_n)$ , such that

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}.$$

– For each  $\omega \in \Omega$  let now

$$(X_{(1)}(\omega), \dots, X_{(n)}(\omega)) = \varphi(X_1(\omega), \dots, X_n(\omega))$$

be the (measurable) permutation, given in (86), of  $(X_1(\omega), \dots, X_n(\omega))$ , such that

$$X_{(1)}(\omega) \leq \dots \leq X_{(n)}(\omega). \quad (87)$$

– The random variables  $X_{(1)}, \dots, X_{(n)} : \Omega \rightarrow \mathbb{R}$  are called the *order statistics* of  $(X_1, \dots, X_n)$ .

- If  $X_i \sim N(\mu, \sigma^2)$  for certain  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ , then it is easy to see that the following representation formula holds for the expectation  $b_i = \mathbb{E} X_{(i)}$  of the order statistics  $X_{(i)}$ :

$$b_i = \mu + \sigma a_i \quad \forall i = 1, \dots, n, \quad (88)$$

– where  $a_i = \mathbb{E} Y_{(i)}$  is the expectation of the  $i$ -th order statistic  $Y_{(i)}$  for  $N(0, 1)$ -distributed sample variables  $Y_1, \dots, Y_n$ .

– The benefit of the representation formula (88) is that the expectations  $a_1, \dots, a_n$  are available in the form of tables or can be determined using Monte Carlo simulation.

- Since the vectors  $(b_1, \dots, b_n)$  and  $(X_{(1)}, \dots, X_{(n)})$  should differ only little under  $H_0$ , the following *empirical correlation coefficient* is considered to verify the null hypothesis  $H_0 : P \in \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}$ :

$$\tilde{T}(X_1, \dots, X_n) = \frac{\sum_{i=1}^n (b_i - \bar{b})(X_{(i)} - \bar{X})}{\sqrt{\sum_{i=1}^n (b_i - \bar{b})^2} \sqrt{\sum_{i=1}^n (X_{(i)} - \bar{X})^2}}, \quad (89)$$

where  $\bar{b} = \sum_{i=1}^n b_i/n$  and  $\bar{X} = \sum_{i=1}^n X_i/n$ .

### 1. Shapiro–Francia test

- Since correlation coefficients are invariant under linear transformations, we are able to replace  $b_i$  in (89) by  $a_i$  for each  $i \in \{1, \dots, n\}$ , where  $\bar{a} = \sum_{i=1}^n a_i/n = 0$  holds.
- Furthermore, it holds that

$$\sum_{i=1}^n (X_{(i)} - \bar{X})^2 = \sum_{i=1}^n (X_i - \bar{X})^2 \quad \text{and} \quad \sum_{i=1}^n a_i \bar{X} = 0,$$

i.e., the definition of  $\tilde{T}(X_1, \dots, X_n)$  in (89) is equivalent to

$$\tilde{T}(X_1, \dots, X_n) = \frac{\sum_{i=1}^n a_i X_{(i)}}{\sqrt{\sum_{i=1}^n a_i^2} \sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}}. \quad (90)$$

- Since it always holds that  $|\tilde{T}(X_1, \dots, X_n)| \leq 1$ , the null hypothesis  $H_0$  is rejected if  $\tilde{T}(X_1, \dots, X_n) < q_{n,\alpha}$ , where  $q_{n,\alpha}$  denotes the  $\alpha$ -quantile of the distribution of  $\tilde{T}(X_1, \dots, X_n)$ .
- This is the so-called *Shapiro–Francia test* for normal distribution, where the quantiles  $q_{n,\alpha}$  of the distribution of  $\tilde{T}(X_1, \dots, X_n)$  can either be taken from a table or be determined using Monte Carlo simulation.

2. *Shapiro-Wilk test*

- In (90) it is possible to consider the linear transformation

$$(a'_1, \dots, a'_n)^\top = \mathbf{K}^{-1}(a_1, \dots, a_n)^\top$$

instead of  $a_1, \dots, a_n$ , where the covariance matrix  $\mathbf{K} = (k_{ij})$  is given by

$$k_{ij} = \mathbb{E}((Y_{(i)} - a_i)(Y_{(j)} - a_j)) \quad \text{with } Y_i \sim N(0, 1).$$

- The test constructed in this way is called *Shapiro-Wilk test*.

## 6 Nonparametric Localization Tests

### 6.1 Two Simple Examples of One-Sample Problems

#### 6.1.1 Binomial Test

- The  $\chi^2$  goodness-of-fit test considered in Section 5.2 can be replaced by the following *binomial test* if  $r = 2$ , i.e., if only two classes are considered (for example when dealing with binary data).

– Then we partition the domain of the (independent and identically distributed) sampling variables  $X_1, \dots, X_n$  into two subsets  $(a_1, b_1]$  and  $(a_2, b_2]$ , such that

$$(a_1, b_1] \cap (a_2, b_2] = \emptyset \quad \text{and} \quad \mathbb{P}(X_1 \in (a_1, b_1] \cup (a_2, b_2]) = 1,$$

and consider the “class size”

$$T(X_1, \dots, X_n) = \#\{i : 1 \leq i \leq n, a_1 < X_i \leq b_1\}.$$

– One can easily see that  $T = T(X_1, \dots, X_n)$  is binomial distributed, i.e.,

$$T \sim \text{Bin}(n, p), \quad \text{where } p = \mathbb{P}(a_1 < X_1 \leq b_1). \quad (1)$$

- To begin with, we consider the problem of testing  $H_0 : p = p_0$  versus  $H_1 : p \neq p_0$ , where  $p_0 \in (0, 1)$  is an arbitrary fixed number.

– Due to (1),  $H_0$  is rejected if  $T \leq t_{\alpha_1}$  or  $T \geq t_{1-\alpha_2}$ ,

– where the “critical values”  $t_{\alpha_1}$  and  $t_{1-\alpha_2}$  for arbitrary  $\alpha_1, \alpha_2 \in (0, 1)$  with  $\alpha_1 + \alpha_2 = \alpha$  are given by

$$\begin{aligned} t_{\alpha_1} &= \max\{t \in \mathbb{R} : \mathbb{P}_{p_0}(T \leq t) \leq \alpha_1\} \\ &= \max\left\{k \in \{0, 1, \dots, n\} : \sum_{i=0}^k \binom{n}{i} p_0^i (1-p_0)^{n-i} \leq \alpha_1\right\} \end{aligned}$$

and

$$\begin{aligned} t_{1-\alpha_2} &= \min\{t \in \mathbb{R} : \mathbb{P}_{p_0}(T \geq t) \leq \alpha_2\} \\ &= \min\left\{k \in \{0, 1, \dots, n\} : \sum_{i=k}^n \binom{n}{i} p_0^i (1-p_0)^{n-i} \leq \alpha_2\right\}. \end{aligned}$$

– For  $p_0 = 0.5$  one usually chooses  $\alpha_1 = \alpha_2 = \alpha/2$ . If  $p_0$  is close to 0 or 1, it is advisable to choose  $\alpha_1$  smaller or greater than  $\alpha_2$ , respectively.

– The quantiles  $t_{\alpha_1}$  and  $t_{1-\alpha_2}$  of the binomial distribution  $\text{Bin}(n, p_0)$  can either be taken from tables or be determined using Monte Carlo simulation.

- The (one-sided) problem of testing  $H_0 : p \leq p_0$  versus  $H_1 : p > p_0$  can be treated in a similar way. Here  $H_0$  is rejected if  $T \geq t_{1-\alpha}$ .

- In an analogous way, one obtains a decision rule for the (one-sided) problem of testing  $H_0 : p \geq p_0$  versus  $H_1 : p < p_0$ , where  $H_0$  is rejected if  $T \leq t_\alpha$ .

#### Remark

- The binomial test described above is also called *sign test* because the generation of 2 classes can be perceived as binarization of the original data.
- In the two one-sided problems of testing, the critical values  $t_{1-\alpha}$  and  $t_\alpha$  are determined for  $p = p_0$  although the null hypothesis is  $H_0 : p \leq p_0$  or  $H_0 : p \geq p_0$ , respectively.

- Considering the values  $t_{1-\alpha}$  and  $t_\alpha$  anyway does here not contradict the fact that for each  $p < p_0$  or  $p > p_0$  the critical value would be smaller than  $t_{1-\alpha}$  or greater than  $t_\alpha$  and that  $H_0$  would thus have to be rejected more often.
- The choice of the critical values  $t_{1-\alpha}$  and  $t_\alpha$  can be explained by the fact that one does not consider a single  $p$  with  $p < p_0$  or  $p > p_0$ , but that  $p$  can be arbitrarily close to  $p_0$  and hence, in particular, also  $p = p_0$  is allowed.
- If the sample size  $n$  is large and if  $p_0$  is close to 0 or 1,
  - the direct computation of the quantiles  $t_{1-\alpha}$  and  $t_\alpha$  of the binomial distribution  $\text{Bin}(n, p_0)$  is difficult.
  - The law of rare events (cf. Section WR-3.2.2) implies that  $t_{1-\alpha}$  and  $t_\alpha$  can be approximated by quantiles of the Poisson distribution  $\text{Poi}(\lambda)$  in this case, where  $\lambda = np_0$  or  $\lambda = n(1 - p_0)$ , respectively.
- Moreover, for arbitrary fixed  $p_0 \in (0, 1)$  the critical values  $t_{1-\alpha}$  and  $t_\alpha$  can be approximated by suitably transformed quantiles of the normal distribution  $N(0, 1)$  if the sample size  $n$  is „sufficiently large”.
  - In this case, it follows from the central limit theorem of DeMoivre–Laplace (cf. Theorem WR-3.6) that the transformed test statistic

$$T' = \frac{T - np_0}{\sqrt{np_0(1 - p_0)}}$$

is approximately  $N(0, 1)$ -distributed, i.e., that

$$\mathbb{P}(T \leq t) = \mathbb{P}(T' \leq t') \approx \Phi(t'), \quad \text{where } t' = \frac{t - np_0}{\sqrt{np_0(1 - p_0)}}$$

and  $\Phi: \mathbb{R} \rightarrow [0, 1]$  is the distribution function of the  $N(0, 1)$ -distribution.

- Therefore, one gets that  $t_\alpha \approx np_0 + z_\alpha \sqrt{np_0(1 - p_0)}$ , where  $z_\alpha$  is the  $\alpha$ -quantile of the  $N(0, 1)$ -distribution.
- Possible criteria for “sufficiently large” which are mentioned in literature are, e.g., the conditions  $n \geq 20$  and  $10 \leq np_0 \leq n - 10$ .
  - When investigating the (two-sided) problem of testing  $H_0: p = p_0$  versus  $H_1: p \neq p_0$ , then  $H_0$  is rejected if

$$T \leq np_0 + z_{\alpha_1} \sqrt{np_0(1 - p_0)} \quad \text{or} \quad T \geq np_0 + z_{1-\alpha_2} \sqrt{np_0(1 - p_0)}.$$

- Similar approximation formulas can be derived for the critical values of the one-sided tests mentioned above.

### Example

- Let the distribution function  $F: \mathbb{R} \rightarrow [0, 1]$  of the sampling variable  $X_1, \dots, X_n$  be continuous and let  $\gamma_p$  be the  $p$ -quantile of  $F$ , i.e., let  $F(\gamma_p) = p$  for  $p \in (0, 1)$ .
- In order to verify the hypothesis  $H_0: \gamma_p = \gamma_p^0$ , one can consider the “coarsened” random sample  $(Y_1, \dots, Y_n)$  with

$$Y_i = \begin{cases} 1, & \text{if } X_i \leq \gamma_p, \\ 0, & \text{if } X_i > \gamma_p. \end{cases}$$

- Then  $Y_i \sim \text{Bin}(1, p)$  for each  $i = 1, \dots, n$  and the hypothesis  $H_0: \gamma_p = \gamma_p^0$  with respect to  $(X_1, \dots, X_n)$  is equivalent to the hypothesis  $H_0: p = p_0$  with respect to  $(Y_1, \dots, Y_n)$ .
- Hence, the binomial test can in particular be used to verify the hypothesis  $H_0: \gamma_{0.5} = 0$ .

### 6.1.2 Run Test for Randomness

- In this section it is *not* assumed that the sampling variables  $X_1, \dots, X_n$  are independent.
  - We merely assume that  $X_1, \dots, X_n$  can only take the values 0 or 1, where the value 0 shall occur  $n_1$  times and the value 1 shall occur  $n_2$  times;  $n_2 = n - n_1$ .
  - Thus, there are altogether  $\binom{n}{n_1}$  possible realizations of the random sample  $(X_1, \dots, X_n)$ .
  - We now want to verify the null hypothesis  $H_0$  that each of these  $\binom{n}{n_1}$  realizations occurs with the same probability.
  - In other words: We want to check whether the localization, i.e., the order in which the  $n_1$  ones and  $n_2$  zeros are arranged, is “purely at random”.
- As a test statistic  $T : \Omega \rightarrow \{0, 1, \dots\}$  we consider the number  $T(\omega)$  of *runs* in the (concrete) sample  $\omega = (x_1, \dots, x_n)$ , i.e., the number of (sub-) sequences of consecutive equal symbols  $\omega = (x_1, \dots, x_n)$ .

#### Example

- Let  $n = 20$  with  $n_1 = 12$  and  $n_2 = 8$ . For

$$\omega = (1, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1) \quad (2)$$

one then gets that  $T(\omega) = 7$ .

- We now investigate the question whether the data given in (2) is compatible with the hypothesis  $H_0$  that the order is “purely at random” or whether  $H_0$  should be rejected.
- For this purpose, we specify the distribution of  $T$  by considering a suitably chosen (Laplace) probability space, cf. Section WR–2.4.1.

**Theorem 6.1** *Assuming that  $H_0$  is true, it holds for each  $i = 1, 2, \dots, \min\{n_1, n_2\}$  that*

$$\mathbb{P}(T = k) = \begin{cases} \frac{2 \binom{n_1 - 1}{i - 1} \binom{n_2 - 1}{i - 1}}{\binom{n}{n_1}}, & \text{wenn } k = 2i, \\ \frac{\binom{n_1 - 1}{i} \binom{n_2 - 1}{i - 1} + \binom{n_1 - 1}{i - 1} \binom{n_2 - 1}{i}}{\binom{n}{n_1}}, & \text{wenn } k = 2i + 1. \end{cases} \quad (3)$$

Furthermore, it holds that

$$\mathbb{E}T = 1 + \frac{2n_1n_2}{n} \quad \text{and} \quad \text{Var}T = \frac{2n_1n_2(2n_1n_2 - n)}{n^2(n - 1)}. \quad (4)$$

#### Proof

- We only prove (3) for the case  $k = 2i$  since the proof for the case  $k = 2i + 1$  proceeds analogously.
  - Hence, let  $k = 2i$ . Then there are at a time  $i$  runs consisting of ones and zeros, respectively.
  - For the decomposition of the  $n_1$  zeros into  $i$  subsets, there are  $\binom{n_1 - 1}{i - 1}$  possibilities.
  - For each of these decompositions, there are  $\binom{n_2 - 1}{i - 1}$  possibilities to divide the  $n_2$  ones into  $i$  subsets.

- If we now additionally assume that the sample  $\omega = (x_1, \dots, x_n)$  can begin either with  $x_1 = 0$  or with  $x_1 = 1$ , we obtain a total of  $2 \binom{n_1-1}{i-1} \binom{n_2-1}{i-1}$  decomposition possibilities.
- Therefore, (3) is proved for the case  $k = 2i$ .
- In order to determine the expectation  $\mathbb{E}T$ , we use the following consideration.
  - For each  $j = 2, \dots, n$  we consider the indicator variable  $Y_j : \Omega \rightarrow \{0, 1\}$  with

$$Y_j = \begin{cases} 1, & \text{if a run starts at the } j\text{-th position,} \\ 0, & \text{else.} \end{cases}$$

- Then it holds that  $\{\omega \in \Omega : Y_j(\omega) = 1\} = \{\omega \in \Omega : X_{j-1}(\omega) \neq X_j(\omega)\}$ , i.e., there are  $2 \binom{n-2}{n_i-1}$  possibilities that a run starts at the  $j$ -th position.
- Therefore, one gets

$$\mathbb{E}Y_j = \mathbb{P}(Y_j = 1) = 2 \frac{\binom{n-2}{n_i-1}}{\binom{n}{n_1}} = 2 \frac{(n-2)!(n-n_1)!n_1!}{(n-n_1-1)!(n_1-1)!n!} = 2 \frac{n_1(n-n_1)}{n(n-1)}.$$

- Together with the identity

$$T = 1 + \sum_{j=2}^n Y_j \tag{5}$$

this implies that

$$\mathbb{E}T = 1 + \sum_{j=2}^n \mathbb{E}Y_j = 1 + 2 \frac{n_1(n-n_1)}{n}.$$

- The variance formula in (4) can be proved in a similar way because (5) implies that

$$\begin{aligned} \text{Var } T &= \mathbb{E} \left( \sum_{j=2}^n \mathbb{E}Y_j \right)^2 - \left( \sum_{j=2}^n \mathbb{E}Y_j \right)^2 \\ &= \sum_{j=2}^n \mathbb{E}Y_j^2 + \sum_{2 \leq j_1, j_2 \leq n, j_1 \neq j_2} \mathbb{E}(Y_{j_1} Y_{j_2}) - \left( \sum_{j=2}^n \mathbb{E}Y_j \right)^2 \\ &= \sum_{j=2}^n \mathbb{E}Y_j + \sum_{2 \leq j_1, j_2 \leq n, j_1 \neq j_2} \mathbb{E}(Y_{j_1} Y_{j_2}) - \left( \sum_{j=2}^n \mathbb{E}Y_j \right)^2, \end{aligned}$$

hence one only needs to specify the moments  $\mathbb{E}(Y_{j_1} Y_{j_2})$ . □

### Remark

- A possible alternative to the null hypothesis  $H_0$  that the localization of the zeros and ones is “purely at random” is their trend to form clumps or clusters.
- As rejection region of  $H_0$  one then chooses the left-hand end of the distribution of  $T$ .
- In other words:  $H_0$  is rejected if  $T \leq r_\alpha(n_1; n_2)$ , where

$$r_\alpha(n_1; n_2) = \max \left\{ r \in \{1, 2, \dots\} : \mathbb{P}(T \leq r) \leq \alpha \right\}$$

is the  $\alpha$ -quantile of the distribution of the test statistic  $T$ .

- The quantiles  $r_\alpha(n_1; n_2)$  can be computed using the formulas for the probabilities  $\mathbb{P}(T = k)$  given in Theorem 6.1. They can be taken from tables in literature.

**Example** (continued) For  $\alpha = 0.1$  and  $n_1 = 12, n_2 = 8$  one obtains that  $r_{0.1}(12; 8) = 7$ . Moreover, it holds for the sample considered in (2) that

$$T(\omega) = 7 \left( \leq r_{0.1}(12; 8) \right),$$

i.e.,  $H_0$  is rejected.

If the (sub-) sample sizes  $n_1$  and  $n_2$  are large, the determination of the quantiles  $r_\alpha(n_1; n_2)$  of  $T = T_{n_1, n_2}$  involves a considerable computational cost. A way out is offered by the following central limit theorem, which we state without proof.

**Theorem 6.2** *If  $n_1, n_2 \rightarrow \infty$  such that  $n_1/(n_1 + n_2) \rightarrow p$  or equivalently  $n_2/(n_1 + n_2) \rightarrow 1 - p$  for a  $p \in (0, 1)$ , then it holds that*

$$\lim_{n_1, n_2 \rightarrow \infty} \frac{1}{n_1 + n_2} \mathbb{E} T_{n_1, n_2} = 2p(1 - p) \quad \text{and} \quad \lim_{n_1, n_2 \rightarrow \infty} \frac{1}{n_1 + n_2} \text{Var} T_{n_1, n_2} = 4p^2(1 - p)^2 \quad (6)$$

as well as

$$\lim_{n_1, n_2 \rightarrow \infty} \mathbb{P} \left( \frac{T_{n_1, n_2} - 2(n_1 + n_2)p(1 - p)}{2\sqrt{n_1 + n_2}p(1 - p)} \leq x \right) = \Phi(x) \quad \forall x \in \mathbb{R}, \quad (7)$$

where  $\Phi : \mathbb{R} \rightarrow [0, 1]$  is the distribution function of the  $N(0, 1)$ -distribution.

**Remark** Theorem 6.2 implies that for large  $n_1, n_2$  the null hypothesis  $H_0$  is rejected if

$$\frac{T_{n_1, n_2} - 2n_1n_2/(n_1 + n_2)}{2n_1n_2/(n_1 + n_2)^{3/2}} \leq z_\alpha, \quad (8)$$

where  $z_\alpha$  is the  $\alpha$ -quantile of the  $N(0, 1)$ -distribution.

## 6.2 Wilcoxon–Rank Test

### 6.2.1 Model Description; Median Test

- We now return to the case that the sampling variables  $X_1, \dots, X_n$  are independent and identically distributed with the distribution function  $F : \mathbb{R} \rightarrow [0, 1]$ .
  - At the end of Section 6.1.1, in the context of the binomial or sign test, we have discussed a *median test* for verifying the hypothesis
 
$$H_0 : \gamma_{0.5} = 0, \quad (9)$$
 where  $\gamma_{0.5}$  is a median of  $F$ , i.e.,  $F(\gamma_{0.5}) = 0.5$ .
    - In this section we consider another (more efficient) approach for testing the hypothesis given in (9).
- In doing so, we assume that the distribution function  $F$  of the sampling variables  $X_1, \dots, X_n$  belongs to the following (nonparametric) class of distribution functions.
  - Let  $G : \mathbb{R} \rightarrow [0, 1]$  be an arbitrary continuous distribution function with the following kind of symmetry with respect to the origin: For each  $x \in \mathbb{R}$  it holds that  $G(-x) = 1 - G(x)$ .
  - This implies in particular that  $G(0) = 1/2$ , i.e., zero is a median of  $G$ .
  - Let the family  $\Delta$  of distribution functions of the sampling variables  $X_1, \dots, X_n$  which is taken into account in the (two-sided) Wilcoxon test be given by  $\Delta = \{F_\delta : F_\delta(x) = G(x - \delta) \forall x, \delta \in \mathbb{R}\}$ .



– Since  $G$  is continuous, one then gets that

$$\mathbb{P}(X_i = x) = \mathbb{P}(X_i = X_j) = 0 \quad (10)$$

for each  $x \in \mathbb{R}$  and for arbitrary  $i, j = 1, \dots, n$  with  $i \neq j$ .

- We discuss the (two-sided) problem of testing  $H_0 : \delta = \delta_0$  vs.  $H_1 : \delta \neq \delta_0$  for some  $\delta_0 \in \mathbb{R}$ .
  - Here we can (w.l.o.g.) set  $\delta_0 = 0$ ; otherwise, the transformed sampling variables  $X'_1, \dots, X'_n$  with  $X'_i = X_i - \delta_0$  can be considered.
  - In a similar way, also the (one-sided) problem of testing  $H_0 : \delta = 0$  vs.  $H_1 : \delta > 0$  can be treated.
- For the verification of the null hypothesis  $H_0 : \delta = 0$  we consider the *ranks*  $R_1, \dots, R_n$  of the random variables  $|X_1|, \dots, |X_n|$  with

$$R_i = \sum_{j=1}^n \mathbb{I}_{\{|X_j| \leq |X_i|\}} \quad \forall i = 1, \dots, n,$$

where the indicator variable  $\mathbb{I}_{\{|X_j| \leq |X_i|\}} : \Omega \rightarrow \{0, 1\}$  is given by

$$\mathbb{I}_{\{|X_j| \leq |X_i|\}}(\omega) = \begin{cases} 1, & \text{if } |X_j(\omega)| \leq |X_i(\omega)|, \\ 0, & \text{else.} \end{cases}$$

- Then we consider the test statistics

$$T_n^+ = \sum_{i=1}^n R_i \mathbb{I}_{\{X_i > 0\}} \quad \text{and} \quad T_n^- = \sum_{i=1}^n R_i \mathbb{I}_{\{X_i < 0\}}. \quad (11)$$

### Remark

- Due to (10) it holds with probability 1 that

$$T_n^- = \sum_{i=1}^n R_i - T_n^+ = \binom{n+1}{2} - T_n^+. \quad (12)$$

- One can show that, assuming that  $H_0 : \delta = 0$  is true, it holds that  $T_n^- \stackrel{d}{=} T_n^+$ ; cf. (18).
- Thus, in the case that  $H_0 : \delta = 0$  is true, the test statistics  $T_n^+$  and  $T_n^-$  should take values that are approximately equal. Because of (12), this means that  $T_n^+ \approx \binom{n+1}{2}/2$ .
- Very small or very large values of  $T_n^+$  hence indicate that the alternative hypothesis  $H_1 : \delta \neq \delta_0$  might be true, i.e.,  $H_0 : \delta = 0$  is rejected if

$$T_n^+ \leq t_{\alpha/2} \quad \text{or} \quad T_n^+ \geq t_{1-\alpha/2}, \quad (13)$$

where the “critical values”  $t_{\alpha/2}$  and  $t_{1-\alpha/2}$  are the  $(\alpha/2)$ -quantile and the  $(1 - \alpha/2)$ -quantile of the distribution of  $T_n^+$ , respectively.

### 6.2.2 Distribution of the Test Statistic $T_n^+$ for Small Sample Sizes

- If the sample size  $n$  is not too large, then the quantiles  $t_{\alpha/2}$  and  $t_{1-\alpha/2}$  in (13) can be determined by combinatorial considerations.
  - Due to (10), the random vector  $\mathbf{R} = (R_1, \dots, R_n)$  of the ranks  $R_1, \dots, R_n$  of  $|X_1|, \dots, |X_n|$  is a (random) permutation of the numbers  $1, \dots, n$ .

– Then the test statistic  $T_n^+$  given in (11) can be represented as follows:

$$T_n^+ = \sum_{i=1}^n i Z_i, \quad \text{where } Z_i = \mathbb{I}_{\{X_{R_i^{-1}} > 0\}} \quad (14)$$

– and  $\mathbf{R}^{-1} = (R_1^{-1}, \dots, R_n^{-1})$  denotes the inverse permutation of  $\mathbf{R}$ , i.e., if  $R_i = j$ , then it holds that  $R_j^{-1} = i$ .

- Moreover, the following lemma is useful to determine the distribution of  $T_n^+$ .

**Lemma 6.1** *Assuming that  $H_0 : \delta = 0$  is true, it holds that:*

- The random vectors  $(\mathbb{I}_{\{X_1 > 0\}}, \dots, \mathbb{I}_{\{X_n > 0\}})$  and  $\mathbf{R} = (R_1, \dots, R_n)$  are independent.
- The components  $Z_1, \dots, Z_n$  of  $(Z_1, \dots, Z_n)$  are independent and identically distributed with  $Z_i \sim \text{Bin}(1, 1/2)$ .

**Proof**

- We first show that the random variables  $\mathbb{I}_{\{X_i > 0\}}$  and  $|X_i|$  are independent for each  $i = 1, \dots, n$ .
  - For each  $x \geq 0$  it holds that

$$\mathbb{P}(\mathbb{I}_{\{X_i > 0\}} = 1, |X_i| \leq x) = \mathbb{P}(0 < X_i \leq x) = G(x) - \frac{1}{2}$$

and

$$\mathbb{P}(\mathbb{I}_{\{X_i > 0\}} = 1) \mathbb{P}(|X_i| \leq x) = \frac{1}{2} (G(x) - G(-x)) = \frac{1}{2} (G(x) - (1 - G(x))) = G(x) - \frac{1}{2}.$$

- Moreover, it obviously holds for each  $x < 0$  that

$$\mathbb{P}(\mathbb{I}_{\{X_i > 0\}} = 1, |X_i| \leq x) = 0 = \mathbb{P}(\mathbb{I}_{\{X_i > 0\}} = 1) \mathbb{P}(|X_i| \leq x).$$

- In the same way, it can be shown that

$$\mathbb{P}(\mathbb{I}_{\{X_i > 0\}} = 0, |X_i| \leq x) = \mathbb{P}(\mathbb{I}_{\{X_i > 0\}} = 0) \mathbb{P}(|X_i| \leq x)$$

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

- Since the independence of the sampling variables  $X_1, \dots, X_n$  implies the independence of the random vectors  $(\mathbb{I}_{\{X_1 > 0\}}, |X_1|), \dots, (\mathbb{I}_{\{X_n > 0\}}, |X_n|)$ ,
  - it follows that  $(\mathbb{I}_{\{X_1 > 0\}}, \dots, \mathbb{I}_{\{X_n > 0\}})$  and  $(|X_1|, \dots, |X_n|)$  are independent random vectors.
  - Since  $\mathbf{R} = (R_1, \dots, R_n)$  is a Borel measurable function of  $(|X_1|, \dots, |X_n|)$ , also the random vectors  $(\mathbb{I}_{\{X_1 > 0\}}, \dots, \mathbb{I}_{\{X_n > 0\}})$  and  $\mathbf{R}$  are independent.
- Therefore, one gets for arbitrary  $i \in \{1, \dots, n\}$  and  $z \in \{0, 1\}$  that

$$\begin{aligned} \mathbb{P}(Z_i = z) &= \mathbb{P}(\mathbb{I}_{\{X_{R_i^{-1}} > 0\}} = z) \\ &= \sum_{\mathbf{r}} \mathbb{P}(\mathbb{I}_{\{X_{R_i^{-1}} > 0\}} = z \mid \mathbf{R} = \mathbf{r}) \mathbb{P}(\mathbf{R} = \mathbf{r}) \\ &= \sum_{\mathbf{r}} \mathbb{P}(\mathbb{I}_{\{X_{r_i} > 0\}} = z \mid \mathbf{R} = \mathbf{r}) \mathbb{P}(\mathbf{R} = \mathbf{r}) \\ &= \sum_{\mathbf{r}} \mathbb{P}(\mathbb{I}_{\{X_{r_i} > 0\}} = z) \mathbb{P}(\mathbf{R} = \mathbf{r}) = \frac{1}{2} \sum_{\mathbf{r}} \mathbb{P}(\mathbf{R} = \mathbf{r}) = \frac{1}{2}, \end{aligned}$$

where the summation extends over all permutations  $\mathbf{r} = (r_1, \dots, r_n)$  of the numbers  $1, \dots, n$ .

- This implies that

$$\begin{aligned}
\mathbb{P}(Z_1 = z_1, \dots, Z_n = z_n) &= \mathbb{P}(\mathbb{1}_{\{X_{R_1^{-1}} > 0\}} = z_1, \dots, \mathbb{1}_{\{X_{R_n^{-1}} > 0\}} = z_n) \\
&= \sum_{\mathbf{r}} \mathbb{P}(\mathbb{1}_{\{X_{R_1^{-1}} > 0\}} = z_1, \dots, \mathbb{1}_{\{X_{R_n^{-1}} > 0\}} = z_n \mid \mathbf{R} = \mathbf{r}) \mathbb{P}(\mathbf{R} = \mathbf{r}) \\
&= \sum_{\mathbf{r}} \mathbb{P}(\mathbb{1}_{\{X_{r_1^{-1}} > 0\}} = z_1, \dots, \mathbb{1}_{\{X_{r_n^{-1}} > 0\}} = z_n \mid \mathbf{R} = \mathbf{r}) \mathbb{P}(\mathbf{R} = \mathbf{r}) \\
&= \sum_{\mathbf{r}} \mathbb{P}(\mathbb{1}_{\{X_1 > 0\}} = z_{r_1}, \dots, \mathbb{1}_{\{X_n > 0\}} = z_{r_n} \mid \mathbf{R} = \mathbf{r}) \mathbb{P}(\mathbf{R} = \mathbf{r}) \\
&= \sum_{\mathbf{r}} \mathbb{P}(\mathbb{1}_{\{X_1 > 0\}} = z_{r_1}, \dots, \mathbb{1}_{\{X_n > 0\}} = z_{r_n}) \mathbb{P}(\mathbf{R} = \mathbf{r}) \\
&= \frac{1}{2^n} = \mathbb{P}(Z_1 = z_1) \dots \mathbb{P}(Z_n = z_n)
\end{aligned}$$

for arbitrary  $\mathbf{z} = (z_1, \dots, z_n) \in \{0, 1\}^n$ .

**Theorem 6.3** Assuming that  $H_0 : \delta = 0$  is true, the distribution of  $T_n^+$  is given by

$$\mathbb{P}(T_n^+ = k) = \frac{a_k}{2^n} \quad \forall k = 0, 1, \dots, n, \quad (15)$$

where

$$a_k = \#\left\{ \mathbf{z} = (z_1, \dots, z_n) \in \{0, 1\}^n : \sum_{i=1}^n i z_i = k \right\}. \quad (16)$$

Moreover, it holds that

$$\mathbb{E} T_n^+ = \frac{n(n+1)}{4} \quad \text{and} \quad \text{Var} T_n^+ = \frac{n(n+1)(2n+1)}{24}. \quad (17)$$

**Proof**

- The representation formula (14) for  $T_n^+$  and Lemma 6.1 imply that

$$\mathbb{P}(T_n^+ = k) = \mathbb{P}\left(\sum_{i=1}^n i Z_i = k\right) = \sum_{\mathbf{z}=(z_1, \dots, z_n) \in \{0,1\}^n : \sum_{i=1}^n i z_i = k} \mathbb{P}(Z_1 = z_1, \dots, Z_n = z_n) = \frac{a_k}{2^n}$$

for each  $k = 0, 1, \dots, n$ .

- Furthermore, one obtains that

$$\mathbb{E} T_n^+ = \mathbb{E}\left(\sum_{i=1}^n i Z_i\right) = \sum_{i=1}^n i \mathbb{E} Z_i = \frac{1}{2} \binom{n+1}{2} = \frac{n(n+1)}{4}$$

and

$$\text{Var} T_n^+ = \text{Var}\left(\sum_{i=1}^n i Z_i\right) = \sum_{i=1}^n i^2 \text{Var} Z_i = \frac{1}{4} \sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{24}. \quad \square$$

**Remark**

- From (12) and (14) it moreover follows with Lemma 6.1 that

$$T_n^- = \binom{n+1}{2} - T_n^+ = \binom{n+1}{2} - \sum_{i=1}^n i Z_i = \sum_{i=1}^n i(1 - Z_i) \stackrel{d}{=} \sum_{i=1}^n i Z_i = T_n^+,$$

i.e., under  $H_0 : \delta = 0$  it holds that

$$T_n^- \stackrel{d}{=} T_n^+. \quad (18)$$

- Thus, (12) implies that for each  $k = 0, 1, \dots, n$

$$\mathbb{P}(T_n^+ = k) = \mathbb{P}\left(T_n^- = \frac{n(n+1)}{2} - k\right) = \mathbb{P}\left(T_n^+ = \frac{n(n+1)}{2} - k\right),$$

i.e., the distribution of  $T_n^+$  is symmetric with respect to the expectation  $\mathbb{E}T_n^+ = n(n+1)/4$ .

- This means that also the quantiles  $t_{\alpha,n} = \max\{t \in \mathbb{R} : \mathbb{P}(T_n^+ \leq t) \leq \alpha\}$  have this property of symmetry, i.e., for each  $\alpha \in (0, 1)$  it holds that

$$t_{\alpha,n} = \frac{n(n+1)}{2} - t_{1-\alpha,n}.$$

- The quantiles  $t_{\alpha,n}$  can either be taken from tables or be determined using Monte Carlo simulation.

### 6.2.3 Asymptotic Distribution

- The direct determination of the quantiles  $t_{\alpha/2}$  and  $t_{1-\alpha/2}$  by using Theorem 6.3 is difficult if the sample size  $n$  is large.
  - Another way to approximatively determine the distribution of the test statistic  $T_n^+$  is based on the representation formula (14).
  - In this context the fact is used that  $T_n^+ = \sum_{i=1}^n i Z_i$  is a sum of independent random variables, which follows from Lemma 6.1.
  - In fact, the central limit theorem for sums of independent (but not necessary identically distributed) random variables implies that  $T_n^+$  is normally distributed.
- For this purpose we consider the following stochastic model: For each  $n \geq 1$  let  $X_{n1}, \dots, X_{nn} : \Omega \rightarrow \mathbb{R}$  be a sequence of independent random variables,
  - where we (w.l.o.g.) assume that for each  $k \in \{1, \dots, n\}$

$$\mathbb{E} X_{nk} = 0, \quad 0 < \sigma_{nk}^2 = \text{Var} X_{nk} < \infty \quad \text{and} \quad \sum_{k=1}^n \sigma_{nk}^2 = 1. \quad (19)$$

- If the random variables  $X_{n1}, \dots, X_{nn}$  do not satisfy the conditions formulated in (19), then we consider the transformed random variables  $X'_{n1}, \dots, X'_{nn}$  with

$$X'_{nk} = \frac{X_{nk} - \mathbb{E} X_{nk}}{\sqrt{n \text{Var} X_{nk}}}. \quad (20)$$

- We denote the distribution function of  $X_{nk}$  by  $F_{nk}$ , where we do not exclude the case that  $F_{nk}$  for each  $k \in \{1, \dots, n\}$  can depend on the number  $n$  of considered random variables  $X_{n1}, \dots, X_{nn}$ .

The following *central limit theorem of Lindeberg* (cf. Theorem WR-5.22) is the basis to show that  $T_n^+$  is asymptotically normally distributed.

#### Lemma 6.2

- For each  $n \in \mathbb{N}$  let  $X_{n1}, \dots, X_{nn} : \Omega \rightarrow \mathbb{R}$  be a sequence of independent random variables, which satisfy the conditions (19).

- If furthermore for each  $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} \sum_{k=1}^n \int_{\mathbb{R} \setminus (-\varepsilon, \varepsilon)} x^2 dF_{nk}(x) = 0, \quad (21)$$

then it holds for each  $x \in \mathbb{R}$  that

$$\lim_{n \rightarrow \infty} P(X_{n1} + \dots + X_{nn} \leq x) = \Phi(x), \quad (22)$$

where  $\Phi : \mathbb{R} \rightarrow [0, 1]$  is the distribution function of the  $N(0, 1)$ -distribution.

**Theorem 6.4** Under  $H_0 : \delta = 0$  it holds that

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\frac{T_n^+ - \mathbb{E} T_n^+}{\sqrt{\text{Var} T_n^+}} \leq x\right) = \Phi(x) \quad \forall x \in \mathbb{R}. \quad (23)$$

**Proof**

- Because of (14) it is sufficient to show that the random variables  $X_{n1}, \dots, X_{nn}$  with

$$X_{nk} = \frac{kZ_k - k \mathbb{E} Z_k}{\sqrt{\text{Var} T_n^+}} \quad (24)$$

satisfy the conditions of Lemma 6.2.

- It follows directly from equation (24) that (19) is fulfilled.
- Therefore, it merely remains to show that the Lindeberg-condition (22) is satisfied.
- For the distribution function  $F_{nk} : \mathbb{R} \rightarrow [0, 1]$  of the random variable  $X_{nk}$ , introduced in (24), it follows from Lemma 6.1 that

$$F_{nk}(x) = \begin{cases} 0, & \text{if } x < \frac{-k}{2\sqrt{\text{Var} T_n^+}}, \\ \frac{1}{2}, & \text{if } \frac{-k}{2\sqrt{\text{Var} T_n^+}} \leq x < \frac{k}{2\sqrt{\text{Var} T_n^+}}, \\ 1, & \text{if } \frac{k}{2\sqrt{\text{Var} T_n^+}} \leq x. \end{cases}$$

- This implies that  $\int_{\mathbb{R} \setminus (-\varepsilon, \varepsilon)} x^2 dF_{nk}(x) = 0$  for each  $k \in \{1, \dots, n\}$  if  $n$  is chosen in such a way that

$$\frac{n^2}{4\text{Var} T_n^+} = \frac{6n^2}{n(n+1)(2n+1)} < \varepsilon^2,$$

where the last equality follows from the formula for  $\text{Var} T_n^+$  in Theorem 6.3.

- Therefore, the validity of the Lindeberg-condition (22) is shown.  $\square$

**Remark**

- Because of Theorem 6.4 the following critical area is considered in the case of the (two-sided) test problem  $H_0 : \delta = 0$  vs.  $H_1 : \delta \neq 0$ .
- For sufficiently large  $n$ ,  $H_0 : \delta = 0$  is rejected if

$$\left| \frac{T_n^+ - \mathbb{E} T_n^+}{\sqrt{\text{Var} T_n^+}} \right| \geq z_{1-\alpha/2}, \quad (25)$$

where  $\mathbb{E} T_n^+$  or  $\text{Var} T_n^+$  are given in Theorem 6.3 and  $z_{1-\alpha/2}$  is the  $(1 - \alpha/2)$ -quantile of the  $N(0, 1)$ -distribution.

- In the literature the condition  $n \geq 20$  is suggested as a possible criterion for “sufficiently large”.

### 6.3 Two-Sample Problems

- In this section we discuss nonparametric tests for the case that two independent random samples  $(X_1, \dots, X_{n_1})$  and  $(Y_1, \dots, Y_{n_2})$  are observed.
- In other words: We assume that the random variables  $X_1, \dots, X_{n_1}, Y_1, \dots, Y_{n_2}$  are completely independent with the (unknown) distribution functions  $F$  and  $G$ , i.e.,

$$F(x) = \mathbb{P}(X_i \leq x) \quad \text{and} \quad G(y) = \mathbb{P}(Y_j \leq y) \quad \forall x, y \in \mathbb{R}, \quad i = 1, \dots, n, \quad j = 1, \dots, m.$$

- Then, a (two-sided) test problem is for example given by

$$H_0 : F(x) = G(x) \quad \forall x \in \mathbb{R} \quad \text{vs.} \quad H_1 : F(x) \neq G(x) \quad \exists x \in \mathbb{R}. \quad (26)$$

- As a one-sided alternative to  $H_0 : F(x) = G(x) \quad \forall x \in \mathbb{R}$  the following hypotheses can be considered:

$$H_1 : F(x) \geq G(x) \quad \forall x \in \mathbb{R} \quad \text{and} \quad F(x) > G(x) \quad \exists x \in \mathbb{R} \quad (27)$$

or

$$H_1 : F(x) \leq G(x) \quad \forall x \in \mathbb{R} \quad \text{and} \quad F(x) < G(x) \quad \exists x \in \mathbb{R}. \quad (28)$$

#### 6.3.1 Run Test of Wald–Wolfowitz

- For the analysis of the test problem, given in (26), one can apply the run test for randomness, which has been discussed in Section 6.1.2.

- For this purpose, we combine the sample variables  $X_1, \dots, X_{n_1}$  and  $Y_1, \dots, Y_{n_2}$  to one random sample

$$(X'_1, \dots, X'_n) = (X_1, \dots, X_{n_1}, Y_1, \dots, Y_{n_2}), \quad \text{where } n = n_1 + n_2,$$

and consider the ordered sample  $(X'_{(1)}, \dots, X'_{(n)})$ .

- Here we assume that the distribution functions  $F$  and  $G$  are continuous, i.e., the mapping

$$(X'_1, \dots, X'_n) \mapsto (X'_{(1)}, \dots, X'_{(n)})$$

is uniquely determined with probability 1.

- Under  $H_0 : F(x) = G(x) \quad \forall x \in \mathbb{R}$  it is to be expected that the  $X_i$ 's and  $Y_j$ 's in  $(X'_{(1)}, \dots, X'_{(n)})$  are “well mixed”,

- since the sample variables  $X'_{(1)}, \dots, X'_{(n)}$  then are independent and identically distributed.
- If the trend for “clumping and clustering” is considered as an alternative, then  $H_0$  is rejected if the number  $T$  of iterations in the (binary) sample  $(Z_1, \dots, Z_n)$  is “too small”, where  $Z_i = 0$  if  $X'_{(i)} = X_j$  for some  $j \in \{1, \dots, n\}$  and  $Z_i = 1$  if  $X'_{(i)} = Y_j$  for some  $j \in \{1, \dots, n\}$ .

#### Examples

- In a medical study the body heights of  $n_1 = 8$  girls and  $n_2 = 10$  boys were analyzed.
- The measurement results are:

$x_i$	117	121	122	124	125	126	128	132		
$y_j$	110	113	114	115	116	118	119	120	123	127

- If we order these measurements by size and assign a 0 to the heights of the boys and a 1 to the heights of the girls, then we obtain the sequence

$$\omega = (0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1) \quad \text{with } T(\omega) = 8. \quad (29)$$

- Otherwise, Theorem 6.1 implies that for the  $\alpha$ -quantile  $r_\alpha(n_1; n_2)$  of the distribution of  $T$  it holds that  $r_{0.05}(8; 10) = 6$  for  $\alpha = 0.05$ .
- In this case  $H_0$  is therefore rejected because  $T(\omega) = 8 > 6 = r_{0.05}(8; 10)$ .

### Remark

- The run test considered in this section is not able to identify alternatives of the type (27) or (28).
- The example given in (29) makes this clear: Since the number of iterations  $T(\omega) = 8$  does not change, if we (in contrast to the previous approach) assign a 1 to the heights of the boys and a 0 to the heights of the girls.
- Also for two-sided alternatives the run test of Wald–Wolfowitz, also called “omnibus–test”, should only be used if the form of the alternative is not specified further.
- For special alternatives, which for example only affect location or variability characteristics, other test methods are more efficient, cf Section 6.3.2.

### 6.3.2 Wilcoxon Rank–Sum Test for Location Alternatives

- We now discuss another nonparametric test for the case that two independent random samples  $(X_1, \dots, X_{n_1})$  and  $(Y_1, \dots, Y_{n_2})$  are observed.
- However, we will here consider more special alternatives as in (26) – (28).
  - We assume that the random variables  $X_1, \dots, X_{n_1}$  and  $Y_1, \dots, Y_{n_2}$  are completely independent with the (unknown) continuous distribution functions  $F$  and  $G$ .
  - Similar as in Section 6.2 it is assumed that there is some  $\delta \in \mathbb{R}$  such that

$$F(x) = G(x + \delta) \quad \forall x \in \mathbb{R}.$$

- A (two-sided) test problem, which is consistent with the above mentioned more general test problem (26), is then given by

$$H_0 : \delta = 0 \quad \text{vs.} \quad H_1 : \delta \neq 0. \quad (30)$$

- The following hypotheses can be considered as one-sided alternatives to  $H_0 : \delta = 0$ :

$$H_1 : \delta > 0 \quad \text{or} \quad H_1 : \delta < 0. \quad (31)$$

- In the same way as in Section 6.3.1 we merge the sample variables  $X_1, \dots, X_{n_1}$  and  $Y_1, \dots, Y_{n_2}$  to one combined random sample  $(X'_1, \dots, X'_n) = (X_1, \dots, X_{n_1}, Y_1, \dots, Y_{n_2})$ , where  $n = n_1 + n_2$ .
  - Furthermore, we consider the (random) vector of the ranks  $\mathbf{R}' = (R'_1, \dots, R'_n)$  of the sample variables  $X'_1, \dots, X'_n$  in the combined sample, where

$$R'_i = \sum_{j=1}^n \mathbb{I}_{\{X'_j \leq X'_i\}} \quad \forall i = 1, \dots, n.$$

- As in Section 6.3.1 it has to be expected under  $H_0 : \delta = 0$  that the  $X_i$ 's and  $Y_j$ 's in the combined sample  $(X'_{(1)}, \dots, X'_{(n)})$  are “well mixed” because then the sample variables  $X'_{(1)}, \dots, X'_{(n)}$  are independent and identically distributed.

– Thus, for the two-sided test problem in (30),  $H_0$  is rejected if the *rank-sum*

$$T_{n_1, n_2} = \sum_{i=1}^{n_1} R'_i \quad (32)$$

is “too small” or “too large”.

- In order to perform the test, the distribution of the test statistic  $T_{n_1, n_2}$ , introduced in (32), has to be determined. For this purpose the following lemma is useful.

**Lemma 6.3**

- Let  $X : \Omega \rightarrow \{\dots, -1, 0, 1, \dots\}$  be a discrete random variable such that  $\mathbb{E}|X| < \infty$  and that for some  $\mu \in \mathbb{R}$  the following symmetry property is fulfilled:

$$\mathbb{P}(X = \mu - k) = \mathbb{P}(X = \mu + k) \quad \forall k \in \{\dots, -1, 0, 1, \dots\}. \quad (33)$$

- Then it holds that  $\mathbb{E}X = \mu$ .

**Proof**

- We can w.l.o.g. assume that  $\mu = 0$  since otherwise the transformed random variable  $X' = X - \mu$  can be considered.
- Then it follows from (33) with  $\mu = 0$  that

$$\mathbb{E}X = \sum_{k=-\infty}^{\infty} k \mathbb{P}(X = k) = - \sum_{k=1}^{\infty} k \mathbb{P}(X = -k) + \sum_{k=1}^{\infty} k \mathbb{P}(X = k) \stackrel{(33)}{=} 0. \quad \square$$

**Theorem 6.5**

- Under  $H_0 : \delta = 0$  the distribution of  $T_{n_1, n_2}$  is given by

$$\mathbb{P}(T_{n_1, n_2} = k) = \frac{a_{k, n_1, n_2}}{\binom{n_1 + n_2}{n_1}} \quad \forall k = \frac{n_1(n_1 + 1)}{2}, \dots, n_1 n_2 + \frac{n_1(n_1 + 1)}{2}, \quad (34)$$

where

$$a_{k, n_1, n_2} = \#\left\{ \mathbf{z} = (z_1, \dots, z_{n_1 + n_2}) \in \{0, 1\}^{n_1 + n_2} : \#\{i : z_i = 1\} = n_1, \sum_{i=1}^{n_1 + n_2} iz_i = k \right\}. \quad (35)$$

- Furthermore, it holds that

$$\mathbb{P}(T_{n_1, n_2} = k) = \mathbb{P}(T_{n_1, n_2} = 2\mu - k) \quad \forall k \in \{\dots, -1, 0, 1, \dots\} \quad (36)$$

and therefore

$$\mathbb{E}T_{n_1, n_2} = \mu, \quad (37)$$

where  $\mu = n_1(n_1 + n_2 + 1)/2$ .

**Proof**

- Under  $H_0 : \delta = 0$  the sample variables  $X'_1, \dots, X'_{n_1 + n_2}$  are independent and identically distributed.



- Hence, each of the  $\binom{n_1+n_2}{n_1}$  partitions of the  $n_1$  variables  $X_1, \dots, X_{n_1}$  into the  $n_1 + n_2$  existing rank spots has the same probability.
- Moreover, it holds for the minimum and maximum value  $t_{\min}$  and  $t_{\max}$ , respectively, of  $T_{n_1, n_2}$  that

$$t_{\min} = \sum_{i=1}^{n_1} i = \frac{n_1(n_1+1)}{2} \quad \text{and} \quad t_{\max} = \sum_{i=n_2+1}^{n_2+n_1} i = n_1 n_2 + \frac{n_1(n_1+1)}{2}.$$

- From this the validity of (34) – (35) is obtained.
- In order to prove (36) we use the following symmetry property.
  - Each  $\mathbf{z} = (z_1, \dots, z_{n_1+n_2}) \in \{0, 1\}^{n_1+n_2}$  with

$$\#\{i : z_i = 1\} = n_1 \quad \text{and} \quad \sum_{i=1}^{n_1+n_2} i z_i = k$$

corresponds to some  $\tilde{\mathbf{z}} = (\tilde{z}_1, \dots, \tilde{z}_{n_1+n_2}) \in \{0, 1\}^{n_1+n_2}$  with

$$\#\{i : \tilde{z}_{n_1+n_2+1-i} = 1\} = n_1 \quad \text{and} \quad \sum_{i=1}^{n_1+n_2} (n_1 + n_2 + 1 - i) \tilde{z}_i = n_1(n_1 + n_2 + 1) - k.$$

- Since the sample variables  $X'_{(1)}, \dots, X'_{(n)}$  are independent and identically distributed, it thus follows for each  $k \in \{\dots, -1, 0, 1, \dots\}$  that

$$\mathbb{P}(T_{n_1, n_2} = k) = \mathbb{P}(T_{n_1, n_2} = n_1(n_1 + n_2 + 1) - k) = \mathbb{P}(T_{n_1, n_2} = 2\mu - k), \quad (38)$$

where  $2\mu = n_1(n_1 + n_2 + 1)$ .

- In order to show (37) it is sufficient to substitute  $k = \mu - i$  in (38).
  - Then (38) implies that

$$\mathbb{P}(T_{n_1, n_2} = \mu - i) = \mathbb{P}(T_{n_1, n_2} = \mu + i) \quad \forall i \in \{\dots, -1, 0, 1, \dots\}.$$

- From this and Lemma 6.3 the validity of (37) is obtained. □

### Remark

- Now, (38) implies the following symmetry property for the quantiles  $t_{\alpha, n_1, n_2}$  of  $T_{n_1, n_2}$ .
  - For each  $\alpha \in (0, 1)$  it holds that

$$t_{\alpha, n_1, n_2} = n_1(n_1 + n_2 + 1) - t_{1-\alpha, n_1, n_2}.$$

- The quantiles  $t_{\alpha, n_1, n_2}$  can either be taken from tables or be determined using Monte Carlo simulation.
- The null hypothesis  $H_0 : \delta = 0$  is rejected in favor of  $H_1 : \delta \neq 0$  if

$$T_{n_1, n_2} \leq t_{\alpha/2, n_1, n_2} \quad \text{or} \quad T_{n_1, n_2} \geq n_1(n_1 + n_2 + 1) - t_{\alpha/2, n_1, n_2}.$$

- Analogously, the null hypothesis  $H_0 : \delta = 0$  is rejected in favor of  $H_1 : \delta < 0$  or  $H_1 : \delta > 0$  if

$$T_{n_1, n_2} \geq n_1(n_1 + n_2 + 1) - t_{\alpha, n_1, n_2} \quad \text{or} \quad T_{n_1, n_2} \leq t_{\alpha, n_1, n_2}.$$

If the sample sizes  $n_1$  and  $n_2$  are large enough, it is difficult to determine the quantiles  $t_{\alpha, n_1, n_2}$  directly via Theorem 6.5. However, then the distribution of the test statistic  $T_{n_1, n_2}$  can approximatively be determined from the following central limit theorem.

**Theorem 6.6** *If  $n_1, n_2 \rightarrow \infty$  such that  $n_1/(n_1 + n_2) \rightarrow p$  or  $n_2/(n_1 + n_2) \rightarrow 1 - p$  for some  $p \in (0, 1)$ , then it holds that*

$$\lim_{n_1, n_2 \rightarrow \infty} \mathbb{P} \left( \frac{T_{n_1, n_2} - \mathbb{E} T_{n_1, n_2}}{\sqrt{\text{Var} T_{n_1, n_2}}} \leq x \right) = \Phi(x) \quad \forall x \in \mathbb{R}, \quad (39)$$

where

$$\mathbb{E} T_{n_1, n_2} = \frac{n_1(n_1 + n_2 + 1)}{2}, \quad \text{Var} T_{n_1, n_2} = \frac{n_1 n_2 (n_1 + n_2 + 1)}{12}$$

and  $\Phi : \mathbb{R} \rightarrow [0, 1]$  is the distribution function of the  $N(0, 1)$ -distribution.

### Remark

- Because of Theorem 6.6 the null hypothesis  $H_0 : \delta = 0$  is rejected for large  $n_1, n_2$  in favor of  $H_1 : \delta \neq 0$  if

$$\left| \frac{T_{n_1, n_2} - n_1(n_1 + n_2 + 1)/2}{\sqrt{n_1 n_2 (n_1 + n_2 + 1)/12}} \right| \geq z_{1-\alpha/2},$$

where  $z_\alpha$  is the  $\alpha$ -quantile of the  $N(0, 1)$ -distribution.

- Analogously,  $H_0 : \delta = 0$  is rejected in favor of  $H_1 : \delta > 0$  or  $H_1 : \delta < 0$  if

$$\frac{T_{n_1, n_2} - n_1(n_1 + n_2 + 1)/2}{\sqrt{n_1 n_2 (n_1 + n_2 + 1)/12}} \geq z_{1-\alpha} \quad \text{or} \quad \frac{T_{n_1, n_2} - n_1(n_1 + n_2 + 1)/2}{\sqrt{n_1 n_2 (n_1 + n_2 + 1)/12}} \leq -z_{1-\alpha}.$$