

Lectures on Mathematics

Peter Tallos

CUB, Department of Mathematics

September, 2020

Preface

This book covers the necessary mathematics intended for students with major in *Economics*. This course provides the fundamental mathematical background for studying *Microeconomics*, *Macroeconomics*, *Statistics* and other important topics in economics, or probabilistic or stochastic disciplines. The main mathematical topics covered are *Mathematical Analysis (Calculus)*, *Probability Theory*, and *Linear Algebra*.

Most of the time, we avoid the rigorous mathematical proofs of our statements. Instead, we rather present "justifications", which are intuitive, but not necessarily precise. However, emphasis is placed on the correct formulations of definitions. Some paragraphs indicated by the word ATTENTION cover somewhat more complicated arguments, their detailed explanations are given in the classroom.

The text is illustrated by a large number of examples. On the one hand, they help the deeper understanding. On the other hand, they give an idea, how to apply them in practical situations. Therefore, the thorough study of examples is a profoundly important homework assignment. Each chapter covers one week of the semester, on a one week – one chapter basis.

In the end of each chapter references are given to the *Textbook*, which should be interpreted the following way.

Textbook-1: K. Sydsaeter and P. Hammond, *Mathematics for Economic Analysis*, Prentice Hall, 1995, ISBN 0–13–583600–X, or any of the later editions.

Textbook-2 R. E. Walpole, R. H. Myers, S. L. Myers and K. Ye *Probability and Statistics for Engineers and Scientists*, Prentice Hall, 2012, ISBN: 978–0–321–62911–1, or any of the later editions.

These textbooks are widely used at most recognized universities worldwide.

Some of the indicated homework exercises refer to the *Textbook*. Most of the midterm quiz or final exam problems are similar or identical to those exercises. More problems and exercises with solutions are posted on my web site. These files are updated regularly.

Special thanks to my colleagues Csaba Puskás, Éva Ernyes and Balázs Fleiner, who read the manuscript, and their valuable remarks significantly improved the quality of the text. My thanks go to my former students as well, their comments in or outside the classroom were extremely helpful for making the text more understandable.

Budapest, September, 2020.

Peter Tallos

Contents

I	First Semester: Differential and Integral Calculus	11
1	Sequences	13
1.1	Limits of sequences	13
1.2	Sequences tending to infinity	14
1.3	Squeezing Theorem	15
1.4	Bounded and monotone sequences	16
1.5	Euler's number e	17
2	Infinite Series	19
2.1	Series	19
2.2	Geometric series	20
2.3	Convergence based on examining the partial sums	20
2.4	Conditions for convergence	21
2.5	Absolute convergence	22
2.6	Quotient-test	24
3	Limits and continuity	27
3.1	Basic concepts	27
3.2	Squeezing theorem	29
3.3	One-sided limits	29
3.4	Continuity	30
4	Differentiation of functions	33
4.1	The derivative	33
4.2	Tangent lines	34
4.3	Rules of differentiation	35
4.4	Composition of functions	36
4.5	Chain-Rule	37

5	The Mean Value Theorem	39
5.1	The inverse function	39
5.2	Differentiability of the inverse function	40
5.3	The exponential and logarithm functions	41
5.4	Necessary condition for an extremum	42
5.5	Lagrange's Mean Value Theorem	43
5.6	L'Hôpital's Rule	44
6	Complete analysis of functions	45
6.1	Monotone functions	45
6.2	Finding extreme points	46
6.3	Higher order derivatives	47
6.4	Second order conditions	48
6.5	Convex and concave functions	50
7	Integration	53
7.1	The indefinite integral	53
7.2	Basic integrals	54
7.3	Initial value problems	55
7.4	Definite integrals	55
7.5	Newton-Leibniz-formula	57
8	Methods of integration	59
8.1	Integration by parts	59
8.2	Integration by parts in definite integrals	60
8.3	Integration by substitution	61
8.4	Substitution in definite integrals	62
8.5	Linear differential equations	62
9	Extension of integration	65
9.1	Improper integrals	65
9.2	Improper integrals on the real line	66
9.3	Integration by parts in improper integrals	68
9.4	Harmonic series revisited	69
10	Power series	71
10.1	Sum of power series	71
10.2	Radius of convergence	72
10.3	Differentiability of power series	73
10.4	Finding the coefficients	75
10.5	Taylor-series of the exponential function	75

<i>CONTENTS</i>	5
11 Functions of two variables	77
11.1 Partial derivatives	77
11.2 Tangent planes	78
11.3 Chain Rule	79
11.4 Local extrema	80
11.5 First order necessary condition	81
12 Constrained extrema	83
12.1 Implicit functions	83
12.2 Constrained minima	85
12.3 Lagrange multipliers	86
12.4 Solving the constrained minimization problem	87
II Second Semester: Probability Theory	89
13 Probability	91
13.1 Experiments	91
13.2 The sample space	91
13.3 Events	92
13.4 Operations with events	92
13.5 Probability space	94
14 Sampling methods	97
14.1 Classical probability spaces	97
14.2 Sampling without replacement	99
14.3 Sampling with replacement	100
14.4 The Bernoulli experiment	101
15 Conditional probability and Bayes' Rule	103
15.1 Conditional probability	103
15.2 Independence	104
15.3 Theorem of Total Probability	105
15.4 Bayes' Rule	107
16 Random variables and distributions	109
16.1 Random variables	109
16.2 Distribution of discrete variables	110
16.3 The cumulative distribution function	111
16.4 The density function	112

17 Mean and variance	115
17.1 Mean of discrete distributions	115
17.2 Mean of infinite distributions	117
17.3 Mean of continuous distributions	118
17.4 Basic properties of the mean	118
17.5 Variance and standard deviation	119
18 Special discrete distributions	121
18.1 Characteristic distribution	121
18.2 Binomial distribution	122
18.3 Hypergeometric distribution	122
18.4 Geometric distribution	123
18.5 Poisson distribution	124
19 Special continuous distributions	127
19.1 Uniform distribution	127
19.2 Exponential distribution	128
19.3 The standard normal distribution	129
19.4 Normal distribution	131
20 Joint distributions	133
20.1 Joint cumulative distribution function	133
20.2 Discrete joint distributions	133
20.3 Continuous joint distributions	135
20.4 Independence	136
20.5 Conditional distributions	138
21 Covariance and correlation	139
21.1 Mean of a sum	139
21.2 Mean of a product	140
21.3 Variance of a sum	141
21.4 Covariance and correlation	142
21.5 Theorem of Total Expectation	143
22 Sums of random variables	145
22.1 Sums of discrete variables	145
22.2 Sums of continuous variables	145
22.3 The Poisson process	147
22.4 Sum of standard normal distributions	148
22.5 Central Limit Theorem	149

23 Law of Large Numbers	151
23.1 Chebyshev's Theorem	151
23.2 Chebyshev's Theorem in equivalent form	153
23.3 Poisson approximation	154
23.4 Law of Large Numbers	154
III Third Semester: Linear algebra	157
25 Vector spaces and subspaces	159
25.1 The vector space \mathbb{R}^n	159
25.2 Subspaces	160
25.3 Generated subspace	161
25.4 Linear independence	162
26 Linear independence and basis	165
26.1 Generating system	165
26.2 Basis	166
26.3 Dimension	167
26.4 Gauss-Jordan-elimination	168
27 Linear mappings and matrices	171
27.1 Linear mappings	171
27.2 Matrix of a linear map	172
27.3 Rank and degree of freedom of a matrix	173
27.4 Multiplication of matrices	175
28 Linear systems	177
28.1 Homogeneous systems	177
28.2 Inhomogeneous systems	178
28.3 Inverse of a matrix	181
28.4 Finding the inverse	182
29 Eigenvalue, eigenvector	185
29.1 Eigenvalue, eigenvector	185
29.2 Eigensubspace	186
29.3 Finding eigenvectors	187
29.4 Independent eigenvectors	188
29.5 Diagonal form of transformations	188

30 Determinant	191
30.1 Permutations	191
30.2 The determinant	192
30.3 Properties of the determinant	192
30.4 Evaluating the determinant	194
30.5 Finding the eigenvalues	195
31 Scalar product	197
31.1 Scalar product	197
31.2 Angle of vectors, perpendicularity	198
31.3 Orthogonal systems	199
31.4 Gram-Schmidt-procedure	200
31.5 Orthogonal complement	200
32 The spectral theorem	203
32.1 Transpose of a matrix	203
32.2 Orthogonal matrices	204
32.3 Symmetric matrices	205
32.4 Spectral theorem of symmetric matrices	206
33 Quadratic forms	209
33.1 Quadratic forms	209
33.2 Symmetric matrix of a quadratic form	210
33.3 Definite quadratic forms	211
33.4 Completing the square	212
33.5 Definite property based on eigenvalues	213
34 Functions with several variables	215
34.1 Partial derivatives	215
34.2 The derivative	216
34.3 Chain-rule	217
34.4 Second order partial derivatives	219
34.5 Young's theorem	220
35 Local extrema	223
35.1 Local extrema	223
35.2 First order necessary condition	223
35.3 Second order necessary condition	224
35.4 Sufficient condition for local extrema	225
35.5 Finding the extreme values	226
35.6 The special case of two variables	227

36 Least squares method, regression	229
36.1 Least squares method	229
36.2 Analytic solution	230
36.3 Algebraic solution	231
36.4 Regression	231

Part I

First Semester: Differential and Integral Calculus

Chapter 1

Sequences

1.1 Limits of sequences

The function

$$a : \mathbb{N} \rightarrow \mathbb{R}$$

defined on the set of natural numbers \mathbb{N} is called a (infinite) sequence.

We use the notation a_n for the n -th element.

Some examples: $a_n = n$, $a_n = \frac{1}{n}$, $a_n = \frac{n+1}{n+2}$.

Definition 1.1 The sequence a_n is said to be *convergent* and tends to A , if for any $\varepsilon > 0$, there exists an index N , such that,

$$|a_n - A| < \varepsilon.$$

whenever $n \geq N$.

If the sequence is convergent then A is called the limit of the sequence a_n and we write

$$\lim_{n \rightarrow \infty} a_n = A.$$

If there is no such real number A , then the sequence is called *divergent*.

Theorem 1.2 *If the sequences a_n and b_n are convergent and $\lim_{n \rightarrow \infty} a_n = A$ and $\lim_{n \rightarrow \infty} b_n = B$ then*

- $\lim_{n \rightarrow \infty} (a_n \pm b_n) = A \pm B$,
- $\lim_{n \rightarrow \infty} (a_n \cdot b_n) = A \cdot B$,
- if $B \neq 0$, then $\lim_{n \rightarrow \infty} \frac{a_n}{b_n} = \frac{A}{B}$

Example 1.3 Let us consider the sequence $a_n = \frac{1}{n}$. For an arbitrary $\varepsilon > 0$ let N be any integer, greater than $1/\varepsilon$. Then if $n \geq N$

$$\frac{1}{n} < \varepsilon,$$

therefore, in view of Definition 1.1

$$\lim_{n \rightarrow \infty} \frac{1}{n} = 0.$$

Example 1.4 In a similar way we can find the limits of other sequences. Let us consider for example the sequence

$$a_n = \frac{2n^2 + 5}{n^2 - 6n + 8}.$$

If we divide both the numerator and the denominator by n^2 , then we have

$$a_n = \frac{2 + 5/n^2}{1 - 6/n + 8/n^2},$$

where the limit of the numerator is 2 and the limit of the denominator is 1. Therefore

$$\lim_{n \rightarrow \infty} a_n = 2.$$

Every irrational number can be written as a limit of a sequence of rational numbers. For example, consider the sequence $a_1 = 1.4$, $a_2 = 1.41$, $a_3 = 1.414$, $a_4 = 1.4142 \dots$ then

$$\lim_{n \rightarrow \infty} a_n = \sqrt{2}$$

Indeed, according to Definition 1.1, if $\varepsilon = 10^{-N}$, then $|a_n - \sqrt{2}| < \varepsilon$ for $n \geq N$.

A typical example for a sequence which has no limit is

$$a_n = (-1)^n.$$

1.2 Sequences tending to infinity

Let us investigate the sequence

$$a_n = 2n + 5.$$

The terms of this sequence are greater than any given number K if n is large enough. In that case, we say, that the limit of the sequence is infinity. We use the symbol ∞ to denote infinity.

Definition 1.5 We say that the sequence a_n approaches $+\infty$ if for any real number K there exists an index N such that for every $n \geq N$ we have $a_n > K$. This is expressed in the formula

$$\lim_{n \rightarrow \infty} a_n = +\infty ,$$

In a completely analogous way we can define the fact that a sequence approaches $-\infty$, that is $\lim_{n \rightarrow \infty} a_n = -\infty$.

1.3 Squeezing Theorem

Often the limit of a sequence can be determined with the aid of other sequences the limits of which are known. Such a situation is described by the Squeezing Theorem.

Theorem 1.6 (Squeezing Theorem) *Let a_n , b_n and c_n be sequences such that for every index n*

$$a_n \leq b_n \leq c_n$$

holds and, moreover, the sequences a_n and c_n converge to the same limit A . Then the sequence b_n is also convergent and $\lim_{n \rightarrow \infty} b_n = A$.

Example 1.7 Let $a > 1$ be a real number and consider the sequence $b_n = \sqrt[n]{a}$. Since $a > 1$, the elements of the sequence can be written in the form

$$\sqrt[n]{a} = 1 + h_n ,$$

where $h_n > 0$ for every n . By the Binomial Theorem we get

$$a = (1 + h_n)^n > 1 + nh_n .$$

where we skipped all other positive terms on the right-hand side. Rearranging the inequality it follows that

$$0 < h_n < \frac{a - 1}{n} .$$

The expression on the right-hand side tends to zero, hence, by the Squeezing Theorem $h_n \rightarrow 0$, that is $\sqrt[n]{a} \rightarrow 1$.

Obviously, if $0 < a \leq 1$, then we can carry out the same argument, by taking reciprocals of the elements of the sequence. This shows that our theorem holds for any constant $a > 0$.

1.4 Bounded and monotone sequences

Clearly, the elements of a sequence approaching infinity cannot stay between two real numbers. We introduce the following definition.

Definition 1.8 The sequence a_n is *bounded from above*, if there is a real number K such that $a_n \leq K$ for every index n . If there is a real number K such that $a_n \geq K$ for every index n , the sequence is said to be *bounded from below*. A sequence is called bounded if it is bounded both from above and from below.

Example 1.9 Decide whether the sequence

$$a_n = \frac{2n}{\sqrt{4n^2 + 5} + 8}$$

is bounded or not? Dividing both the numerator and the denominator by $2n$ we get

$$a_n = \frac{1}{\sqrt{1 + 5/4n^2} + 8/2n},$$

hence $0 \leq a_n \leq 1$. Thus the sequence is bounded. It is also clear that the smallest upper bound of the sequence is 1, while 0 is a lower bound, but not the greatest one.

Monotone sequences have special importance.

Definition 1.10 We say that the sequence a_n is *monotone increasing*, if $a_n \leq a_{n+1}$ for every index n . A decreasing sequence is defined similarly. A sequence that is either increasing or decreasing is called monotone.

Example 1.11 Consider the sequence

$$a_n = \frac{2n - 1}{n + 2}.$$

We have

$$a_n = \frac{2n + 4 - 5}{n + 2} = 2 - \frac{5}{n + 2}.$$

The value of the fraction subtracted from 2 decreases if n increases, therefore the sequence a_n is increasing. It is also clear that the sequence is bounded from above and its smallest upper bound is 2. Moreover,

$$\lim_{n \rightarrow \infty} a_n = 2.$$

Our next theorem states that this property is characteristic for bounded monotone sequences.

Theorem 1.12 *An increasing sequence which is bounded from above is convergent.*

An analogous statement holds for decreasing sequences that are bounded from below.

We do not prove this theorem, but we note it is based on the property of real numbers that we always have a least upper bound (among the infinitely many upper bounds) which turns out to be the limit of the sequence.

Analogous theorem applies for monotone decreasing and bounded from below sequences.

1.5 Euler's number e

In many applications of mathematics the sequence

$$a_n = \left(1 + \frac{1}{n}\right)^n. \quad (1.1)$$

appears frequently. We can show that this sequence is monotone increasing, bounded from above, and consequently convergent.

To verify these properties we exploit the inequality between the arithmetic and geometric means. In particular, if x_1, \dots, x_n are positive numbers, then

$$x_1 \dots x_n \leq \left(\frac{x_1 + \dots + x_n}{n}\right)^n$$

for every integer n . Equality holds if and only if $x_1 = \dots = x_n$ that is, all the numbers are equal.

Proposition 1.13 *The sequence (1.1) is strictly monotone increasing and bounded from above.*

Proof. Let n be a given integer. Consider the $n + 1$ pieces of positive numbers

$$x_1 = 1 + \frac{1}{n}, \dots, x_n = 1 + \frac{1}{n}, x_{n+1} = 1$$

which are not all equal. Using the inequality for the arithmetic and geometric means, we have

$$\left(1 + \frac{1}{n}\right)^n < \left(\frac{n+1+1}{n+1}\right)^{n+1} = \left(1 + \frac{1}{n+1}\right)^{n+1}$$

which exactly says that the sequence is strictly monotone increasing.

Second, consider the $n + 2$ pieces of positive numbers

$$x_1 = 1 + \frac{1}{n}, \dots, x_n = 1 + \frac{1}{n}, x_{n+1} = \frac{1}{2}, x_{n+2} = \frac{1}{2}$$

which are not all equal. Using the inequality again, we have

$$\frac{1}{4} \cdot \left(1 + \frac{1}{n}\right)^n < \left(\frac{n+1+1}{n+2}\right)^{n+2} = 1$$

Rearranging the inequality we obtain $a_n < 4$, that means that the sequence is bounded from above. Consequently, the sequence (1.1) is convergent. \square

We use the notation e for the limit of this sequence. More elaborate computations show that e is irrational, and

$$e = 2.7182\dots$$

Proposition 1.14 *Let α be an arbitrarily given real number. Then*

$$\lim_{n \rightarrow \infty} \left(1 + \frac{\alpha}{n}\right)^n = e^\alpha$$

Example 1.15 Consider the sequence

$$a_n = \left(\frac{2n+1}{2n+3}\right)^n$$

Then, by rewriting the sequence we get

$$a_n = \left(\frac{2n+1}{2n+3}\right)^n = \frac{\left(1 + \frac{1/2}{n}\right)^n}{\left(1 + \frac{3/2}{n}\right)^n} \rightarrow \frac{e^{1/2}}{e^{3/2}}$$

and hence $\lim_{n \rightarrow \infty} a_n = e^{-1}$.

Study at home:

1. Careful study of Mathematical Analysis Exercises.
2. Study the exercises below.
3. Textbook-1, Chapter 1 and Section 6.4.

Chapter 2

Infinite Series

2.1 Series

Definition 2.1 Let a_k be a real infinite sequence and compose the formal sum

$$\sum_{k=1}^{\infty} a_k. \quad (2.1)$$

This symbol is called an *infinite series* (or just simply a series).

The meaning of this expression should be clarified, because only the addition of finitely many real numbers was defined so far.

For any natural number n define the n -th *partial sum* of the series (2.1) as follows:

$$S_n = \sum_{k=1}^n a_k \quad (2.2)$$

This way we created a real sequence S_n .

Definition 2.2 The infinite series (2.1) is said to be *convergent* and its sum is S , if the sequence S_n is convergent and its limit is S . In this case we use the notation:

$$S = \sum_{k=1}^{\infty} a_k$$

Otherwise the series is said to be *divergent*.

Please note that the infinite series is divergent if the sequence S_n has no limit or its limit is not finite. For instance, if $a_k = (-1)^k$ for all k then

$$S_n = \sum_{k=1}^n (-1)^k = 0 \text{ if } n \text{ is even and } S_n = \sum_{k=1}^n (-1)^k = -1 \text{ if } n \text{ is odd}$$

therefore, the sequence S_n has no limit and the series is obviously divergent.

2.2 Geometric series

Example 2.3 (Geometric series) Let r be a real number and consider the infinite geometric series with common ratio r :

$$\sum_{k=0}^{\infty} r^k$$

The n th partial sum of this series is

$$S_n = \sum_{k=0}^{n-1} r^k = \begin{cases} \frac{1-r^n}{1-r} & \text{if } r \neq 1 \\ n & \text{if } r = 1 \end{cases}$$

It is well known about the sequence $a_n = r^n$ that $r^n \rightarrow 0$ if $|r| < 1$, $r^n \rightarrow 1$ if $r = 1$ and otherwise the sequence is divergent. Therefore, we get that the geometric series is convergent if and only if $|r| < 1$ and then its sum is given by

$$S = \sum_{k=0}^{\infty} r^k = \frac{1}{1-r}$$

2.3 Convergence based on examining the partial sums

Example 2.4 Consider the series

$$\sum_{k=2}^{\infty} \frac{1}{k(k-1)}. \quad (2.3)$$

The terms of this series can be rewritten in this form:

$$\frac{1}{k(k-1)} = \frac{1}{k-1} - \frac{1}{k}$$

Observe that the n -th partial sum will be given like:

$$S_n = (1 - 1/2) + (1/2 - 1/3) + \dots + (1/(n-1) - 1/n) = 1 - 1/n$$

The limit of this sequence is obviously 1 (the negative and positive identical terms cancel each other) and we conclude that the series is convergent and its sum is $S = 1$.

Example 2.5 Try to apply the above argument for the series

$$\sum_{k=2}^{\infty} \frac{1}{k^3 - k}$$

and by eliminating the terms that cancel each other, find the sum of the series.

2.4 Conditions for convergence

Theorem 2.6 (Necessary condition for convergence) *Assume that the series*

$$\sum_{k=1}^{\infty} a_k$$

is convergent. Then $\lim_{k \rightarrow \infty} a_k \rightarrow 0$.

Example 2.7 This theorem formulates a necessary condition which may not be sufficient. For instance, we can show that the series

$$\sum_{k=1}^{\infty} \frac{1}{k}$$

fulfills the necessary condition but it is divergent. This series is called the *Harmonic series*.

Indeed, let an integer n be given, and consider the 2^n -th partial sum of the Harmonic series. Rearrange the terms in the following way

$$S_{2^n} = 1 + \frac{1}{2} + \left(\frac{1}{3} + \frac{1}{4}\right) + \left(\frac{1}{5} + \dots + \frac{1}{8}\right) + \dots + \left(\frac{1}{2^{n-1} + 1} + \dots + \frac{1}{2^n}\right),$$

where every expression within the parentheses goes up to the next power of 2. The sum of terms inside the parentheses is always bigger than $1/2$, and we have exactly n pairs of parentheses, hence

$$S_{2^n} > 1 + \frac{1}{2}n.$$

That tells us that sequence of partial sums is not bounded and therefore, the series is divergent.

Theorem 2.8 (Sufficient condition for convergence) *Let us suppose that for each index k we have $a_k \geq 0$ and the series*

$$\sum_{k=1}^{\infty} a_k$$

is convergent. If for every index k we have $0 \leq b_k \leq a_k$, then the series

$$\sum_{k=1}^{\infty} b_k$$

is also convergent.

Indeed, on the one hand the sequence of partial sums $S_n = \sum_{k=1}^n b_k$ is monotone increasing, and on the other hand it is also bounded. Consequently, the series is convergent.

In an analogous way we may formulate a sufficient condition for divergence: if all terms of a series are bigger than the nonnegative terms of a divergent series, then it is divergent as well.

Example 2.9 As an application consider the series $\sum_{k=1}^{\infty} 1/k^2$. Since for every $k > 1$

$$\frac{1}{k^2} < \frac{1}{k(k-1)}$$

then for the n -th partial sums we get

$$S_n = \sum_{k=1}^n \frac{1}{k^2} < 1 + \sum_{k=2}^n \frac{1}{k(k-1)}$$

According to the sufficient condition we conclude that this series is convergent, and for its sum we have $S < 2$.

In general, it can be verified that the series $\sum_{k=1}^{\infty} 1/k^\alpha$ is divergent, if $\alpha \leq 1$, and it is convergent if $\alpha > 1$ (see more details in Chapter 9).

2.5 Absolute convergence

In this section we examine series that may contain positive and terms as well. Consider the series

$$\sum_{k=1}^{\infty} a_k \tag{2.4}$$

where the terms a_k are not necessarily all nonnegative.

Definition 2.10 We say that the series (2.4) is *absolutely convergent*, if the series

$$\sum_{k=1}^{\infty} |a_k|$$

is convergent.

Theorem 2.11 *If a series is absolutely convergent, then it is convergent as well.*

We do not go into the details of the proof. As a justification we note the following. If S_n denotes the sum of the absolute values of the first n terms, then by our condition it is convergent and

$$\lim \sum_{k=1}^n |a_k| = \lim S_n = S.$$

Let R_n and T_n denote the sum of the negative and positive terms respectively from the first n terms of the series $\sum_{k=1}^{\infty} a_k$. Then R_n is monotone decreasing, while T_n is monotone increasing, and both sequences are bounded, since

$$R_n \geq -S \quad \text{and} \quad T_n \leq S.$$

Therefore both sequences are convergent, in notation: $\lim R_n = R$, and $\lim T_n = T$. Thus the limit of S_n can be given as:

$$\lim S_n = \lim \sum_{k=1}^n a_k = \lim(T_n + R_n) = T + R,$$

and we deduce that the series is really convergent.

The example below shows that the converse of our previous theorem is not necessarily true.

Example 2.12 Consider the following series with alternating signs:

$$\sum_{k=1}^{\infty} \frac{(-1)^{k-1}}{k}$$

Clearly, this series is not absolutely convergent, since the series with the absolute values of the terms is identical to the Harmonic series, which is divergent.

We show however, that the series above is convergent. Indeed, the sum of the terms with even indices:

$$\begin{aligned} S_{2n} &= \left(1 - \frac{1}{2}\right) + \left(\frac{1}{3} - \frac{1}{4}\right) + \dots + \left(\frac{1}{2n-1} - \frac{1}{2n}\right) = \\ &= \frac{1}{2} + \frac{1}{12} + \dots + \frac{1}{2n(2n-1)}. \end{aligned}$$

In view of Example 2.3, this sequence is monotone increasing and bounded from above, because $S_{2n} < 2$. Hence, it is convergent. Denote its limit by

$$\lim S_{2n} = S.$$

On the other hand, for the sum of the terms with odd indices we have

$$S_{2n-1} = S_{2n} + \frac{1}{2n}$$

therefore, $\lim S_{2n-1} = S$, which means that $\lim S_n = S$. This implies that the series is convergent.

2.6 Quotient-test

In this section we formulate a very useful sufficient condition for the convergence or divergence of infinite series. Create the absolute values of the quotients of the consecutive terms of the series

$$\sum_{k=1}^{\infty} a_k$$

and suppose that the limit

$$\lim_{k \rightarrow \infty} \left| \frac{a_{k+1}}{a_k} \right| = \alpha$$

exists.

Theorem 2.13 (Quotient-test)

- If $\alpha < 1$, then the series is absolutely convergent.
- If $\alpha > 1$, then the series is divergent.
- If $\alpha = 1$, then both cases can occur.

Proof. If $\alpha < 1$, then choose a real number β with $\alpha < \beta < 1$. Then from a certain index N we have

$$\left| \frac{a_{k+1}}{a_k} \right| < \beta$$

for every $k \geq N$. Then goin step-by-step backward we get

$$|a_{k+1}| < \beta |a_k| < \beta^2 |a_{k-1}| < \dots < \beta^{k-N+1} |a_N|$$

So, for the $n + 1$ -th partial sum

$$S_{n+1} = \sum_{k=0}^n |a_{k+1}| < \sum_{k=0}^{N-1} |a_{k+1}| + |a_N| \cdot \sum_{k=N}^n \beta^{k-N+1}$$

where the last sum is the partial sum of a convergent series (in view of $0 < \beta < 1$), and consequently bounded if $n \rightarrow \infty$. This proves the statement.

If $\alpha > 1$, then the the proof can be carried out similarly, with a choice of $1 < \beta < \alpha$ we can come up with an estimate with a divergent geometric series.

□

Example 2.14 In this example we demonstrate that in the case of $\alpha = 1$ nothing can be stated about the convergence of the series.

Indeed, if the divergent Harmonic series is considered, then for $a_k = 1/k$ we have

$$\frac{a_{k+1}}{a_k} = \frac{k}{k+1} \rightarrow 1 \quad \text{if } k \rightarrow \infty.$$

However, if we take the convergent series, where $a_k = 1/k^2$, then

$$\frac{a_{k+1}}{a_k} = \left(\frac{k}{k+1}\right)^2 \rightarrow 1 \quad \text{if } k \rightarrow \infty,$$

which demonstrates that both cases can occur.

Example 2.15 Find out if the series

$$\sum_{k=1}^{\infty} \frac{k^2 \cdot 2^k}{k!}$$

is convergent or not. Use the Quotient-rule:

$$\frac{a_{k+1}}{a_k} = \frac{(k+1)^2 2^{k+1}}{(k+1)!} \cdot \frac{k!}{k^2 2^k} = 2 \left(\frac{k+1}{k}\right)^2 \cdot \frac{1}{k+1} \rightarrow 0$$

Thus $\alpha = 0 < 1$, which tells us that the series is convergent.

Study at home:

1. Review the "Mathematical Analysis Exercises"
2. Additional homework: check the exercises below
3. Textbook-1, Section 6.5.

Chapter 3

Limits and continuity

3.1 Basic concepts

In the subsequent chapter we study the limit of functions $f : \mathbb{R} \rightarrow \mathbb{R}$. Let x_0 be a point (possibly equal to $\pm\infty$) for which there exists a sequence x_n in the domain of f such that $x_n \neq x_0$ and $x_n \rightarrow x_0$.

Definition 3.1 The limit of the function f at the point x_0 is said to be A (which can be $\pm\infty$) and in notation

$$\lim_{x \rightarrow x_0} f(x) = A$$

if for any sequence x_n from the domain of f whenever $x_n \rightarrow x_0$, $x_n \neq x_0$, then $f(x_n) \rightarrow A$.

ATTENTION!

Please note that the limit of f at x_0 has nothing to do with $f(x_0)$. The function may not even be defined at x_0 . However, in some cases the limit may be equal to $f(x_0)$.

Theorem 3.2 *If the functions f and g have limits at x_0 and $\lim_{x \rightarrow x_0} f(x) = A$ and $\lim_{x \rightarrow x_0} g(x) = B$ then*

- $\lim_{x \rightarrow x_0} (f \pm g)(x) = A \pm B$,
- $\lim_{x \rightarrow x_0} (f \cdot g)(x) = A \cdot B$,
- if $B \neq 0$ then $\lim_{x \rightarrow x_0} \frac{f}{g}(x) = \frac{A}{B}$,
- if $A \neq 0$ and $B = 0$ then $\lim_{x \rightarrow x_0} \frac{f}{g}(x) = \pm\infty$.

Example 3.3 Determine the limit

$$\lim_{x \rightarrow 2} \frac{x^2 - 4}{x - 2}.$$

This function is not defined for $x = 2$ but it is equal to $x + 2$ at any point $x \neq 2$. Therefore it is easily seen that

$$\lim_{x \rightarrow 2} \frac{x^2 - 4}{x - 2} = \lim_{x \rightarrow 2} (x + 2) = 4.$$

Example 3.4 Consider the function $f(x) = 1/x$. This function is not defined at $x = 0$. On the other hand, for any sequence $x_n > 0$, $x_n \rightarrow 0$ from the domain we have $f(x_n) \rightarrow +\infty$ while $f(-x_n) \rightarrow -\infty$. Thus this function has no limit at $x = 0$, that is

$$\lim_{x \rightarrow 0} \frac{1}{x}$$

does not exist.

Example 3.5 Consider the following limit:

$$\lim_{x \rightarrow +\infty} \frac{2x^4 - 5x^3 + x - 8}{8x^3 - x^2 + 12}$$

Dividing both the numerator and denominator by x^3 we get the expression

$$\frac{2x - 5 + 1/x^2 - 8/x^3}{8 - 1/x + 12/x^3}.$$

Now for any sequence $x_n \rightarrow +\infty$ the limit of the numerator is $+\infty$, while the limit of the denominator equals 8, thus the fraction tends to $+\infty$.

Very similarly, we can show that the limit of the fraction is $-\infty$, if $x \rightarrow -\infty$.

Example 3.6 Show that

$$\lim_{x \rightarrow +\infty} (\sqrt{1 + x^2} - x) = 0.$$

Indeed,

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

and the expression on the right hand side approaches 0 if $x \rightarrow +\infty$.

3.2 Squeezing theorem

In this section we formulate the Squeezing theorem for limits of functions.

Theorem 3.7 (Squeezing Theorem) *Let f , g and h be real functions such that for any x*

$$f(x) \leq g(x) \leq h(x)$$

and furthermore, $\lim_{x \rightarrow x_0} f(x) = \lim_{x \rightarrow x_0} h(x) = A$. Then the limit of the function g at x_0 exists, and

$$\lim_{x \rightarrow x_0} g(x) = A.$$

Example 3.8 Find the limit

$$\lim_{x \rightarrow 0} \frac{\sin x}{x}$$

This is an even function, therefore it is enough to consider positive values of x . A geometric interpretation (open the Figures file!) shows that for all $0 < x < \pi/2$

$$\sin x < x < \tan x.$$

Dividing this inequality by $\sin x$, we get

$$1 < \frac{x}{\sin x} < \frac{1}{\cos x}.$$

By taking the reciprocals, we obtain

$$\cos x < \frac{\sin x}{x} < 1$$

for every $0 < x < \pi/2$. In view of the Squeezing Theorem we receive

$$\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1$$

3.3 One-sided limits

In some situations the limit of a function at a given point does not exist, but we still can speak about a one-sided limit.

Definition 3.9 We say that the *right-hand limit* of f at the point x_0 exists and is equal to

$$\lim_{x \rightarrow x_0^+} f(x) = A$$

if for any sequence $x_n \rightarrow x_0$, $x_n > x_0$ from the domain of f we have $f(x_n) \rightarrow A$. The *left-hand* limit is defined analogously.

It is obvious from the definition that if at a point the limit exists, then both one-sided limits exist, and they are equal.

Example 3.10 Consider the function:

$$f(x) = \frac{2x + 1}{x - 2}$$

It is easy to see that if x_n approaches 2 from the right then $f(x_n) \rightarrow +\infty$, while if x_n approaches 2 from the left then $f(x_n) \rightarrow -\infty$. Therefore

$$\lim_{x \rightarrow 2^-} f(x) = -\infty \quad \text{and} \quad \lim_{x \rightarrow 2^+} f(x) = +\infty .$$

We can say that the limit of a function at a point exists if and only if both one-sided limits exist, and they are equal (the common value is the limit).

3.4 Continuity

Definition 3.11 Consider a function f that is defined on an interval. We say that the function f is *continuous* at a point x_0 of its domain if

$$\lim_{x \rightarrow x_0} f(x) = f(x_0) .$$

If f is not continuous at a point x_0 of its domain, then it is said that the function has a discontinuity there.

A function is simply called continuous, if it is continuous at every point of the domain.

ATTENTION!

Continuity is defined only at points in the domain of the function. For instance the function $f(x) = 1/x$ is continuous at each point of its domain, that is at each $x \neq 0$. The point $x_0 = 0$ is not in the domain of f , so we cannot speak of discontinuity here.

On the other hand, f cannot be defined at $x_0 = 0$ so that it becomes continuous, as the limit of the function does not exist there.

Functions obtained from continuous function by composition or by elementary operations (addition, subtraction, multiplication, division) are also continuous except maybe at points, where the denominator of the fraction equals zero.

Example 3.12 For instance, consider the following function:

$$f(x) = \begin{cases} \frac{1-\cos x}{x^2} & \text{if } x \neq 0 \\ \frac{1}{2} & \text{if } x = 0 \end{cases}$$

It is clear that this function is continuous for all $x \neq 0$, furthermore

$$\frac{1-\cos x}{x^2} = \frac{1-\cos^2 x}{(1+\cos x)x^2} = \left(\frac{\sin x}{x}\right)^2 \cdot \frac{1}{1+\cos x}.$$

This shows that the limit of the function at 0 equals 1/2. Thus, this function is continuous on the whole real line.

We think of a continuous function as one whose graph can be drawn by an unbroken curve (without lifting the pencil from the paper). This is expressed in Bolzano's theorem.

Theorem 3.13 (Bolzano) *Let f be a continuous function on the finite interval $[a, b]$, and suppose that $f(a)$ and $f(b)$ have different signs. Then there exists a point $c \in (a, b)$ such that $f(c) = 0$.*

We do not prove the theorem, but note that a simple idea would be bisecting the interval, and selecting the part where f has opposite signs at the endpoints. If we keep doing this infinitely many times, we receive a sequence of intervals, so that each one is the half of the preceding interval. We think that the intersection of the intervals reduces to a single point, which is necessarily a zero of the function.

Example 3.14 *Prove that the equation*

$$2x^5 - 18x^4 + 3x^3 + 20x - 13 = 0$$

has at least one real solution. The expression on the left side of the equation defines a continuous function f for which

$$\lim_{x \rightarrow +\infty} f(x) = +\infty \quad \text{and} \quad \lim_{x \rightarrow -\infty} f(x) = -\infty.$$

Therefore f is positive for sufficiently large values of x and takes negative values if x is small enough. Therefore, by the Bolzano-theorem the equation has at least one real solution.

The following property of continuous functions is of fundamental importance for extremum problems and optimization.

Theorem 3.15 (Weierstrass) *Let f be a continuous function on the finite interval $[a, b]$. Then f takes its maximum and minimum on this interval.*

We do not prove this theorem, but note that the function has to be bounded, and there exists a lowest upper bound. It can be shown that the lowest upper bound is the maximum of the function. A similar argument applies for the minimum.

For example, the function

$$f(x) = \begin{cases} x & \text{if } 0 \leq x \leq 1 \\ 3 - x & \text{if } 1 < x \leq 2 \end{cases}$$

does not reach its maximum value on the interval $[0, 2]$, but as we see, it is not continuous at 1.

Study at home:

- Textbook-1, read Sections 6.1, 6.2, 6.7, 7.1 and 7.2.
- Textbook-1, Exercises on pages 171–172, 177–178, 198, 202 and 205.
- Thorough study of "Mathematical Analysis Exercises" on my web site.

Chapter 4

Differentiation of functions

4.1 The derivative

Let f be a real function defined on an interval, and suppose that x_0 is an interior point of the interval.

Definition 4.1 We say that f is *differentiable* at x_0 if the following limit exists and it is finite:

$$\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$$

This limit is called the *derivative* of f at the point x_0 , its notation is $f'(x_0)$. We say that the function f is differentiable in an interval, if it is differentiable at every interior point of the interval.

The quotient above is called the *difference quotient* of f at the point x_0 .

Example 4.2 Consider the function $f(x) = x^2$ on the real line. The difference quotient at x_0 is:

$$\frac{f(x_0 + h) - f(x_0)}{h} = \frac{(x_0 + h)^2 - x_0^2}{h} = 2x_0 + h$$

whose limit is $2x_0$, if $h \rightarrow 0$. Consequently

$$f'(x_0) = 2x_0 .$$

In a very similar way we can show that in the case of $f(x) = x^n$ (where n is an integer),

$$f'(x_0) = nx_0^{n-1} .$$

In fact, use the identity

$$(x_0 + h)^n - x_0^n = h((x_0 + h)^{n-1} + (x_0 + h)^{n-2} \cdot x_0 + \dots + x_0^{n-1})$$

Theorem 4.3 *If f differentiable at x , then it is continuous at x .*

Proof. Let $h_n \rightarrow 0$, $h_n \neq 0$ be a sequence, then by the differentiability

$$\lim_{n \rightarrow \infty} \frac{f(x + h_n) - f(x)}{h_n} = f'(x),$$

which is finite. This is only possible if $\lim_{n \rightarrow \infty} (f(x + h_n) - f(x)) = 0$, that is $\lim_{n \rightarrow \infty} f(x + h_n) = f(x)$. This exactly means that f is continuous at x . \square

ATTENTION! The converse statement is not true in general, as it is demonstrated by the following example

Example 4.4 Consider the function $f(x) = |x|$ on the real line, and examine its difference quotient at $x_0 = 0$. It is clear that

$$\frac{f(h) - f(0)}{h} = \frac{|h|}{h} = \begin{cases} 1 & \text{if } h > 0 \\ -1 & \text{if } h < 0 \end{cases}$$

and therefore, the limit does not exist when $h \rightarrow 0$, since the right-hand limit is $+1$, while the left-hand limit is -1 . Thus the function f is not differentiable at $x = 0$

However, f is differentiable at any other point, in particular $f'(x) = 1$, if $x > 0$, and $f'(x) = -1$, if $x < 0$.

4.2 Tangent lines

Geometric interpretation (see Figures.pdf) shows that $f'(x_0)$ is the slope of the tangent line to the graph of f at x_0 .

By using this observation, we can give the equation of the tangent line to the graph of f that passes through the point $P(x_0, f(x_0))$:

$$y = f'(x_0)(x - x_0) + f(x_0).$$

For instance, the equation of the tangent line to the graph of $f(x) = x^3$ at $x_0 = 1$ is

$$y = 3(x - 1) + 1$$

Example 4.5 Find the equation of the tangent line to the graph of $f(x) = \sin x$ at $x_0 = 0$. On the one hand, the tangent line passes through the origin, on the other hand, the slope is:

$$f'(0) = \lim_{h \rightarrow 0} \frac{f(h) - f(0)}{h} = \lim_{h \rightarrow 0} \frac{\sin h}{h} = 1.$$

Therefore, the equation is $y = x$ that intersects the graph at the origin.

4.3 Rules of differentiation

Consider the functions f and g , and assume that both are differentiable at x . The rules below follow from the basic properties of limits.

Derivative of a sum If α and β real numbers, then $\alpha f(x) + \beta g(x)$ is differentiable at x and

$$(\alpha f(x) + \beta g(x))' = \alpha f'(x) + \beta g'(x),$$

Derivative of a product $f(x) \cdot g(x)$ is differentiable at x and

$$(f(x) \cdot g(x))' = f'(x) \cdot g(x) + f(x) \cdot g'(x),$$

Derivative of a quotient if $g(x) \neq 0$, then $f(x)/g(x)$ is differentiable at x , and

$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{g(x)^2}.$$

As an example let us see how we can prove the differentiability of the product:

$$\begin{aligned} \frac{f(x+h) \cdot g(x+h) - f(x) \cdot g(x)}{h} &= \\ \frac{f(x+h) \cdot g(x+h) - f(x+h) \cdot g(x)}{h} + & \\ \frac{f(x+h) \cdot g(x) - f(x) \cdot g(x)}{h} &= \\ f(x+h) \frac{g(x+h) - g(x)}{h} + g(x) \frac{f(x+h) - f(x)}{h} & \end{aligned}$$

Here the limit of the first factor is $f(x)g'(x)$ based on the continuity of f , while the limit of the second factor is $f'(x)g(x)$, if $h \rightarrow 0$. That completes the proof. The proofs of the other rules can be carried out in a very similar way.

Example 4.6 The tangent line to the graph of $f(x) = 1/x$ taken at any point encloses a triangle with the coordinate axes. (See Figures.pdf.) Show that the area of this triangle is the same, no matter at what point the tangent line is taken.

Because of the symmetry, it is enough to focus to points $x_0 > 0$. By the Quotient-rule

$$f'(x_0) = -\frac{1}{x_0^2}$$

hence, the equation of the tangent line taken at x_0 is:

$$y = -\frac{1}{x_0^2}(x - x_0) + \frac{1}{x_0}$$

The intersection points with the coordinate axes are:

if $x = 0$, then the intersection point on the y -axis is $b = 2/x_0$

and similarly

if $y = 0$, then the intersection point on the x -axis is $a = 2x_0$.

Thus, the area of the enclosed right triangle is

$$A = \frac{1}{2}ab = \frac{1}{2} \cdot 2x_0 \cdot \frac{2}{x_0} = 2$$

which is independent of the choice of x_0 .

4.4 Composition of functions

Let f and g be both $\mathbb{R} \rightarrow \mathbb{R}$ functions so that the range of g lies inside (subset) the domain of f . Then the function

$$x \rightarrow f(g(x))$$

is called the *composition* of f and g . For this function we use the notation $f \circ g$, that is:

$$f \circ g(x) = f(g(x)) .$$

For instance if $f(x) = \sqrt{x}$ and $g(x) = 1 + x^2$, then

$$f \circ g(x) = \sqrt{1 + x^2} .$$

Attention, the order is important!

In general $f \circ g \neq g \circ f$. If we consider the example above, then

$$g \circ f(x) = 1 + x$$

but this function is defined only for $x \geq 0$!

It may even turn out that $f \circ g$ is defined on the nonnegative half line, but $g \circ f$ is not defined anywhere. For instance, if

$$f(x) = -1 - x^4 \quad \text{and} \quad g(x) = \sqrt{x} ,$$

then $f \circ g(x) = -1 - x^2$, if $x \geq 0$, but $g \circ f(x) = \sqrt{-1 - x^4}$ is not defined for any real number.

4.5 Chain-Rule

Our theorem on the differentiability of composition functions is a very powerful tool for calculating the derivatives of more complicated functions.

Theorem 4.7 (Chain-Rule) *Suppose that g is differentiable at x , and f is differentiable at $g(x)$, then $f \circ g$ is differentiable at x , and its derivative is given by*

$$(f \circ g)'(x) = f'(g(x)) \cdot g'(x)$$

If we introduce the notation $k = g(x+h) - g(x)$, then the difference quotient of the composition function $f \circ g$ at x can be written like:

$$\begin{aligned} \frac{f(g(x+h)) - f(g(x))}{h} &= \\ \frac{f(g(x) + k) - f(g(x))}{k} \cdot \frac{g(x+h) - g(x)}{h} \end{aligned}$$

provided $g(x+h) - g(x) \neq 0$. In the case of $h \rightarrow 0$, in view of the continuity of g , we have $k \rightarrow 0$, and consequently, the limit of the expression on the right-hand is:

$$f'(g(x)) \cdot g'(x)$$

Unfortunately, this idea does not work when $k = 0$. In that case the proof is somewhat more complicated, we do not go into the details of that situation.

Example 4.8 For example, consider the function

$$F(x) = (1 + 3x - x^2)^6 .$$

We can find the derivative without expanding the 6-th power, if we notice that with the notations $f(x) = x^6$ and $g(x) = 1 + 3x - x^2$ we can write $F = f \circ g$. Therefore, by the Chain-Rule:

$$F'(x) = 6(1 + 3x - x^2)^5 \cdot (3 - 2x) .$$

Example 4.9 Now find the derivative of

$$F(x) = \left(\frac{2x+3}{5+x^2} \right)^3 \quad x \in \mathbb{R}$$

Then by using the notations

$$g(x) = \frac{2x+3}{5+x^2} \quad \text{and} \quad f(x) = x^3$$

we get $F = f \circ g$. Keep in mind that g is a quotient (use the Quotient-rule!), so we obtain

$$F'(x) = f'(g(x)) \cdot g'(x) = 3 \left(\frac{2x+3}{5+x^2} \right)^2 \cdot \frac{2(5+x^2) - 2x(2x+3)}{(5+x^2)^2}$$

that form can still be further simplified if we wish.

Study at home:

1. Review the "Mathematical Analysis Exercises"
2. Review the Exercises below
3. Textbook-1, Chapter 4, Sections 5.2 and 5.6.

Chapter 5

The Mean Value Theorem

5.1 The inverse function

Consider a function $f : \mathbb{R} \rightarrow \mathbb{R}$ that is one-to-one on a given interval. In the case of a continuous function this means that it is either strictly monotone decreasing or strictly monotone increasing (in view of Bolzano's theorem, see Theorem 3.13).

Definition 5.1 The *inverse* of f is the function f^{-1} whose domain is the range of f , its range is the domain of f , and further

$$f^{-1} \circ f(x) = x$$

at every point in the domain of f .

This “reverse” correspondence can be obtained by taking the equality

$$y = f(x)$$

and isolate x as the function of y :

$$x = f^{-1}(y) .$$

For instance, if $f(x) = (2x + 5)^3$, then we get

$$f^{-1}(y) = \frac{\sqrt[3]{y} - 5}{2} .$$

Geometrically this means that the graphs of f^{-1} and of f are symmetric with respect to the straight line $y = x$ (that bisects the right angle at the origin).

5.2 Differentiability of the inverse function

Theorem 5.2 *Assume that f is continuous and strictly monotone on a given interval, and it is differentiable at an interior point x . Also suppose that $f'(x) \neq 0$. Then f^{-1} is differentiable at $y = f(x)$, and*

$$(f^{-1})'(y) = \frac{1}{f'(x)}.$$

Roughly, the situation is the following. Consider the difference quotient:

$$\frac{f^{-1}(y+h) - f^{-1}(y)}{h}$$

Let x and $x+k$ be points in the domain of f such that $y = f(x)$ and $y+h = f(x+k)$. Then the difference quotient can be written in the following form:

$$\frac{x+k-x}{f(x+k)-f(x)} = \frac{1}{\frac{f(x+k)-f(x)}{k}}$$

If here $h \rightarrow 0$, then $k \rightarrow 0$ (ATTENTION, this is not trivial! It means the continuity of f^{-1}), and hence, the limit of the fraction on the right-hand side is really $1/f'(x)$.

Example 5.3 Find the derivative of the function

$$g(x) = \sqrt[n]{x}$$

at a point $x > 0$. As we see, g is the inverse of the power function $f(x) = x^n$ on the non-negative half line, that is $g(y) = f^{-1}(y)$. Thus,

$$g'(y) = \frac{1}{f'(x)} = \frac{1}{nx^{n-1}} = \frac{1}{n} \cdot y^{\frac{1}{n}-1}$$

since $y = x^n$ and consequently

$$x^{n-1} = y^{\frac{n-1}{n}}$$

In view of this example we conclude that for every rational exponent r the function $F(x) = x^r$ is differentiable at every point $x > 0$, and its derivative is:

$$F'(x) = rx^{r-1}.$$

Example 5.4 Calculate the derivative of the function

$$F(x) = \sqrt{1+x^4}$$

Set $f(x) = \sqrt{x}$ and $g(x) = 1+x^4$, with these notations we have $F = f \circ g$. Making use of the Chain-Rule we get

$$F'(x) = f'(g(x)) \cdot g'(x) = \frac{4x^3}{2\sqrt{1+x^4}}$$

5.3 The exponential and logarithm functions

Consider the exponential function with base e on the real line, and its inverse, which is the logarithm function with base e (that is denoted by the symbol \ln):

$$f(x) = e^x \quad f^{-1}(x) = \ln x \quad (x > 0).$$

They are called the *natural* exponential function, and the *natural* logarithm function, respectively. Below we find their derivatives. We start with the equality

$$\lim_{x \rightarrow \pm\infty} \left(1 + \frac{1}{x}\right)^x = e.$$

Find the derivative of the natural logarithm function at $x_0 = 1$.

$$\frac{\ln(1+h) - \ln 1}{h} = \ln(1+h)^{1/h}$$

whose common right-hand limit and left-hand limit at zero is $\ln e$. (Here we supposed the continuity of the logarithm function.) Therefore, the derivative is 1.

The derivative of $f(x) = e^x$ at the point 0 can be determined by exploiting our theorem about the differentiability of the inverse function:

$$f'(0) = \lim_{h \rightarrow 0} \frac{e^h - 1}{h} = \frac{1}{(\ln)'(1)} = 1.$$

This enables us to get the derivative of the exponential function at an arbitrary point x :

$$f'(x) = \lim_{h \rightarrow 0} \frac{e^{x+h} - e^x}{h} = e^x \cdot \lim_{h \rightarrow 0} \frac{e^h - 1}{h} = e^x$$

Using the differentiability of the inverse again, we obtain the derivative of the logarithm function at any given point $x > 0$:

$$(f^{-1})'(x) = \frac{1}{e^{\ln x}} = \frac{1}{x}.$$

Example 5.5 As a straightforward application, find the derivative of the function

$$f(x) = x^\alpha$$

at any given point $x > 0$, where α is an arbitrary real exponent. First we write:

$$f(x) = x^\alpha = e^{\alpha \ln x}$$

Then, in view of the Chain-Rule we get

$$f'(x) = \alpha \frac{1}{x} e^{\alpha \ln x} = \alpha \frac{1}{x} x^\alpha = \alpha x^{\alpha-1}$$

This tells us that the differentiation can be carried out the same way as in the case of rational exponents.

5.4 Necessary condition for an extremum

Consider a function $f : \mathbb{R} \rightarrow \mathbb{R}$.

Definition 5.6 We say that a point x_0 in the domain of f is a (global) minimum point, if $f(x_0) \leq f(x)$ for every point $x \neq x_0$ in the domain of f .

We say that a point x_0 in the domain of f is a local minimum point, if there exists a positive number $\varepsilon > 0$ such that $f(x_0) \leq f(x)$ at every point in the domain x with $0 < |x - x_0| < \varepsilon$.

In both cases we strict minimum points if strict inequalities apply.

We can formulate analogous definitions for maximum points.

It is obvious that a global minimum point is also a local minimum point. The converse statement however, is not true in general, as it is shown in the following example. For instance, the function

$$f(x) = \begin{cases} (x+1)^2 & \text{ha } x < 0 \\ (x-1)^2 & \text{ha } x \geq 0 \end{cases}$$

admits a local maximum at $x = 0$ (here the function is continuous, but not differentiable, check it!) but this function does not have a global maximum, since it is not bounded from above.

For differentiable functions we can present the following characterization of local extreme (minimum or maximum) points.

Theorem 5.7 *Let us suppose that f is defined on an interval, at it is differentiable at an interior point x_0 . If x_0 is a local minimum point of f , then $f'(x_0) = 0$.*

Proof. Indeed, consider the difference quotient:

$$\frac{f(x_0 + h) - f(x_0)}{h}.$$

If $h > 0$, then the difference quotient for small values of h is non-negative, and consequently, the right-hand limit is non-negative. On the other hand, if $h < 0$, then similarly, the left-hand limit is non-positive. By the differentiability assumption the difference quotient has a limit when $h \rightarrow 0$, which therefore, can only be zero. Thus $f'(x_0) = 0$.

This theorem formulates only a necessary condition for minimum, which is not sufficient! For example, the function $f(x) = x^3$ has no extreme point at $x = 0$, but $f'(0) = 0$.

In the case of a differentiable function, those points x_0 where $f'(x_0) = 0$, are called *critical* (or sometimes stationary) points. Using this vocabulary, we may say that the extreme points of a function are critical, the converse statement is not necessarily true.

5.5 Lagrange's Mean Value Theorem

Based on the geometric interpretation, the Mean Value Theorem formulates a very illustrative statement.

Theorem 5.8 *Let f be continuous on the finite closed interval $[a, b]$, and differentiable in the interior of the interval. Then there exists a point $\xi \in (a, b)$ so that*

$$f'(\xi) = \frac{f(b) - f(a)}{b - a}$$

Proof. Introduce the function

$$g(x) = f(x) - f(a) - \frac{f(b) - f(a)}{b - a}(x - a)$$

According to the assumptions, this function is continuous on the interval $[a, b]$, hence, by Weierstrass' Theorem (see Theorem 3.15) it achieves its minimum and maximum in $[a, b]$ interval. At least one of the extreme points (either the minimum, or the maximum) is in the interior of the interval, because

$$g(a) = g(b) = 0.$$

If this interior extreme point is $\xi \in (a, b)$, then by our previous theorem $g'(\xi) = 0$. This exactly means that

$$f'(\xi) - \frac{f(b) - f(a)}{b - a} = 0.$$

Please observe, that the continuity assumption in our theorem is vital! Sketch a figure to show that!

5.6 L'Hôpital's Rule

The procedure below makes it possible to compute complicated limits relatively easily.

Let both f and g be differentiable, and their derivatives f' and g' are continuous in a neighborhood of a point x_0 , and suppose that $f(x_0) = g(x_0) = 0$. We want to find the limit

$$\lim_{x \rightarrow x_0} \frac{f(x)}{g(x)}$$

which is of the form $0/0$ so "undefined".

By the Mean Value Theorem

$$\frac{f(x)}{g(x)} = \frac{\frac{f(x)-f(x_0)}{x-x_0}}{\frac{g(x)-g(x_0)}{x-x_0}} = \frac{f'(\xi)}{g'(\eta)}$$

where ξ and η are points between x and x_0 . Now if $x \rightarrow x_0$, then both $\xi \rightarrow x_0$ and $\eta \rightarrow x_0$. Therefore, by the continuity of the derivative functions we get

$$\lim_{x \rightarrow x_0} \frac{f(x)}{g(x)} = \lim_{x \rightarrow x_0} \frac{f'(x)}{g'(x)}$$

This equality is called L'Hôpital's Rule. If the resulting limit still has the form $0/0$, then apply L'Hôpital's Rule again until a "decent" limit is received.

Example 5.9 Find the following limit by using L'Hôpital's Rule:

$$\lim_{x \rightarrow 0} \frac{2 \sin x}{1 - \sqrt{1+x}}$$

Taking the limits of the derivatives, we have:

$$\lim_{x \rightarrow 0} \frac{2 \sin x}{1 - \sqrt{1+x}} = \frac{2 \cos 0}{-\frac{1}{2\sqrt{1+0}}} = -4$$

Study at home:

1. Review of the exercises in "Mathematical Analysis Exercises"
2. Textbook-1, Sections 5.1, 5.4, 7.5 and 7.6, Chapter 8.

Chapter 6

Complete analysis of functions

6.1 Monotone functions

Definition 6.1 We say that f is monotone increasing on an interval, if for any two points of the interval with $x_1 < x_2$ we have $f(x_1) \leq f(x_2)$. An analogous definition applies for monotone decreasing functions.

We say that the function is strictly monotone (in either case), if we have strict inequalities in the definition.

Theorem 6.2 *Let f be continuous on a finite closed interval $[a, b]$, and differentiable in its interior. If we have $f'(x) > 0$ at every interior point of the interval, then f is strictly monotone increasing on $[a, b]$.*

Indeed, if $x_1 < x_2$ are two arbitrary points of the interval $[a, b]$, then by the Lagrange's Mean Value Theorem there exists a point $x_1 < \xi < x_2$, such that

$$f(x_2) - f(x_1) = f'(\xi)(x_2 - x_1) .$$

By our assumption the right-hand side is positive, therefore

$$f(x_2) - f(x_1) > 0$$

that means f is strictly monotone increasing on the interval.

Now, let us examine a function that is monotone increasing and differentiable in an interval. For any two different points x and $x + h$ in the interval we have:

$$\frac{f(x+h) - f(x)}{h} \geq 0$$

regardless of $h > 0$ or $h < 0$. Passing to the limit $h \rightarrow 0$ we obtain $f'(x) \geq 0$. Thus, we can formulate the following theorem.

Theorem 6.3 *Let f be continuous on the interval $[a, b]$, and differentiable in its interior. Then f is monotone increasing on the interval if and only if $f'(x) \geq 0$ at each interior point of the interval.*

A completely similar statement can be formulated for monotone decreasing functions.

However, the assertion that if f is strictly monotone increasing, then we would have $f'(x) > 0$ for every interior point x is NOT TRUE. For example, the function $f(x) = x^3$ is strictly monotone increasing on the entire real line, but $f'(0) = 0$.

6.2 Finding extreme points

Consider a function $f : \mathbb{R} \rightarrow \mathbb{R}$ and pick an interior point x_0 in the domain. Suppose that f is differentiable at x_0 .

As we have seen, the necessary condition for x_0 for being an extreme point is $f'(x_0) = 0$. The question is, how we can formulate a sufficient condition for really having an extremum at x_0 . It is easy to see that if there exists a positive number $\varepsilon > 0$ so that f is monotone decreasing on the interval $[x_0 - \varepsilon, x_0]$, moreover f is monotone increasing on the interval $[x_0, x_0 + \varepsilon]$, then x_0 is definitely a local minimum point of f .

For differentiable functions we can summarize this observation in the following theorem.

Theorem 6.4 *Assume that f is differentiable in an interval, and x_0 is an interior point of the interval. If there exists a positive number $\varepsilon > 0$, so that*

- $f'(x) \leq 0$, if $x \in (x_0 - \varepsilon, x_0)$
- $f'(x) \geq 0$, if $x \in (x_0, x_0 + \varepsilon)$

then x_0 is a local minimum point of f .

Obviously, an analogous statement can be formulated for the case of local maximum as well.

Example 6.5 Find the extreme points and the intervals of monotonicity of the function

$$f(x) = x^2 e^{-x}$$

By the Product-Rule, the derivative is:

$$f'(x) = (2x - x^2)e^{-x}$$

whose sign depends exclusively on the first factor (the second is positive). Consequently:

- If $x \in (-\infty, 0)$, then $f'(x) < 0$, so f is monotone decreasing.
- If $x = 0$, then $f'(0) = 0$, this is a critical point.
- If $x \in (0, 2)$, then $f'(x) > 0$, so f is monotone increasing.
- If $x = 2$, then $f'(2) = 0$, this is another critical point.
- If $x \in (2, +\infty)$, then $f'(x) < 0$, so f is monotone decreasing.

By the changing the signs of f' we can conclude that $x = 0$ is a minimum point (global), while $x = 2$ is a local maximum point.

Example 6.6 Consider the following function on the real line:

$$f(x) = x + \sin x$$

Since $f'(x) = 1 + \cos x$, it is clear that function has critical points at

$$x = (2k + 1)\pi \quad k = 0, \pm 1, \pm 2, \dots$$

However, none of them is an extremum:

$$x \neq (2k + 1)\pi \quad \text{then} \quad f'(x) > 0,$$

because $\cos x > -1$. This means that the derivative does not change its sign. In fact, this function is strictly monotone increasing on the entire real line.

6.3 Higher order derivatives

If a function f is differentiable in a given interval, then the correspondence $x \rightarrow f'(x)$ is called the derivative function of f . If f' is again differentiable at a given point x_0 , then we say that f is twice differentiable at this point. Instead of using the complicated notation $(f')'(x_0)$, we use the brief formula

$$f''(x_0)$$

and this is called the second derivative of f at x_0 .

In a completely similar way, if n is a given integer, we can define the n -th derivative of the function f at x_0 , and its notation is

$$f^{(n)}(x_0).$$

For instance, for the function $f(x) = 1/x$ at any given point $x_0 \neq 0$ we have

$$f''(x_0) = \frac{2}{x_0^3} \quad \text{and} \quad f^{(n)}(x_0) = \frac{(-1)^n n!}{x_0^{n+1}}$$

for every integer n .

Example 6.7 Consider the function $f(x) = \sin x$, and find its derivative function.

$$\begin{aligned} f'(x) &= \lim_{h \rightarrow 0} \frac{\sin(x+h) - \sin x}{h} = \lim_{h \rightarrow 0} \frac{\sin x \cos h + \cos x \sin h - \sin x}{h} \\ &= \sin x \cdot \lim_{h \rightarrow 0} \frac{\cos h - 1}{h} + \cos x \cdot \lim_{h \rightarrow 0} \frac{\sin h}{h} \end{aligned}$$

In view of Example 3.12 the first limit is 0, and in view of Example 3.8 the second limit is 1. Therefore,

$$f'(x) = \cos x$$

By using the identity $\cos x = \sin(x + \pi/2)$ and the Chain-Rule, we have

$$(\cos x)' = \cos(x + \pi/2) = -\sin x$$

Therefore, the higher order derivatives of $f(x) = \sin x$ can be given in terms of the divisibility by 4:

$$f^{(n)}(x) = \begin{cases} \cos x & \text{if } n = 4k + 1 \\ -\sin x & \text{if } n = 4k + 2 \\ -\cos x & \text{if } n = 4k + 3 \\ \sin x & \text{if } n \text{ is divisible by } 4 \end{cases}$$

6.4 Second order conditions

It may happen that we analyze a function, where the sign of its derivative is not easy to determine (for instance a higher degree polynomial). In a case like that, the second order (sufficient) condition proves to be useful.

Theorem 6.8 *Let f be differentiable in an interval, and suppose that f is twice differentiable at an interior point x_0 .*

If $f'(x_0) = 0$ and $f''(x_0) > 0$, then x_0 is a local minimum point of f .

Proof. Indeed, by examining the different quotient we get

$$\begin{aligned} f''(x_0) &= \lim_{h \rightarrow 0} \frac{f'(x_0 + h) - f'(x_0)}{h} = \\ &= \lim_{h \rightarrow 0} \frac{f'(x_0 + h)}{h} > 0 \end{aligned}$$

This means that the quotient $f'(x_0 + h)/h$ is positive for $0 < |h| < \varepsilon$ for some $\varepsilon > 0$. This implies that

- if $x \in (x_0 - \varepsilon, x_0)$, then $f'(x) < 0$,
- if $x \in (x_0, x_0 + \varepsilon)$, then $f'(x) > 0$.

Making use of Theorem 6.4 we conclude that x_0 is really a local minimum point.

We can formulate an analogous second order sufficient condition for the case of local maximum.

By using proof by contradiction, we get the second order necessary condition for an extremum point.

Theorem 6.9 *Assume that f is twice differentiable in an interval, and let x_0 be an interior point of the interval.*

- If x_0 is a local minimum point, then $f'(x_0) = 0$, and $f''(x_0) \geq 0$.
- If x_0 is a local maximum point, then $f'(x_0) = 0$, and $f''(x_0) \leq 0$.

Example 6.10 For $x > 0$ consider the function

$$f(x) = x \ln x$$

Then $f'(x) = 1 + \ln x$, therefore, the only critical point of f is $x = 1/e$. On the other hand $f''(x) = 1/x$, so we have

$$f''(1/e) = e > 0,$$

Thus, $x = 1/e$ is a local minimum point of f . (It is not hard to verify that this is a global minimum point as well.)

Please observe that our theorems provide no information for a critical point x_0 with

$$f''(x_0) = 0.$$

The reason that in this “marginal” situation anything can happen. For example, examine the behavior of the functions

$$f(x) = x^n \quad (n \geq 3)$$

at the critical point $x_0 = 0$. On the one hand, here $f'(0) = 0$ and $f''(0) = 0$. On the other hand

- if n is even, then $x_0 = 0$ is (global) minimum point,
- if n is odd, then $x_0 = 0$ is not an extremum point (so-called saddle point).

Very similarly, if n is even, then $x_0 = 0$ is a (global) maximum point of $-f$.

6.5 Convex and concave functions

Definition 6.11 The function f is said to be convex on the interval $[a, b]$, if for any two points x_1 and x_2 from the interval, and for any real number $0 \leq \alpha \leq 1$

$$f(\alpha x_1 + (1 - \alpha)x_2) \leq \alpha f(x_1) + (1 - \alpha)f(x_2) .$$

The geometric meaning of this definition is that any cord to the graph (i.e. a segment that connects two points on the graph) can nowhere be below the graph of the function.

Concave functions are defined by the opposite inequality.

We now give a simple characterization of convexity for twice differentiable functions.

Theorem 6.12 *Assume that f is continuous on a closed interval, and twice differentiable in the interior. The necessary and sufficient condition for the convexity of f is:*

$$f''(x) \geq 0$$

at every interior point of the interval.

In particular, this means that for convex functions the slope of the tangent line (i.e. the derivative) is monotone increasing. Geometrically this can be illustrated by the fact that the graph of the function is nowhere below the tangent line.

Example 6.13 Give a complete analysis of the function

$$f(x) = \frac{x}{1 + x^2}$$

First we calculate the derivative:

$$f'(x) = \frac{1 - x^2}{(1 + x^2)^2} .$$

By examining the sign of the derivative, we come up with the following summary:

- f is strictly monotone decreasing on the interval $(-\infty, -1)$
- $x = -1$ is a (global) minimum point
- f is strictly monotone increasing on the interval $(-1, 1)$
- $x = 1$ is a (global) maximum point
- f is strictly monotone decreasing on the interval $(1, +\infty)$.

The convexity is investigated by specifying the sign of the second derivative:

$$f''(x) = \frac{2x^3 - 6x}{(1 + x^2)^3}$$

Obviously, the denominator is positive, so it is enough to find the sign of the numerator:

$$2x^3 - 6x = 2x(x^2 - 3)$$

By examining the factors we come up with the following summary:

- f is concave on the interval $(-\infty, -\sqrt{3})$
- f is convex on the interval $(-\sqrt{3}, 0)$
- f is concave on the interval $(0, \sqrt{3})$
- f is convex on the interval $(\sqrt{3}, +\infty)$.

Please notice that we have $f''(-\sqrt{3}) = f''(0) = f''(\sqrt{3}) = 0$, and the second derivative changes the sign at those points. In other words those points separate the convex and concave segments of the function. Such points are called the *points of inflection* of f . At a point of inflection the tangent line intersects the graph of the function.

Probably the most important property of convex function is that every local minimum point is a global minimum point as well.

Theorem 6.14 *Consider a twice differentiable convex function f on an interval, and let x_0 be an interior point of the interval. If x_0 is a local minimum point, then it is a global minimum point.*

Proof. Indeed, on the one hand $f'(x_0) = 0$, on the other hand f' is monotone increasing. Therefore, at every interior point x_0 :

- if $x < x_0$, then $f'(x) \leq 0$, and hence, $f(x) \geq f(x_0)$,
- if $x > x_0$, then $f'(x) \geq 0$, and hence, $f(x) \geq f(x_0)$.

This proves our statement.

A completely analogous theorem can be formulated for concave functions and maximum points.

Example 6.15 Define the function f for $x > 0$ on the positive half line:

$$f(x) = ax + 2 \ln x$$

where a is an unspecified parameter. For what value of a will f possess a global maximum point at $x = 6$?

By the necessary condition for an extremum

$$f'(x) = a + \frac{2}{x} = 0$$

that yields $x = -2/a$. By the condition $x = 6$, we get $a = -1/3$. The second derivative of f is:

$$f''(x) = -\frac{1}{x^2} < 0,$$

therefore the function is concave on the whole domain. Consequently, for the parameter $a = -1/3$ the function f has a global maximum point at $x = 6$.

Study at home:

1. Careful review of the "Mathematical Analysis Exercises"
2. Textbook-1: Chapter 9.

Chapter 7

Integration

7.1 The indefinite integral

Definition 7.1 Let f be a function defined on an interval I . A differentiable function F defined on I is called the *indefinite integral* of f , or sometimes its *primitive function*, if

$$F'(x) = f(x)$$

for every $x \in I$.

It is clear that taking the indefinite integral is the reverse operation of differentiation. It is important to note that the indefinite integral is not unique! Indeed, if F is the indefinite integral of a function f , then by adding a constant C to F we again have an indefinite integral:

$$(F(x) + C)' = F'(x) = f(x)$$

for every $x \in I$.

We show that this is the only way to create other indefinite integrals.

Theorem 7.2 *If F is an indefinite integral of f on the interval I , then any indefinite integral of f can be given in the form $F + C$, where C is a constant.*

Proof. Indeed, if the differentiable function G is an indefinite integral of f on the interval I , then at every point $x \in I$ we have

$$(F(x) - G(x))' = f(x) - f(x) = 0$$

This means that the derivative of $F - G$ is zero on I . By the Mean Value Theorem we get that $F - G$ is constant on the interval.

In view of our theorem, we use the following notation for indefinite integrals:

$$\int f(x) dx = F(x) + C$$

For instance, by simple differentiation we can verify

$$\int \cos x dx = \sin x + C$$

or very similarly

$$\int x^\alpha dx = \frac{x^{\alpha+1}}{\alpha+1} + C \quad (\alpha \neq -1)$$

where C is an arbitrary constant. If a function has an indefinite integral on an interval, then there are infinitely many of them.

7.2 Basic integrals

The following rule can be useful for finding indefinite integrals:

Theorem 7.3 $\int (\alpha f(x) + \beta g(x)) dx = \alpha \int f(x) dx + \beta \int g(x) dx$

This rule can be extended to any sums with finitely many terms.

ATTENTION: Not all functions have indefinite integrals. For example, a function f with a point of discontinuity, where the one-sided limits exist, they are finite, but not equal, cannot possess an indefinite integral. The following theorem formulates a useful sufficient condition for the existence of the indefinite integral.

Theorem 7.4 *If f is continuous on the interval I , then it has an indefinite integral.*

We can easily create rules for finding indefinite integrals by reversing the differentiation rules. By taking the opposites of differentiation rules for elementary functions, we obtain rules for finding indefinite integrals.

In general, any formula for an indefinite integral can be verified by direct differentiation. For example:

$$\begin{aligned} \int \sin x dx &= -\cos x + C \\ \int (2x^2 - 5x + 8) dx &= \frac{2}{3}x^3 - \frac{5}{2}x^2 + 8x + C \\ \int e^{2x-1} dx &= \frac{1}{2}e^{2x-1} + C \\ \int \frac{2x}{1+x^2} dx &= \ln(1+x^2) + C \end{aligned}$$

7.3 Initial value problems

As we have seen, a function can have infinitely many indefinite integrals (if any exists), and they differ only in an additive constant. However, if fix a point in the coordinate system, and looking only for a definite integral that passes through the given point, then the solution of the problem may be unique.

Example 7.5 Find the function F for which

$$F'(x) = 2e^{-x} \quad \text{and} \quad F(0) = 1$$

In this case we are looking for a specific indefinite integral

$$F(x) = 2 \int e^{-x} dx = -2e^{-x} + C$$

so that $F(0) = 1$. The condition implies $C = 3$, and this is the only solution.

7.4 Definite integrals

In this section we briefly outline how Bernhard Riemann, professor of mathematics at University of Göttingen (Germany) introduced the concept of integration in the 19-th century. The idea is based on the two-sided approximation developed Archimedes, the ancient greek mathematician. This idea is a fundamental element of human thinking, and this is how Archimedes determined the area of the circle in Syracuse, using the areas of approximating polygons from inside and outside.

Let f be a continuous function on the finite interval $[a, b]$, and consider the partition of the interval into n subintervals by using the points

$$a = x_0 < x_1 < \dots < x_n = b$$

On every subinterval $[x_{k-1}, x_k]$ let m_k denote minimum value of f , and let M_k denote the maximum value of f . Those extreme values exist by virtue of Weierstrass' theorem (see Theorem 3.15). Create the sum

$$s_n = \sum_{k=1}^n m_k (x_k - x_{k-1})$$

that we call *lower sum*, and the sum

$$S_n = \sum_{k=1}^n M_k (x_k - x_{k-1})$$

that we call *upper sum*. The sum the areas of these rectangles approximate the area below the graph of f from below, and from above, respectively. Check Figures.pdf for details!

We can easily see that by inserting a new node point s_n cannot decrease, and S_n cannot increase. It can be shown that if the density of the partion gets higher then the lowest upper bound of the lower sums coincides with the highest lower bound of the upper sums. Following Riemann's idea, this common value S is called the definite integral of f on the interval $[a, b]$. The notation is:

$$S = \int_a^b f(x) dx$$

which means the (signed!) area below the graph of f .

ATTENTION!

The area above the x -axis comes with positive sign, the area below the x -axis comes with negative sign, respectively.

Based on this geometric interpretation, the following properties of the definite integral are intuitively obvious.

Theorem 7.6 *Let f and g be functions that have definite integrals on $[a, b]$. Then*

1. *if $f(x) \leq g(x)$ on the interval $[a, b]$, then*

$$\int_a^b f(x) dx \leq \int_a^b g(x) dx$$

2. *in particular, $|\int_a^b f(x) dx| \leq \int_a^b |f(x)| dx$.*

3. *If $f(x) \leq M$ on the interval $[a, b]$ (M is a constant), then*

$$\int_a^b f(x) dx \leq M(b - a)$$

4. *If f is continuous on the interval $[a, b]$, then there exists a point $\bar{x} \in [a, b]$, for which $\int_a^b f(x) dx = f(\bar{x})(b - a)$.*

5. *By definition: $\int_b^a f(x) dx = -\int_a^b f(x) dx$ if $a \leq b$.*

6. $\int_a^c f(x) dx = \int_a^b f(x) dx + \int_b^c f(x) dx$

Create a picture, and interpret the above statements geometrically!

7.5 Newton-Leibniz-formula

In this section we show how a definite integral can be evaluated by using the indefinite integral (primitive function). Our main result is sometimes called the "Fundamental Theorem of Calculus" (in the English literature).

Theorem 7.7 (Newton-Leibniz-formula) *If F is a primitive function of the continuous function f on the finite interval $[a, b]$, then*

$$\int_a^b f(x) dx = F(b) - F(a)$$

Justification (not a proof!): It is easy to see that our statement is independent of the choice of the indefinite integral. Indeed, if G is another primitive function of f , then

$$G(x) = F(x) + C$$

on $[a, b]$ with a constant C (see Theorem 7.2), and therefore,

$$\int_a^b f(x) dx = [G(x)]_a^b = G(b) - G(a) = (F(b) + C) - (F(a) + C) = F(b) - F(a).$$

On the other hand, fix a point $x \in [a, b]$ and consider the integral

$$F(x) = \int_a^x f(t) dt$$

Then $F(a) = 0$, since the length of the path of integration is zero. It would be enough to show that this F is an indefinite integral of f .

In view of Theorem 7.6 for any $a < x < b$ and $h \neq 0$ with $x + h \in [a, b]$, there exists a point \bar{x} between x and $x + h$ with the following property:

$$\frac{1}{h}(F(x+h) - F(x)) = \frac{1}{h} \int_x^{x+h} f(t) dt = \frac{1}{h} f(\bar{x}) \cdot h$$

Now, if we pass to the limit $h \rightarrow 0$, then $\bar{x} \rightarrow x$, and by the continuity of f we also have $f(\bar{x}) \rightarrow f(x)$ that is

$$\lim_{h \rightarrow 0} \frac{1}{h}(F(x+h) - F(x)) = F'(x) = \lim_{h \rightarrow 0} f(\bar{x}) = f(x)$$

This means that F is really a primitive function of f . □

It was an amazing achievement by Newton and Leibniz, and the mathematics of their time, to find the beautiful relationship between the derivative and the geometry of definite integrals, as it is described in our theorem.

This discovery is so fundamental that it cannot be overestimated. First, it triggered a very rapid development in physics and chemistry, and somewhat later it gave a massive boost to the evolution of sciences like biology, economics and others. Summing up, we may say today that the theory of differentiation and integration provides the precise scientific language and vocabulary in all branches of sciences.

For convenience, sometimes we use the following notation:

$$\int_a^b f(x) dx = [F(x)]_a^b = F(b) - F(a)$$

As a consequence of the Newton-Leibniz-formula, we can formulate the following statement.

Consequence 7.8 *If f is continuous on an interval, then it has a primitive function on that interval.*

Proof. In view of the proof of the Newton-Leibniz-formula, we get that the function

$$F(x) = \int_a^x f(t) dt$$

is really a primitive function of f on the given interval. \square

Example 7.9 Evaluate the definite integral below.

$$\int_1^2 \left(2x^3 + 1 + \frac{1}{x^2} \right) dx = \left[\frac{x^4}{2} + x - \frac{1}{x} \right]_1^2 = 9$$

Some more examples:

$$\int_0^{\pi/2} \sin x dx = [-\cos x]_0^{\pi/2} = 1$$

$$\int_0^1 e^x dx = [e^x]_0^1 = e - 1$$

$$\int_0^4 \sqrt{x} dx = \left[\frac{2}{3} \cdot x^{3/2} \right]_0^4 = \frac{16}{3}$$

Study at home:

1. Careful review of "Mathematical Analysis Exercises"
2. Textbook-1, Chapter 10.

Chapter 8

Methods of integration

8.1 Integration by parts

If f and g are continuously differentiable functions on an interval I , then by the Product-Rule we have:

$$\int f'(x)g(x) dx = f(x)g(x) - \int f(x)g'(x) dx$$

This formula is called *integration by parts*. For example, consider the integral

$$\int xe^{-x} dx$$

then by using the allocation $f'(x) = e^{-x}$ and $g(x) = x$ (could we do it the other way?):

$$\int xe^{-x} dx = -xe^{-x} + \int e^{-x} dx = -xe^{-x} - e^{-x} + C$$

Example 8.1 Use integration by parts in the integral

$$\int x^n \ln x dx$$

(where $n \neq -1$). Introduce the notation $f'(x) = x^n$ and $g(x) = \ln x$, then (what do we get in the opposite way?)

$$\int x^n \ln x dx = \frac{x^{n+1}}{n+1} \ln x - \int \frac{x^n}{n+1} dx = \frac{x^{n+1}}{n+1} \ln x - \frac{x^{n+1}}{(n+1)^2} + C$$

In particular, for $n = 0$ we have:

$$\int \ln x dx = x \ln x - x + C = x(\ln x - 1) + C$$

8.2 Integration by parts in definite integrals

We can use integration by parts in definite integrals in the following way:

$$\int_a^b f'(x)g(x) dx = [f(x)g(x)]_a^b - \int_a^b f(x)g'(x) dx$$

For instance, by setting $f'(x) = \sin x$ and $g(x) = x$ (would the opposite way successful?):

$$\begin{aligned} \int_0^\pi x \sin x dx &= [-x \cos x]_0^\pi + \int_0^\pi \cos x dx \\ &= \pi + [\sin x]_0^\pi = \pi \end{aligned}$$

This procedure is faster than first computing the indefinite integral and then substituting the bounds. Further, it may minimize the chance of miscalculation.

Example 8.2 Sometimes we need to carry out integration by parts more times in a row. Consider the integral

$$\int x^2 e^{-\lambda x} dx$$

where $\lambda > 0$ is a given parameter. Introduce the notations $f'(x) = e^{-\lambda x}$, and $g(x) = x^2$, then

$$\int x^2 e^{-\lambda x} dx = -\frac{1}{\lambda} x^2 e^{-\lambda x} + \frac{1}{\lambda} \int 2x e^{-\lambda x} dx$$

The last integral can be evaluated by a repeated integration by parts.

Attention! We stick to the notations $f'(x) = e^{-\lambda x}$ and $g(x) = x$. In the opposite situation we come to an absolutely useless identity. Give it a try!

$$\int x^2 e^{-\lambda x} dx = -\frac{1}{\lambda} x^2 e^{-\lambda x} - \frac{2}{\lambda^2} x e^{-\lambda x} - \frac{2}{\lambda^3} e^{-\lambda x} + C$$

Example 8.3 Find the definite integral below:

$$\int_0^\pi e^x \sin x dx$$

Apply the setting $f'(x) = e^x$ and $g(x) = \sin x$, then by two consecutive integrations by parts:

$$\begin{aligned} \int_0^\pi e^x \sin x dx &= [e^x \sin x]_0^\pi - \int_0^\pi e^x \cos x dx \\ &= -[e^x \cos x]_0^\pi - \int_0^\pi e^x \sin x dx \end{aligned}$$

Isolate the original integral on the left-hand side:

$$2 \int_0^{\pi} e^x \sin x \, dx = -[e^x \cos x]_0^{\pi}$$

which means

$$\int_0^{\pi} e^x \sin x \, dx = \frac{1}{2}(e^{\pi} + 1)$$

8.3 Integration by substitution

From the differentiation of a composition of functions (i.e. the Chain-Rule) we derive the following identity:

$$\int f(g(t))g'(t) \, dt = \int f(x) \, dx$$

where $x = g(t)$ is a continuously differentiable function on an interval. This formula is called the *integration by substitution*.

Example 8.4 Calculate the following indefinite integral:

$$\int 5t^3 \sqrt{2+t^4} \, dt$$

Observe that by introducing the substitution $x = g(t) = t^4$, the integral can be rewritten in this form:

$$\int 5t^3 \sqrt{2+t^4} \, dt = \frac{5}{4} \int \sqrt{2+x} \, dx = \frac{5}{4} \cdot \frac{2}{3} (2+x)^{3/2} + C$$

By performing the backsubstitution:

$$\int 5t^3 \sqrt{2+t^4} \, dt = \frac{5}{6} (2+t^4)^{3/2} + C$$

Example 8.5 Consider an example, where the converse approach is useful:

$$\int e^x \sqrt{1+e^x} \, dx$$

Introduce the substitution $x = g(t) = \ln t$, then $g'(t) = 1/t$, and we obtain:

$$\int e^x \sqrt{1+e^x} \, dx = \int t \sqrt{1+t} \frac{1}{t} \, dt = \frac{2}{3} (1+t)^{3/2} + C$$

By the backsubstitution $t = e^x$ we get:

$$\int e^x \sqrt{1+e^x} \, dx = \frac{2}{3} (1+e^x)^{3/2} + C$$

8.4 Substitution in definite integrals

When substitution is applied in definite integrals, instead of backsubstitution, it is much more efficient to change the bounds of the integral according to the substitution:

$$\int_a^b f(g(t))g'(t) dt = \int_{g(a)}^{g(b)} f(x) dx$$

Example 8.6 In the example below we use the setting $x = g(t) = \cos t$, then $g'(t) = -\sin t$, and

$$\begin{aligned} \int_0^{\pi/2} \frac{\sin 2t}{1 + \cos^2 t} dt &= \int_0^{\pi/2} \frac{2 \sin t \cos t}{1 + \cos^2 t} dt \\ &= - \int_1^0 \frac{2x}{1 + x^2} dx = \int_0^1 \frac{2x}{1 + x^2} dx \\ &= [\ln(1 + x^2)]_0^1 = \ln 2 \end{aligned}$$

Example 8.7 Apply this rule to evaluate the following celebrated integral:

$$\int_0^1 \sqrt{1 - x^2} dx$$

Introduce the substitution $x = g(t) = \sin t$, then $g'(t) = \cos t$ and (please observe the change of the bounds of the integral!):

$$\begin{aligned} \int_0^1 \sqrt{1 - x^2} dx &= \int_0^{\pi/2} \cos^2 t dt \\ &= \frac{1}{2} \int_0^{\pi/2} (1 + \cos 2t) dt = \frac{1}{2} \left[t + \frac{\sin 2t}{2} \right]_0^{\pi/2} = \frac{\pi}{4} \end{aligned}$$

The geometric interpretation of this example is as follows. We determined the area of the first quadrant of the unit circle with center at the origin!

8.5 Linear differential equations

By a differential equation we mean an equation in which the unknown function and its derivative appear. Several problems and models in micro and macroeconomics lead to such equations. A typical equation like that is the linear differential equation.

Let a and b be given real numbers, and we are looking for the unknown differentiable function y for which

$$\begin{aligned}y' &= ay + b \\ y(0) &= y_0\end{aligned}\tag{8.1}$$

where y_0 is an "a priori" given real number.

The equality $y(0) = y_0$ is called the initial condition. We say that the differentiable function y is a solution to the above problem, if for any $t \in \mathbb{R}$ we have $y'(t) = ay(t) + b$, moreover $y(0) = y_0$. The question is, how to find the solution of this problem?

Let us suppose that y is a solution. Multiply both sides of the equation by the expression e^{-at} , then after rearranging the terms, we get

$$y'(t)e^{-at} - ay(t)e^{-at} = be^{-at}$$

for every real number t . valós számra. Observe that on the left-hand side we have precisely the derivative of the product $y(t)e^{-at}$. Therefore, by integrating both sides from 0 to t -ig (and changing the variable of the integration from t to s)

$$\int_0^t (y'(s)e^{-as} - ay(s)e^{-as}) ds = [y(s)e^{-as}]_0^t = \int_0^t be^{-as} ds$$

By plugging in the bounds we receive

$$y(t)e^{-at} - y(0) = \int_0^t be^{-as} ds.$$

Rearranging and multiplying both sides by the expression e^{at} we can formulate our result in the following theorem.

Theorem 8.8 (Cauchy-formula) *The solution to problem (8.1) is given by*

$$y(t) = e^{at} \left(y_0 + \int_0^t be^{-as} ds \right)$$

on the entire real line.

Recall that without prescribing the initial condition $y(0) = y_0$ the linear differential equation (8.1) would possess infinitely many solutions.

Example 8.9 For instance, if we are looking for the solution of the linear differential equation

$$\begin{aligned}y' &= 2y + 5 \\ y(0) &= 3\end{aligned}$$

then by the Cauchy-formula we conclude that

$$y(t) = e^{2t} \left(3 + \int_0^t 5e^{-2s} ds \right) = e^{2t} \left(3 - \frac{5}{2} [e^{-2s}]_0^t \right) = \frac{11}{2} e^{2t} - \frac{5}{2}$$

for each $t \in \mathbb{R}$.

Verify that this is the correct solution, by direct substitution!

Study at home:

1. Careful review of "Mathematical Analysis Exercises"
2. Textbook-1, Sections 11.1 and 11.2.

Chapter 9

Extension of integration

9.1 Improper integrals

Assume that f is continuous on an infinite interval $[a, +\infty)$. Then for every $b \geq a$ the integral $\int_a^b f(x) dx$ exists.

Definition 9.1 We say that the *improper integral* of f exists (or convergent) on the infinite interval $[a, \infty)$, if the limit $\lim_{b \rightarrow \infty} \int_a^b f(x) dx$ exists and it is finite. The value of the improper integral is defined by

$$\int_a^\infty f(x) dx = \lim_{b \rightarrow \infty} \int_a^b f(x) dx$$

If the limit above is not finite, or does not exist, then we say that the improper integral does not exist (or not convergent).

We define the improper integral

$$\int_{-\infty}^a f(x) dx$$

in a completely analogous way.

Example 9.2 Investigate the improper integral

$$\int_1^\infty \frac{1}{x} dx$$

By the definition

$$\int_1^b \frac{1}{x} dx = [\ln x]_1^b = \ln b$$

Passing to the limit $b \rightarrow \infty$ we see that the limit of $\ln b$ is not finite, therefore this improper integral is not convergent.

However, the improper integral

$$\int_1^{\infty} \frac{1}{x^2} dx$$

does exist, since

$$\lim_{b \rightarrow \infty} \int_1^b \frac{1}{x^2} dx = \lim_{b \rightarrow \infty} \left[-\frac{1}{x} \right]_1^b = 1$$

and the value of the improper integral is 1.

By applying the same argument, we see that the improper integral

$$\int_1^{\infty} \frac{1}{x^\alpha} dx$$

is convergent if and only if $\alpha > 1$, and its value is

$$\int_1^{\infty} \frac{1}{x^\alpha} dx = \frac{1}{\alpha - 1} \quad (9.1)$$

since the limit at the upper bound is zero.

Example 9.3 Consider the following important example (density function of the exponential distribution):

$$\int_0^{\infty} \lambda e^{-\lambda x} dx$$

where $\lambda > 0$ is a given constant. Then for any $b > 0$ we have:

$$\int_0^b \lambda e^{-\lambda x} dx = [-e^{-\lambda x}]_0^b = 1 - e^{-\lambda b}$$

Consequently

$$\int_0^{\infty} \lambda e^{-\lambda x} dx = \lim_{b \rightarrow \infty} (1 - e^{-\lambda b}) = 1$$

for any given constant $\lambda > 0$.

9.2 Improper integrals on the real line

Definition 9.4 We say that improper integral of f on the real line exists, if the integrals

$$\int_{-\infty}^0 f(x) dx \quad \text{and} \quad \int_0^{\infty} f(x) dx$$

are convergent. Then the value of $\int_{-\infty}^{\infty} f(x) dx$ is given by the sum of the two integrals.

Example 9.5 For instance, the improper integral

$$\int_{-\infty}^{\infty} \frac{2x}{1+x^2} dx$$

does not exist, although for any given $b > 0$ we get

$$\int_{-b}^b \frac{2x}{1+x^2} dx = 0$$

because the integrand is an odd function. However,

$$\int_0^b \frac{2x}{1+x^2} dx = \ln(1+b^2)$$

and its limit is $+\infty$, when $b \rightarrow \infty$ and according to the definition the integral is not convergent. The same can be said about the integral on $(-\infty, 0]$.

Example 9.6 Evaluate the following improper integral:

$$I = \int_0^{\infty} x e^{-cx^2} dx$$

where $c > 0$ is a given constant. Here for every $b > 0$ we obtain

$$\int_0^b x e^{-cx^2} dx = \left[-\frac{1}{2c} e^{-cx^2} \right]_0^b$$

This implies that $I = 1/2c$. On the other hand, the integrand is an odd function, thus,

$$\int_{-\infty}^{\infty} x e^{-cx^2} dx = 0.$$

Note that it was important to verify that the integral is convergent!

Example 9.7 (Gauss-integral) The following integral is important in probability theory:

$$I = \int_{-\infty}^{\infty} e^{-x^2} dx$$

(density function of the normal distribution). The evaluation of this improper integral needs some sophisticated calculations, we skip the details here. The

reason why this problem is hard is that the primitive function cannot be given explicitly.

ATTENTION! That does not mean there is no primitive function! The integrand is continuous, which implies that the primitive function exists (see the Chapter 7). The main difficulty is that this primitive function cannot be expressed in terms of elementary functions.

It can be shown that

$$\int_0^{\infty} e^{-x^2} dx = \frac{\sqrt{\pi}}{2}$$

and therefore $I = \sqrt{\pi}$, since the integrand is an even function.

By applying the substitution $x = t\sqrt{2}$, we also see that

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2}} dx = 1 \quad (9.2)$$

This equality will play an important role in probability theory.

9.3 Integration by parts in improper integrals

In the upcoming examples we use integration by parts in improper integrals. For simplicity, instead of passing to the limit $b \rightarrow +\infty$, we briefly indicate the upper bound $+\infty$. (But we should know what it means!)

Example 9.8 Suppose that λ is a positive constant, and evaluate the improper integral

$$\int_0^{\infty} \lambda x e^{-\lambda x} dx$$

By setting $f'(x) = \lambda e^{-\lambda x}$ and $g(x) = x$ (this way we make sure that the multiplier x will disappear in the second integral), we get

$$\begin{aligned} \int_0^{\infty} \lambda x e^{-\lambda x} dx &= [-x e^{-\lambda x}]_0^{\infty} - \int_0^{\infty} -e^{-\lambda x} dx \\ &= -\left[\frac{e^{-\lambda x}}{\lambda}\right]_0^{\infty} = \frac{1}{\lambda}. \end{aligned}$$

Observe that the expression within the brackets is zero! It is a consequence of L'Hôpital's Rule.

Example 9.9 Suppose again that λ is a positive constant, and now evaluate the improper integral

$$\int_0^{\infty} \lambda x^2 e^{-\lambda x} dx$$

Applying again the setting $f'(x) = \lambda e^{-\lambda x}$ and $g(x) = x^2$ (this way we make sure that the degree of the multiplier x^2 decreases), by two consecutive integrations by parts (with the same setting) we obtain

$$\begin{aligned} \int_0^{\infty} \lambda x^2 e^{-\lambda x} dx &= [-x^2 e^{-\lambda x}]_0^{\infty} - \int_0^{\infty} -2x e^{-\lambda x} dx \\ &= \left[\frac{-2x e^{-\lambda x}}{\lambda} \right]_0^{\infty} - \int_0^{\infty} -2 \frac{e^{-\lambda x}}{\lambda} dx \\ &= \left[-2 \frac{e^{-\lambda x}}{\lambda^2} \right]_0^{\infty} = \frac{2}{\lambda^2} \end{aligned}$$

In this example we needed two integrations by parts in a row to eliminate the multiplier x^2 . In view of the L'Hôpital-Rule, the expressions inside the brackets are zero.

Example 9.10 Use integration by parts to evaluate the improper integral

$$\int_{-\infty}^{\infty} x^2 e^{-x^2/2} dx$$

By allocating the roles among the factors in a smart way, we conclude:

$$\int_{-\infty}^{\infty} (-x) \cdot (-x e^{-x^2/2}) dx = [-x e^{-x^2/2}]_{-\infty}^{\infty} + \int_{-\infty}^{\infty} e^{-x^2/2} dx = \sqrt{2\pi}$$

where we relied on formula (9.2). Indeed, making use of L'Hôpital's Rule, we see that both limits of the expression within the brackets are zero, hence

$$\int_{-\infty}^{\infty} x^2 e^{-x^2/2} dx = \sqrt{2\pi}. \quad (9.3)$$

9.4 Harmonic series revisited

As we have seen in Chapter 2, for a given exponent $\alpha > 0$ the infinite series

$$\sum_{k=1}^{\infty} \frac{1}{k^\alpha} \quad (9.4)$$

is divergent if $\alpha \leq 1$, and it is convergent if $\alpha \geq 2$. However, we were unable to find the answer when $1 < \alpha < 2$. Now we give a complete solution by using improper integrals. Consider n -th partial sum of the series

$$S_n = \sum_{k=1}^n \frac{1}{k^\alpha}$$

and sketch the graph of the function

$$f(x) = \frac{1}{x^\alpha}$$

on the positive part of the real line. Take the values of the functions at the integers $1, \dots, n$, then by examining the graph we can easily see that

$$S_n < 1 + \int_1^n \frac{1}{x^\alpha} dx$$

since the function f is strictly monotone decreasing.

ATTENTION! Check Figures.pdf for the details!

On the other hand f is positive, and for $\alpha > 1$ its improper integral on the interval $[1, \infty)$ is convergent, see the equality (9.1). Therefore

$$S_n < 1 + \int_1^n \frac{1}{x^\alpha} dx < 1 + \int_1^\infty \frac{1}{x^\alpha} dx = 1 + \frac{1}{\alpha - 1} = \frac{\alpha}{\alpha - 1}.$$

We conclude that S_n is bounded from above, and it is clearly strictly monotone increasing, hence it is convergent. We summarize this result in the following theorem.

Theorem 9.11 *The infinite series (9.4) is convergent if and only if $\alpha > 1$, and in this case*

$$\sum_{k=1}^{\infty} \frac{1}{k^\alpha} < \frac{\alpha}{\alpha - 1}$$

Study at home

1. Careful review of "Mathematical Analysis Exercises"
2. Textbook-1, Sections 11.3 and 11.4.

Chapter 10

Power series

10.1 Sum of power series

If $-1 < x < 1$ is a given real number, then the geometric series

$$\frac{1}{1-x} = \sum_{k=0}^{\infty} x^k.$$

is convergent. It is an interesting question if a given function f can be given in the form

$$f(x) = \sum_{k=0}^{\infty} a_k x^k \tag{10.1}$$

with appropriate coefficients a_k . In this case we say that f can be expanded in a power series.

Definition 10.1 The series on the right-hand side of the equality (10.1) is called a *power series*, the function f on the left-hand side is called the *sum* of the power series.

In this chapter we examine two interesting questions.

1. For what values x is the power series convergent, and what is its sum f .
2. Conversely, if a function f is given, how can we find the power series whose sum is precisely f (if possible).

A power series is obviously convergent for $x = 0$ and its sum is a_0 . The set of all values of x for which the power series is convergent is called the *set of convergence*.

10.2 Radius of convergence

The set of convergence of a power series is always an interval that is symmetric about the origin. This fact is formulated in the following theorem.

Theorem 10.2 (Cauchy-Hadamard-theorem) *For the power series (10.1) there exists a nonnegative number R (maybe $R = 0$ or infinity) so that the series is convergent in the open interval $-R < x < R$, and it is divergent outside the closed interval $[-R, R]$.*

Proof. We just restrict our attention to the case when the limit

$$\lim_{k \rightarrow \infty} \left| \frac{a_{k+1}}{a_k} \right| = r$$

exists. Introduce the notation:

$$R = \begin{cases} 1/r & \text{if } 0 < r < +\infty \\ +\infty & \text{if } r = 0 \\ 0 & \text{if } r = \infty \end{cases}$$

In view of the Quotient Test the series is convergent, if

$$\lim_{k \rightarrow \infty} \left| \frac{a_{k+1}x^{k+1}}{a_kx^k} \right| < 1$$

which exactly means that $|x| < R$.

A completely analogous argument shows that the series is divergent when $|x| > R$. \square

ATTENTION!

This theorem says nothing about the boundary of the interval! At $|x| = R$ the series may or may not be convergent. This cannot be decided by our theorem, further analysis is needed.

Definition 10.3 The number R above is called the *radius of convergence* of the power series.

Example 10.4 Consider the power series

$$\sum_{k=0}^{\infty} \frac{x^k}{k!}$$

Here we have

$$\lim_{k \rightarrow \infty} \frac{a_{k+1}}{a_k} = \lim_{k \rightarrow \infty} \frac{k!}{(k+1)!} = \lim_{k \rightarrow \infty} \frac{1}{k+1} = 0$$

and hence $R = \infty$. This means that the power series is convergent on the whole real line.

Another example is the power series

$$\sum_{k=1}^{\infty} \frac{x^k}{k}$$

Then we get

$$\lim_{k \rightarrow \infty} \frac{a_{k+1}}{a_k} = \lim_{k \rightarrow \infty} \frac{k}{k+1} = 1$$

and hence $R = 1$. We conclude that the series is convergent in the open interval $(-1, 1)$, and it is divergent outside the closed interval $[-1, 1]$.

On the other hand, we see that for $x = 1$ we obtain the divergent harmonic series, and further for $x = -1$ we get a convergent series with alternating signs, see Example 2.12. Thus, the interval of convergence of this power series is the interval

$$[-1, 1)$$

closed from the left and open from the right. Please observe that on the boundary anything can happen!

10.3 Differentiability of power series

Consider a power series whose radius of convergence is $R > 0$ and its sum function is f that is

$$\sum_{k=0}^{\infty} a_k x^k = f(x)$$

for every $-R < x < R$.

Theorem 10.5 *The sum f of the power series is differentiable, in particular*

$$f'(x) = \sum_{k=1}^{\infty} k a_k x^{k-1}$$

in the open interval $(-R, R)$.

We do not prove this theorem (it is technical), just note that it is based on the so-called "uniform convergence" principle. Some consequences however, can easily be derived from this statement.

- The derivative of the sum is obtained from differentiating the power series term by term. This is not obvious, since the sum rule (in general) is not true for infinitely many terms. FIND COUNTEREXAMPLES!

- Observe that the radius of convergence of the derivative power series is still R . VERIFY!
- As we see that f' is the sum of a power series in the same interval, by repeated applications of the theorem, we deduce that f is infinitely many times differentiable in the open interval $(-R, R)$.

Example 10.6 Consider the geometric series in the open interval $-1 < x < 1$

$$\sum_{k=0}^{\infty} x^k = \frac{1}{1-x}$$

Note that the first term is 1, whose derivative is zero. Making use of our theorem

$$\sum_{k=1}^{\infty} kx^{k-1} = \frac{1}{(1-x)^2}$$

for every $-1 < x < 1$.

Example 10.7 Find the function f that is given by the following power series:

$$f(x) = \sum_{k=1}^{\infty} (-1)^{k-1} \frac{x^k}{k}$$

A simple calculation shows that the radius of convergence is $R = 1$. On the one hand $f(0) = 0$, on the other hand, by the differentiability of the power series

$$f'(x) = \sum_{k=1}^{\infty} k \frac{(-x)^{k-1}}{k} = \sum_{k=1}^{\infty} (-x)^{k-1} = \frac{1}{1+x}$$

for each $-1 < x < 1$. This implies

$$f(x) = f(0) + \int_0^x \frac{1}{1+t} dt = [\ln(1+t)]_0^x = \ln(1+x)$$

in the open interval $(-1, 1)$. Moreover, by Example 2.12 the original series is convergent at $x = 1$, which leads to the celebrated identity

$$1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \frac{1}{5} - \dots = \ln 2$$

However, the series is divergent at $x = -1$.

10.4 Finding the coefficients

Suppose that a function f can be given as the sum of a power series in the interval of convergence. Then necessarily f is infinitely many times differentiable in the interval. How could we determine the coefficients of the power series?

By successively taking the derivatives of both sides of equality (10.1), the coefficients a_k can be computed step by step. Indeed, observe that

$$f(0) = a_0, \quad f'(0) = a_1, \quad f''(0) = 2a_2, \quad \dots$$

and in general, for any given index k we get:

$$f^{(k)}(0) = k! \cdot a_k$$

If we substitute these expressions for a_k in the power series, then we have

$$f(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(0)}{k!} x^k$$

This form is called the *Taylor-series* (or Taylor expansion) of f .

10.5 Taylor-series of the exponential function

In this section we consider the exponential function $f(x) = e^x$. If this function is the sum of a power series, then the coefficients can only be

$$a_k = \frac{1}{k!}$$

for every k . Indeed, any derivative of e^x is e^x , which takes the value 1 at $x = 0$. Therefore, the Taylor-series associated with the function e^x is:

$$\sum_{k=0}^{\infty} \frac{x^k}{k!}$$

and we have seen that this series is convergent on the entire real line.

The reason why we did not write equality is that it is not yet clear at the moment that the sum of this series is really e^x .

To overcome this difficulty, consider the function

$$f(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

on the real line, which is yet to be determined. Clearly $f(0) = 1$. On the other hand, in view of the differentiability theorem:

$$f'(x) = \sum_{k=1}^{\infty} k \frac{x^{k-1}}{k!} = \sum_{k=1}^{\infty} \frac{x^{k-1}}{(k-1)!} = f(x)$$

for every $-\infty < x < \infty$. This is a simple linear differential equation for the unknown f , whose only solution is

$$f(x) = e^x$$

on the real line. As a consequence, we deduce the celebrated identity

$$e = 1 + \frac{1}{1!} + \frac{1}{2!} + \dots + \frac{1}{n!} + \dots$$

by substituting $x = 1$.

Study at home

1. Careful review of "Mathematical Analysis Exercises"
2. Textbook-1, Section 6.5.

Chapter 11

Functions of two variables

11.1 Partial derivatives

Consider a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ of two variables. Fix the coordinate $y = b$ and examine the function

$$x \rightarrow f(x, b)$$

of only one variable. Assume that this function is differentiable at a point a , and determine its derivative.

Definition 11.1 The derivative above is called the *partial derivative* of the function f with respect to the variable x at the point (a, b) . We denote it by

$$\frac{\partial f}{\partial x}(a, b) = f'_1(a, b)$$

Sometimes the notation $f'_x(a, b)$ is also used.

Example 11.2 Consider for instance the function $f(x, y) = (x + 2y)e^{x+3y-1}$ and find its partial derivative with respect to x at the point $(1, 1)$.

Then $f(x, 1) = (x + 2)e^{x+2}$, whose derivative at any x is

$$f'_1(x, 1) = e^{x+2} + (x + 2)e^{x+2} = (x + 3)e^{x+2}$$

Substituting $x = 1$ we obtain $f'_1(1, 1) = 4e^3$.

Example 11.3 Principally, we could also calculate the partial derivative of the function f with respect to the variable x with an arbitrarily selected and

fixed y , and substitute the values $x = a$ and $y = b$. This is of course good, but not always convenient, as shown in the following example. Take

$$f(x, y) = \sqrt{x^2 + y^2 + 5} \cdot e^{-2x+y} \cdot \cos(y + \pi/2)$$

and find the partial derivative with respect to x at the point $(1, 0)$. Then the above way would give you the right answer, but it requires a long calculation and very time consuming. However, if we follow the definition, then we see that

$$f(x, 0) = 0$$

for every x , and therefore $f'_1(1, 0) = 0$.

The correspondence

$$x \rightarrow \frac{\partial f}{\partial x}(x), \quad x \in \mathbb{R}$$

is called the *partial derivative function* of f with respect to the variable x .

11.2 Tangent planes

Partial derivatives (similarly to the one variable case) can be given a nice geometric interpretation. Consider a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ with two variables. The graph of this function is a surface in the three dimensional space. Pick a point

$$P(a, b, f(a, b))$$

on the surface. If this surface is intersected by the plane $y = b$ passing through the point P , then we get a curve lying on the surface. The slope of the tangent line to this curve at P is exactly the partial derivative $f'_1(a, b)$. We can give an analogous interpretation for the slope of the tangent line that lies in the plane $x = a$. The plane spanned by the two tangent lines has the following normal vector (perpendicular):

$$v = (f'_1(a, b), f'_2(a, b), -1)$$

By using the notation $c = f(a, b)$ the equation of this plane is

$$f'_1(a, b)(x - a) + f'_2(a, b)(y - b) - (z - c) = 0.$$

This plane is called the *tangent plane* to the surface at the point P .

Example 11.4 Find the value of the parameter p if the tangent plane to the function

$$f(x, y) = px\sqrt{x^2 + y^2 + 1} - 7$$

at the point $a = 2$, $b = 2$, $c = f(2, 2)$ passes through the point $Q(2, -1, 6)$.

Simple substitution shows that $f(2, 2) = 6p - 7$, which means that we are looking for the equation of the tangent plane at the point $P(2, 2, 6p - 7)$. Calculate the partial derivatives:

$$\frac{\partial f}{\partial x}(2, 2) = \frac{13}{3}p \quad \text{and} \quad \frac{\partial f}{\partial y}(2, 2) = \frac{4}{3}p$$

Hence, the equation of the tangent plane at P is:

$$\frac{13}{3}p(x - 2) + \frac{4}{3}p(y - 2) - (z - 6p + 7) = 0.$$

If the tangent plane passes through the point Q , then its coordinates satisfy the equation of the plane. This gives us the following equation for the unknown parameter p :

$$-4p = 13 - 6p.$$

The only solution is $p = 13/2$.

11.3 Chain Rule

Consider now the functions $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}^2$ where for every $t \in \mathbb{R}$ we use the notation

$$g(t) = (g_1(t), g_2(t))$$

Suppose that the range of g lies in the domain of f . Then we may examine the composition

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

We want to give a condition on the differentiability of $f \circ g$.

Theorem 11.5 (Chain Rule) *If both g_1 and g_2 are differentiable at t , and the partial derivative functions of f are continuous at $g(t)$, then $f \circ g$ is differentiable at t , and*

$$(f \circ g)'(t) = \frac{\partial f}{\partial x}(g(t))g_1'(t) + \frac{\partial f}{\partial y}(g(t))g_2'(t)$$

Our theorem is very similar to the Chain Rule with one variable (see Chapter 4). Its proof (skipped) would follow the same ideas, but technically a bit more involved.

Example 11.6 Take for instance $f(x, y) = x^2 - xy + y^2$, and

$$x = g_1(t) = \cos t \quad y = g_2(t) = \sin t$$

and consider the composition function $F(t) = (f \circ g)(t)$. Making use of the Chain Rule

$$\begin{aligned} F'(t) &= (f \circ g)'(t) = \frac{\partial f}{\partial x}(g(t))g'_1(t) + \frac{\partial f}{\partial y}(g(t))g'_2(t) \\ &= (2 \cos t - \sin t)(-\sin t) + (-\cos t + 2 \sin t) \cos t = \sin^2 t - \cos^2 t \end{aligned}$$

for every $t \in \mathbb{R}$.

Example 11.7 Consider the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ and suppose that its partial derivatives exist and are continuous. Take the vector $v = (v_1, v_2) \in \mathbb{R}^2$ in the plane, and let a point $P(a, b) \in \mathbb{R}^2$ be given. Then the equation of the straight line in the direction v and passing through the point $P(a, b)$ is:

$$g(t) = (a, b) + tv = (a + tv_1, b + tv_2).$$

Using these notations we have $g'_1(t) = v_1$, $g'_2(t) = v_2$. Further, take the composition function

$$F(t) = f((a, b) + tv)$$

then by the Chain Rule, its derivative is given by:

$$F'(t) = \frac{\partial f}{\partial x}((a, b) + tv)v_1 + \frac{\partial f}{\partial y}((a, b) + tv)v_2$$

In particular for $t = 0$ we obtain:

$$F'(0) = \frac{\partial f}{\partial x}(a, b)v_1 + \frac{\partial f}{\partial y}(a, b)v_2$$

11.4 Local extrema

The absolute value (or the distance from the origin) of a vector $v = (x, y)$ in the two dimensional plane is defined by:

$$\|v\| = (x^2 + y^2)^{1/2}$$

that is called the *norm* of the vector v .

Definition 11.8 In the plane \mathbb{R}^2 the set

$$B = \{v \in \mathbb{R}^2 : \|v\| \leq 1\}$$

is called the *unit disk* (with center at the origin and radius equals 1). Consequently, a disk with center at the point $(a, b) \in \mathbb{R}^2$ and radius $r > 0$ is given by

$$(a, b) + rB = \{v \in \mathbb{R}^2 : \|v - (a, b)\| \leq r\}$$

(i.e. the set of points, whose distance from the center is at most r).

Consider a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. We say that a point $P(a, b)$ in the domain is a local minimum point of f , if there exists a $\varepsilon > 0$ such that

$$f(x, y) \geq f(a, b)$$

for all points (x, y) in the domain of f , where $(x, y) \in (a, b) + \varepsilon B$, that is $\|(x, y) - (a, b)\| \leq \varepsilon$.

The local maximum is defined analogously. For global minimum or maximum the inequality must hold on the entire domain of f .

11.5 First order necessary condition

In this section we suppose that partial derivatives of the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ exist and are continuous.

Theorem 11.9 *If the point $(a, b) \in \mathbb{R}^2$ is a local minimum point of f , then $f'_1(a, b) = 0$ and $f'_2(a, b) = 0$.*

Proof. Take a non zero vector $v \in \mathbb{R}^n$ arbitrarily, and consider the composition function

$$F(t) = f((a, b) + tv).$$

In view of our assumption the function F has a local minimum at $t = 0$. On the other hand, F is differentiable, namely

$$F'(t) = \frac{\partial f}{\partial x}((a, b) + tv)v_1 + \frac{\partial f}{\partial y}((a, b) + tv)v_2$$

Applying Theorem 5.7 we get $F'(0) = 0$ for every vector v , in other words

$$\frac{\partial f}{\partial x}((a, b))v_1 + \frac{\partial f}{\partial y}((a, b))v_2 = 0$$

for all real numbers v_1 and v_2 . This is only possible if

$$\frac{\partial f}{\partial x}((a, b)) = 0 \quad \text{and} \quad \frac{\partial f}{\partial y}((a, b)) = 0$$

and this is exactly that we wanted to prove. \square

Analogous theorem applies for the case of local maximum.

This theorem tells us that local extrema can only be at points where both partial derivatives are zero. In other words, local extrema can be only be found in the solution set of the system of equations with both partial derivatives being zero. This is however, just a necessary condition (just like in the one-variable case), and by no means sufficient! For example, in the case of the function

$$f(x, y) = x^3y^2$$

we have the necessary condition $f'_1(x, y) = f'_2(x, y) = 0$. A solution to this system is $(x, y) = (0, 0)$, and at this point

$$f(0, 0) = 0$$

But this is neither a minimum nor a maximum. It is easy to see that the function has both positive and negative values in any disk around the origin (with whatever positive radius). Thus, the origin cannot be a local extreme point.

Example 11.10 Consider the function

$$f(x, y) = \frac{1}{x} + \frac{1}{y} + \frac{xy}{8}$$

on the plane, where $x \neq 0$ and $y \neq 0$, and try to find its local extreme points. Find the zeros of the partial derivatives!

$$\begin{aligned} \frac{\partial f}{\partial x} &= -\frac{1}{x^2} + \frac{y}{8} = 0 \\ \frac{\partial f}{\partial y} &= -\frac{1}{y^2} + \frac{x}{8} = 0 \end{aligned}$$

The only solution to the simultaneous equations is

$$x = 2 \quad \text{and} \quad y = 2,$$

therefore f can only have a local extremum (minimum or maximum) at this point.

A comprehensive method for deciding whether or not a critical point is a local extremum will be discussed in the Linear Algebra course (third semester, sophomore year). We note here that $P(2, 2)$ is in fact a local minimum point of f (see the "Mathematical Analysis Exercises" for more details).

Study at home

1. Careful review of "Mathematical Analysis Exercises"
2. Textbook-1, Sections 15.3, 15.4, 15.6, 16.1 and 16.2.

Chapter 12

Constrained extrema

12.1 Implicit functions

A problem often encountered in microeconomics is the following. If an equation

$$F(x, y) = 0$$

is given, can we uniquely express the variable y from the equation as a function of x ? In other words: can we find a unique function $y = g(x)$ such that the identity

$$F(x, g(x)) = 0$$

holds at every point x ?

Such a function does not necessarily exist. For example, in the case of the equation

$$F(x, y) = x^2 + y^2 - 1 = 0$$

(equation of the unit circle) the variable y cannot be expressed uniquely as a function of x . Geometrically this means that the set of points on the plane that satisfy the equation $F(x, y) = 0$ cannot be the graph of a function. The reason for this is that some vertical lines (parallel to the y -axis) intersect this curve twice.

It may even happen that the variable y cannot be expressed from the equation by algebraic manipulations. Such an example is the equation

$$F(x, y) = e^{x+y} - 2 \cos y + 1 = 0$$

It is easy to see that the point $(x, y) = (0, 0)$ satisfies the equation, but the variable y cannot be isolated on one side.

We also raise the following question. If F is differentiable, then can we express the variable y from the equation as a differentiable function of x ? This

question is answered by the following theorem.

Theorem 12.1 (Implicit function theorem) *Assume that the at the point (x_0, y_0) we have*

$$F(x_0, y_0) = 0$$

moreover the partial derivatives of F are continuous in a neighborhood of this point, and

$$F'_2(x_0, y_0) \neq 0$$

Then there exists a unique continuously differentiable function g in a neighborhood of the point x_0 such that

- $g(x_0) = y_0$
- $F(x, g(x)) = 0$ at every point x
- $g'(x) = -F'_1(x, g(x))/F'_2(x, g(x))$

We point out that from the continuity of the partial derivatives we get that $F'_2(x, g(x)) \neq 0$ in a neighborhood of the point x_0 .

The geometric interpretation of our theorem is that if the tangent line to the planar curve with equation $F(x, y) = 0$ at the point (x_0, y_0) is not parallel to the y -axis (i.e. "the curve cannot turn back"), then y can be expressed (locally) as a differentiable function of x .

Example 12.2 Consider the implicit equation

$$F(x, y) = e^{x+y} + x + y - 1 = 0$$

The point $(0, 0)$ satisfies the equation. On the other hand, at this point

$$F'_2(0, 0) = 2$$

Hence, F fulfills the conditions of the Implicit function theorem: there exist a unique differentiable function $y = g(x)$ with

$$\begin{aligned} g'(x) &= -F'_1(x, g(x))/F'_2(x, g(x)) \\ &= -\frac{1}{e^{x+g(x)} + 1} \cdot (e^{x+g(x)} + 1) = -1 \end{aligned}$$

at every point x . Since $g(0) = 0$, this implies

$$g(x) = -x$$

and this is the only solution.

Example 12.3 A slightly more complicated example is

$$F(x, y) = e^{x+y} - 2 \cos y + 1 = 0$$

The point $(0, 0)$ satisfies the equation. On the other hand, at this point

$$F'_2(0, 0) = 1$$

and hence, the conditions of the Implicit function theorem are fulfilled. We conclude that the equation uniquely determines a differentiable function g so that $F(x, g(x)) = 0$ at every x . However, this function cannot be expressed explicitly by using algebraic manipulations.

12.2 Constrained minima

Consider the functions f and F that are both $\mathbb{R}^2 \rightarrow \mathbb{R}$ and suppose that their partial derivatives are continuous. By a *constrained minimum* problem we mean the following problem:

$$\begin{aligned} f(x, y) &\rightarrow \min & (12.1) \\ F(x, y) &= c \end{aligned}$$

where c is a given real constant. In other words, we look for the minimum (or sometimes maximum) of f on the set

$$H = \{(x, y) \in \mathbb{R}^2 : F(x, y) = c\}$$

This equality is called the *constraint*.

Definition 12.4 We say that the point $(x_0, y_0) \in H$ is the solution of the constrained minimization problem (12.1) if

$$f(x_0, y_0) \leq f(x, y)$$

for every $(x, y) \in H$ esetén. An analogous definition applies for maximum problems.

Example 12.5 The example below illustrates that for constrained minimization problems the usual necessary conditions for extrema do not work. Consider the constrained minimization problem

$$f(x, y) = x^2 + 2y, \quad F(x, y) = x + y = 0 \quad \text{i.e.} \quad c = 0$$

From the constraint $x + y = 0$ we get $y = -x$, and consequently $f(x, y) = x^2 - 2x$ on the set H . This function achieves its minimum at the point $x = 1$ and in H this necessarily means $y = -1$. Thus, the constrained minimum is at the point

$$(x_0, y_0) = (1, -1)$$

However, at this point none of the equalities

$$\frac{\partial f}{\partial x} = 0, \quad \frac{\partial f}{\partial y} = 0$$

is true. Verify this!

This example also exhibits that a constrained extremum problem can be transformed into a non-constrained extremum problem by expressing the variable y as a function of x from the constraint $F(x, y) = c$. In more complicated problems this may not be possible by algebraic manipulations. This is the point where we need the Implicit function theorem.

12.3 Lagrange multipliers

Consider the constrained minimization problem (12.1). By using the Implicit function theorem we make sure that the variable y can be expressed from the constraint $F(x, y) = c$, and that way we can solve the problem. This procedure is described below.

Definition 12.6 The *Lagrange-function* (or Lagrangian) of the problem (12.1) is defined by

$$\mathcal{L}(x, y, \lambda) = f(x, y) - \lambda(F(x, y) - c)$$

λ is an arbitrary real number.

Theorem 12.7 (Lagrange-method) *Let us suppose that (x_0, y_0) is the solution of the problem (12.1), and assume that the partial derivatives of f and F are continuous in a neighborhood of this point. If*

$$F'_2(x_0, y_0) \neq 0, \tag{12.2}$$

there exists a unique real number λ such that

$$\frac{\partial \mathcal{L}}{\partial x}(x_0, y_0, \lambda) = 0, \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial y}(x_0, y_0, \lambda) = 0$$

Proof. In view of (12.2) the conditions of the Implicit function theorem are fulfilled. Thus, there exists a unique continuously differentiable function g with

- $g(x_0) = y_0$, and
- $F(x, g(x)) = c$ in a neighborhood of x_0 , furthermore

$$\bullet g'(x_0) = -F'_1(x_0, y_0)/F'_2(x_0, y_0).$$

If (x_0, y_0) is the solution of problem (12.1), then the function $x \rightarrow f(x, g(x))$ achieves its minimum at x_0 , therefore, its derivative at this point is zero. Applying the Chain Rule, the derivative can be given in this form:

$$f'_1(x_0, y_0) + f'_2(x_0, y_0)g'(x_0) = f'_1(x_0, y_0) - \frac{f'_2(x_0, y_0)}{F'_2(x_0, y_0)}F'_1(x_0, y_0) = 0.$$

Introduce the notation:

$$\lambda = \frac{f'_2(x_0, y_0)}{F'_2(x_0, y_0)}.$$

Using this notation, the above derivative can be rewritten:

$$\frac{\partial \mathcal{L}}{\partial x}(x_0, y_0, \lambda) = f'_1(x_0, y_0) - \lambda F'_1(x_0, y_0) = 0.$$

The second equality of the theorem is trivial by simply substituting λ . Indeed:

$$\frac{\partial \mathcal{L}}{\partial y}(x_0, y_0, \lambda) = f'_2(x_0, y_0) - \lambda F'_2(x_0, y_0) = 0. \square$$

Our theorem could be formulated analogously for the case of maximum.

12.4 Solving the constrained minimization problem

The procedure of solving the constrained minimization problem (12.1) is as follows.

1. Find the Lagrange-function of the problem.
2. Find the partial derivatives with respect to x and y , and make them equal zero.
3. Take into account that $F(x_0, y_0) = c$.
4. Solve the system of three equations for x , y and λ .

The point (x_0, y_0) obtained that way satisfies the necessary condition for an extremum. The solution λ is called the *Lagrange multiplier* associated with the problem.

Example 12.8 Now solve the constrained minimization problem in Example 12.5 by using the Lagrange-method. The Lagrange-function of the problem is:

$$\mathcal{L}(x, y, \lambda) = x^2 + 2y - \lambda(x + y).$$

The system of equations is of the form:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial x}(x_0, y_0, \lambda) &= 2x_0 - \lambda = 0 \\ \frac{\partial \mathcal{L}}{\partial y}(x_0, y_0, \lambda) &= 2 - \lambda = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda}(x_0, y_0, \lambda) &= x_0 + y_0 = 0\end{aligned}$$

The only solution to this system is $\lambda = 2$, $x_0 = 1$ and $y_0 = -1$.

Example 12.9 The following type of problem frequently appears in microeconomics. Find the constrained maximum of consumer demand:

$$\begin{aligned}x^\alpha y^\beta &\rightarrow \max \\ px + y &= m\end{aligned}\tag{12.3}$$

where α , β , p and m are given positive real numbers. In this problem

$$f(x, y) = x^\alpha y^\beta \quad \text{and} \quad F(x, y) = px + y,$$

Therefore, the Lagrange-function of the problem is:

$$\mathcal{L}(x, y, \lambda) = x^\alpha y^\beta - \lambda(px + y - m).$$

The system of equation that comes from the Lagrange-method:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial x}(x_0, y_0, \lambda) &= \alpha \cdot x_0^{\alpha-1} y_0^\beta - \lambda p = 0 \\ \frac{\partial \mathcal{L}}{\partial y}(x_0, y_0, \lambda) &= \beta \cdot x_0^\alpha y_0^{\beta-1} - \lambda = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda}(x_0, y_0, \lambda) &= px_0 + y_0 - m = 0.\end{aligned}$$

This system admits the following single solution:

$$px_0 = \frac{\alpha}{\alpha + \beta} m \quad \text{and} \quad y_0 = \frac{\beta}{\alpha + \beta} m,$$

The Lagrange multiplier λ can be then calculated from the second equation.

Study at home

1. Careful review of "Mathematical Analysis Exercises"
2. Textbook-1, Sections 16.3, 18.1, 18.2, 18.3, 18.4, 18.5 and 18.6.

Part II

Second Semester: Probability Theory

Chapter 13

Probability

13.1 Experiments

In the sequel we deal with experiments that have chance outcomes. In other words, the experiments have outcomes that cannot be predicted.

1. Toss a playing die and check the number that comes out.
2. Toss a pair of dice.
3. Toss a die, then flip a coin as many times as the number on the die.
4. Keep tossing a die until 6 comes out for the first time.
5. Pick a point randomly on the unit disc (with radius 1).

More complicated examples:

- The number of cars that pass an intersection between 10 am and 11 am.
- The number of calls received by a call center between 8 am and 9 am.
- The length of time period between two successive calls
- The price of a stock at the stock exchange at closing time.
- The waiting time at a customer service desk.

13.2 The sample space

Definition 13.1 Let Ω denote the set of all possible outcomes in an experiment. The set Ω is called the *sample space* associated with the experiment.

Specify the sample spaces that are associated with the previous experiments. Then in the same order:

1. $\Omega = \{1, 2, 3, 4, 5, 6\}$
2. $\Omega = \{(1, 1), (1, 2), (2, 1), (1, 3), \dots, (6, 6)\}$
3. $\Omega = \{1H, 1T, 2HH, 2HT, 2TH, 2TT, \dots\}$ (Question: how many elements are in the sample space?)
4. Ω consists of all finite sequences whose last digit is 6, and all previous digits are any of the numbers 1,2,3,4,5.
5. $\Omega = \{(x, y) \in \mathbb{R}^2 : x^2 + y^2 \leq 1\}$

13.3 Events

Definition 13.2 The subsets of the sample space are called *events*.

Take some examples in the sample spaces above.

1. Let A denote the event that the outcome is even. Then $A = \{2, 4, 6\}$.
2. Let A denote the event that the sum of the two numbers is 7. Then $A = \{(1, 6), (6, 1), (2, 5), (5, 2), (3, 4), (4, 3)\}$.
3. Let A denote the event that we have no Tail (all of them are Head). Then $A = \{1H, 2HH, 3HHH, 4HHHH, 5HHHHH, 6HHHHHH\}$.
4. Let A denote the event that we needed at most two tosses. Then $A = \{6, 16, 26, 36, 46, 56\}$.
5. Let A denote the event that the distance of the point from the center is less than $1/2$. Then $A = \{(x, y) : x^2 + y^2 < 1/4\}$.

13.4 Operations with events

We say that the event $A \subset \Omega$ occurs, if the experiment results in an outcome $\omega \in \Omega$ such that $\omega \in A$.

The impossible event has no elements, notation: \emptyset (empty set). The certain event is: Ω (the whole sample space).

1. $A \cap B$ occurs if and only if both A and B occur. We say that A and B are mutually exclusive, if $A \cap B = \emptyset$.

2. $A \cup B$ occurs if and only if either A or B occurs (or both).
3. \bar{A} (the complement of A) occurs if and only if A does not occur.

We say that A implies B (or B is a consequence of A), if $A \subset B$.

Theorem 13.3 (De Morgan Rules)

1. $\overline{A \cup B} = \bar{A} \cap \bar{B}$
2. $\overline{A \cap B} = \bar{A} \cup \bar{B}$

These identities hold true for an arbitrary number of events as well.

Proof. We demonstrate the first identity. Let $x \in \overline{A \cup B}$ be selected arbitrarily. Then

$$x \in \overline{A \cup B} \Rightarrow x \notin A \cup B \Rightarrow x \notin A \text{ and } x \notin B \Rightarrow x \in \bar{A} \text{ and } x \in \bar{B} \Rightarrow x \in \bar{A} \cap \bar{B}$$

This proves that $\overline{A \cup B} \subset \bar{A} \cap \bar{B}$. The opposite direction (and hence the equality) follows from the fact that each implication can be reversed (i.e. they are equivalences). The second identity can be verified in a completely analogous way. \square

When we carry out an experiment, some possible outcomes may not be observable. For instance, if we toss a pair of completely identical (indistinguishable) dice, we cannot decide whether the outcome is $(1, 2)$ or $(2, 1)$. We can only claim that the event $\{(1, 2), (2, 1)\}$ occurred.

Definition 13.4 Let \mathcal{A} denote the collection of *observable events*. We assume that they possess the following properties.

- If $A \in \mathcal{A}$, then $\bar{A} \in \mathcal{A}$ and $\Omega \in \mathcal{A}$.
- If $A_1, A_2, \dots \in \mathcal{A}$, then $A_1 \cup A_2 \cup \dots \in \mathcal{A}$.

Proposition 13.5 *If A and B are observable, then so is $A \cap B$.*

Proof. Indeed, if A and B are observable, then

$$A \cap B = \overline{\overline{A \cap B}} \in \mathcal{A}$$

in view of the De Morgan Rules. \square

By the De Morgan Rules, this proposition remains true for any countable number of events.

Definition 13.6 In the following, by an *experiment* we mean the couple $\mathcal{K} = (\Omega, \mathcal{A})$.

13.5 Probability space

Suppose that we perform an experiment \mathcal{K} n times in a row, and every time we observe whether or not a given event $A \in \mathcal{A}$ occurs. If A occurs k_n times out of n trials, then the relative frequency of A is:

$$\frac{k_n}{n}$$

Experience shows that by raising n , the relative frequency exhibits a dumping oscillation around a specific number. This number can be regarded as the probability of A .

Instead of using this experimental approach, below we develop an axiomatic introduction of probability. From the axioms we can derive the above experimental fact.

Definition 13.7 (Axioms of Probability) Consider an experiment $\mathcal{K} = (\Omega, \mathcal{A})$. By the *probability* we mean a function

$$P : \mathcal{A} \rightarrow [0, 1]$$

that satisfies the following two axioms:

1. $P(\Omega) = 1$
2. If $A_1, A_2, \dots \in \mathcal{A}$ are pairwise mutually exclusive events, then

$$P\left(\bigcup_{k=1}^{\infty} A_k\right) = \sum_{k=1}^{\infty} P(A_k)$$

In this case the triple (Ω, \mathcal{A}, P) is called a *probability space*.

This axiomatic approach is due to A. N. Kolmogorov (1933), and this can be regarded as the origin of modern probability theory.

From the axioms we can easily derive the following properties of probability spaces.

Theorem 13.8

1. For any $A \in \mathcal{A}$ we have

$$P(\bar{A}) = 1 - P(A)$$

and consequently $P(\emptyset) = 0$.

2. If $A, B \in \mathcal{A}$ and $A \subset B$, then

$$P(A) \leq P(B)$$

3. If $A, B \in \mathcal{A}$, then

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Proof. 1. Since $A \cup \bar{A} = \Omega$, moreover A and \bar{A} are exclusive events, the statement follows immediately from the axioms.

2. If $A \subset B$, then $A \cup (B \cap \bar{A}) = B$, moreover A and $B \cap \bar{A}$ are exclusive events, therefore, by the axioms

$$P(B) = P(A) + P(B \cap \bar{A}) \geq P(A)$$

because $P(B \cap \bar{A}) \geq 0$.

The 3. statement is proven the following way. We divide the event $A \cup B$ into disjoint pieces like this:

$$A \cup B = (A \cap \bar{B}) \cup (\bar{A} \cap B) \cup (A \cap B).$$

Then, using the axioms, we get:

$$\begin{aligned} P(A \cup B) &= P(A \cap \bar{B}) + P(\bar{A} \cap B) + P(A \cap B) \\ &= P(A) - P(A \cap B) + P(B) - P(A \cap B) + P(A \cap B) \end{aligned}$$

and the statement ensues. \square

Example 13.9

In a Freshman class the probability that a randomly selected student passed the mathematics exam is 0.72, passed the microeconomics exam is 0.66, and passed both is 0.54. Find the probability that a randomly selected student

- (a) passed at least one of those exams,
- (b) passed the microeconomics exam, but did not pass the mathematics exam,
- (c) passed none of the exams.

Let A denote the event that a randomly selected student passed the mathematics exam, and B is the event that the student passed the microeconomics

exam. Then $P(A) = 0.72$, $P(B) = 0.66$ and $P(A \cap B) = 0.54$. Using the events A and B , the desired probabilities can be given the following way.

$$(a) P(A \cup B) = P(A) + P(B) - P(A \cap B) = 0.84$$

$$(b) P(\overline{A} \cap B) = P(B) - P(A \cap B) = 0.12$$

$$(c) P(\overline{A} \cap \overline{B}) = P(\overline{A \cup B}) = 1 - P(A \cup B) = 0.16$$

Recitation and Exercises

1. Reading: Textbook-2, Sections 2.1, 2.2, 2.3, 2.4, 2.5.
2. Homework: Textbook-2, Exercises 2.11, 2.19, 2.32, 2.33, 2.37, 2.38, 2.54, 2.58, 2.59, 2.61, 2.110, 2.112.
3. Review: Highschool Combinatorics and Binomial Theorem (Textbook-2, Section 2.3), and "Probability Exercises"

Chapter 14

Sampling methods

14.1 Classical probability spaces

Definition 14.1 Consider a probability space (Ω, \mathcal{A}, P) . It is called a *classical probability space*, if

- Ω is a finite set,
- for every $\omega \in \Omega$ we have $\{\omega\} \in \mathcal{A}$,
- every singleton subset of Ω has the same probability.

Obviously, if Ω contains exactly n elements, then for every $\omega \in \Omega$ we get

$$P(\{\omega\}) = \frac{1}{n}$$

In particular, if the event $A \subset \Omega$ consists of k elements, then

$$P(A) = \frac{k}{n}$$

This observation can be interpreted as the probability of A can be given like:

$$P(A) = \frac{\text{number of favorable outcomes}}{\text{total number of outcomes}} \quad (14.1)$$

The formula (14.1) will be called the classical formula.

Example 14.2 A regular playing die is tossed twice in a row. What is the probability that the sum of the two numbers is exactly 7?

Let A denote the event that the sum is 7. Clearly, the sample space Ω contains 36 elements (total number of outcomes), while A is a subset of 6 elements containing the pairs (1, 6), (6, 1), (2, 5), (5, 2), (3, 4), (4, 3) (favorable outcomes). Consequently

$$P(A) = \frac{6}{36} = \frac{1}{6}$$

by making use of the classical formula (14.1).

Example 14.3 From a deck of 52 playing cards we draw 5 cards at random. Find the probability that either all 5 cards are clubs, or at least one of them is an Ace?

Introduce the following notations:

$$A = \{\text{all 5 cards are clubs}\} \quad B = \{\text{at least one of them is Ace}\}$$

Obviously we are looking for $P(A \cup B)$. Since the draws of any 5 cards are equally likely, therefore:

$$P(A) = \frac{\binom{13}{5}}{\binom{52}{5}} \quad P(B) = 1 - \frac{\binom{48}{5}}{\binom{52}{5}}$$

and further:

$$P(A \cap B) = \frac{\binom{12}{4}}{\binom{52}{5}}$$

By using the additive rule

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

Example 14.4 On a seasonal sale in a supermarket there are 10 different pairs of shoes in a basket. A thief quickly grabs 4 pieces of shoes from the basket at random and runs away. What is the probability that he gets at least 1 complete pair?

Below we outline two approaches, but only one of them is correct.

- First select one pair, the other two pieces of shoes can be taken arbitrarily, another pair, or any two of the remaining shoes, i.e.:

$$\frac{10 \binom{18}{2}}{\binom{20}{4}}$$

- Find the probability of not having a complete pair at all. This can be done by selecting a single shoe, and then putting its matching pair aside. Keep in mind that the order of the selection does not count. Then passing to the complement event, we obtain

$$1 - \frac{\frac{20 \cdot 18 \cdot 16 \cdot 14}{4!}}{\binom{20}{4}}$$

Check out that the two probabilities do not coincide! Which one is correct (if any)?

Example 14.5 Keep tossing a die until 6 comes out for the first time. What is the probability that we need an even number of tosses?

Let A stand for the event that we need an even number of tosses and A_k is the event that we need k tosses, respectively. Then we have (verify!)

$$P(A_k) = \left(\frac{5}{6}\right)^{k-1} \cdot \frac{1}{6}$$

for every $k = 1, 2, \dots$. The event A can be expressed like this:

$$A = A_2 \cup A_4 \cup \dots = \bigcup_{k=1}^{\infty} A_{2k}$$

On the right hand side the events mutually exclude each other, hence

$$P(A) = \sum_{k=1}^{\infty} P(A_{2k}) = \sum_{k=1}^{\infty} \left(\frac{5}{6}\right)^{2k-1} \cdot \frac{1}{6} = \frac{5}{11}$$

14.2 Sampling without replacement

Consider a set of N objects so that m of them are defective. Select a sample of n objects from the whole set at random, without replacement ($n \leq m$). Denote by A_k the event, that the sample contains exactly k defective objects ($0 \leq k \leq n$). Then

$$P(A_k) = \frac{\binom{m}{k} \cdot \binom{N-m}{n-k}}{\binom{N}{n}}$$

which we call the formula of sampling without replacement.

Example 14.6 From a deck of 52 playing cards we draw 5 cards at random without replacement. Find the probability that we selected exactly 2 diamonds.

Let A denote the given event. Making use of our formula we get

$$P(A) = \frac{\binom{13}{2} \cdot \binom{39}{3}}{\binom{52}{5}}.$$

In this argument the diamonds are the "defective objects".

Example 14.7 Determine the probability that in Hungarian lottery (5 winners out of 90) we have at least 2 winning numbers on a lottery ticket filled in at random.

Denote by A the event that we have 2 winning numbers, and by A_k the event that we have exactly k winning numbers on our ticket. Clearly, the events A_k are mutually exclusive for $k = 2, \dots, 5$. On the other hand $A = A_2 \cup A_3 \cup A_4 \cup A_5$, and this implies

$$P(A) = \sum_{k=2}^5 P(A_k) = \sum_{k=2}^5 \frac{\binom{5}{k} \cdot \binom{85}{5-k}}{\binom{90}{5}}$$

since the probability of the disjoint union is the sum of the probabilities.

Example 14.8 From a deck of 52 playing cards we select 5 cards at random, without replacement. What is the probability that all 4 suits (clubs, diamonds, hearts, spades) are represented in the sample?

Examine the following argument. Let A denote the event that all 4 suits appear in the sample of 5 cards. Since the choice of any 5 cards is equally likely, we deal with a classical probability space.

In order to find out the number of favorable outcomes, take into account that we have 13 options for each suit. Once one card from each suit has been taken, then any card can be chosen from the remaining 48 cards.

The total number of outcomes: as many as the number of selections of 5 cards out of 52. So:

$$P(A) = \frac{13^4 \cdot 48}{\binom{52}{5}}$$

Is this the correct solution? If not, how could it be fixed?

14.3 Sampling with replacement

Consider again the set of N objects so that m of them are defective. Select n objects at random from the whole set, consecutively one after another with replacement. Let A_k denote the event that the sample contains exactly k defective items.

Examine the draws of different orders. Since the selection of k defectives and $n - k$ non-defectives in any order admits the probability

$$\frac{m^k \cdot (N - m)^{n-k}}{N^n} = \left(\frac{m}{N}\right)^k \left(1 - \frac{m}{N}\right)^{n-k}$$

and we have exactly $\binom{n}{k}$ options for such selections, moreover they mutually exclude each other, we receive

$$P(A_k) = \binom{n}{k} \left(\frac{m}{N}\right)^k \left(1 - \frac{m}{N}\right)^{n-k}$$

This equality is called the formula of sampling with replacement.

Example 14.9 Take 5 cards out of a deck of 52 cards at random, successively with replacement. (The card taken at a time is always put back.) Find the probability that this way

- (a) exactly 2 diamonds are selected,
- (b) at least 2 diamonds are selected.

Introduce the event A_k which means that exactly k diamonds are selected. Then

$$(a) \quad P(A_2) = \binom{5}{2} \left(\frac{1}{4}\right)^2 \left(1 - \frac{1}{4}\right)^3$$

and

$$(b) \quad P(A_2 \cup \dots \cup A_5) = \sum_{k=2}^5 \binom{5}{k} \left(\frac{1}{4}\right)^k \left(\frac{3}{4}\right)^{5-k}$$

because the events A_2, \dots, A_5 are mutually exclusive.

14.4 The Bernoulli experiment

The argument above can be generalized the following way. Suppose that the probability of an event A in a given experiment is a specific number $0 \leq p \leq 1$.

Let us assume that we carry out this experiment n times in a row (independently of each other) and every time we observe whether or not A occurs. This procedure is called the Bernoulli experiment.

Let $0 \leq k \leq n$ be a given integer. Denote by A_k the event that A occurs exactly k times out of the n trials.

Following the reasoning, analogous to the previous section, we immediately get

$$P(A_k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

for every integer $k = 0, 1, \dots, n$.

Example 14.10 In the Hungarian lottery we say that a lottery ticket is a winning ticket, if it contains at least two winning numbers. Suppose we purchase 20 tickets and fill in them at random (independently of each other). Find the probability that we will have at least 5 winning tickets.

For just one ticket the probability of being a winning ticket is:

$$p = \sum_{k=2}^5 \frac{\binom{5}{k} \cdot \binom{85}{5-k}}{\binom{90}{5}}$$

Since this is true for every ticket, and the tickets are filled in independently from each other, this problem can be regarded as a Bernoulli experiment, with the parameter p specified above. Therefore, applying our formula:

$$\sum_{k=5}^{20} \binom{20}{k} p^k (1-p)^{20-k}$$

where p is the probability given above.

Recitation and Exercises

1. Reading: Textbook-2, Sections 2.1, 2.2, 2.3, 2.4, 2.5.
2. Homework: Textbook-2, Exercises 2.20, 2.39, 2.42, 2.48, 2.64, 2.71, 2.72, 2.113, 2.114, 2.115, 2.116.
3. Review: Highschool Combinatorics and Binomial Theorem (Textbook-2, Section 2.3), and "Probability Exercises"

Chapter 15

Conditional probability and Bayes' Rule

15.1 Conditional probability

In several problems we need to find the probability of the event A under the a priori condition that a certain event B occurred. In such problems we take into account only those elements of the sample space, which also belong to B .

This actually means that the sample space Ω is reduced to the subset B , and we calculate the (conditional) probability of A with respect to B .

Definition 15.1 Consider the probability space (Ω, \mathcal{A}, P) and an event $B \in \mathcal{A}$ so that $P(B) \neq 0$. The conditional probability of the event $A \in \mathcal{A}$ with respect to B (read: probability of A given B) is defined by the equality:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Example 15.2 We toss a pair of dice, but we cannot see the outcome. Someone tells us that one of them is a 5. What is the probability that other one is 6?

ATTENTION! The answer is not $1/6$ for the following reason!

Let A and B denote the following events:

$$B = \{\text{one of the tosses is 5}\} \quad A = \{\text{the other one is 6}\}$$

On the one hand $P(B) = 11/36$ since there are 11 pairs that contain 5. On the other hand $A \cap B = \{(5, 6), (6, 5)\}$, and hence $P(A \cap B) = 2/36$. Therefore:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{2/36}{11/36} = \frac{2}{11}$$

Example 15.3 We are looking for a friend in the university main building. He can be in 5 rooms equally likely. The probability that he is in fact in the building is $0 < p < 1$. We have checked 4 of the 5 rooms, and he was in none of them. What is the probability that he is in the fifth room?

Let A_k denote the event that our friend is in room number k ($k = 1, \dots, 5$), which means $P(A_1 \cup \dots \cup A_5) = p$. Since the events A_k are mutually exclusive, this implies that $P(A_k) = p/5$ for every index k . Therefore, in view of the De Morgan Rule we obtain:

$$\begin{aligned} P(A_5 | \overline{A_1} \cap \dots \cap \overline{A_4}) &= P(A_5 | \overline{A_1 \cup \dots \cup A_4}) \\ &= \frac{P(A_5 \cap (\overline{A_1 \cup \dots \cup A_4}))}{P(\overline{A_1 \cup \dots \cup A_4})} \end{aligned}$$

Obviously (think about it!):

$$A_5 \subset \overline{A_1 \cup \dots \cup A_4}$$

and hence

$$P(A_5 \cap (\overline{A_1 \cup \dots \cup A_4})) = P(A_5)$$

Consequently, the desired conditional probability is:

$$\begin{aligned} P(A_5 | \overline{A_1} \cap \dots \cap \overline{A_4}) &= \frac{P(A_5 | \overline{A_1 \cup \dots \cup A_4})}{P(\overline{A_1 \cup \dots \cup A_4})} \\ &= \frac{P(A_5 \cap (\overline{A_1 \cup \dots \cup A_4}))}{P(\overline{A_1 \cup \dots \cup A_4})} \\ &= \frac{P(A_5)}{P(\overline{A_1 \cup \dots \cup A_4})} = \frac{p/5}{1 - 4p/5} = \frac{p}{5 - 4p} \end{aligned}$$

15.2 Independence

Consider the following simple example. Toss a die twice in a row, and we cannot see the result. Someone tells us that the first outcome is an odd number. Find the probability that the sum of the two numbers is 7.

Introduce the events A and B the following way:

$$A = \{\text{the sum is 7}\} \quad B = \{\text{the first outcome is odd}\}$$

Then, by the definition of the conditional probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{3/36}{18/36} = \frac{1}{6}$$

In view of one of a previous example this means

$$P(A|B) = P(A)$$

that is "the occurrence of B has no impact on the probability of A ". This fact is expressed like "the event A is independent of the event B ".

In case of $P(B) \neq 0$ the condition $P(A|B) = P(A)$ is equivalent to the equality:

$$P(A \cap B) = P(A) \cdot P(B) \quad (15.1)$$

Since we figure that independence is a symmetric relation (i.e. if A is independent of B , then B is also independent of A) and the above equality is visibly symmetric, relation (15.1) can serve as a comfortable definition for independence.

Definition 15.4 Let (Ω, \mathcal{A}, P) be a probability space, and $A, B \in \mathcal{A}$ are observable events. We say that A and B are *independent*, if they fulfill the condition (15.1).

Example 15.5 From a deck of 52 cards we draw 2 cards in succession with replacement. Find the probability that the first draw is a diamond, and the second draw is an Ace.

Introduce the following events:

$$A = \{\text{first draw is a diamond}\} \quad B = \{\text{second draw is an Ace}\}$$

Then

$$P(A \cap B) = \frac{13 \cdot 4}{52^2} = \frac{13}{52} \cdot \frac{4}{52} = P(A) \cdot P(B)$$

that tells us that the events A and B are independent.

ATTENTION! We NEVER argue like: since the events A and B are "visibly" independent, therefore $P(A \cap B) = P(A) \cdot P(B)$. On the contrary: we conclude the independence of events by verifying this equality!

15.3 Theorem of Total Probability

Example 15.6 There are 3 identical envelopes on our desk,

1. the first contains 2 of 1000 Ft bills and 3 of 2000 Ft bills (banknotes),
2. the second contains 5 of 1000 Ft bills and 2 of 2000 Ft bills,
3. the third contains 5 of 2000 Ft bills.

We select one of the envelopes at random and draw one of the bills from the envelope. What is the probability that we take a 2000 Ft bill?

Let A denote the event that we draw a 2000 Ft bill. The probability $P(A)$ would be easy to determine if we knew, which envelope is selected. In particular, if B_k stands for the event that envelope k is selected, then the conditional probabilities $P(A|B_k)$ are $3/5$, $2/7$ and 1 respectively.

This observation immediately gives an idea of how to solve the problem. The events B_k are mutually exclusive and their union is the certain event. Thus:

$$A = A \cap \Omega = A \cap (B_1 \cup B_2 \cup B_3) = (A \cap B_1) \cup (A \cap B_2) \cup (A \cap B_3)$$

Since the events on the right-hand side are exclusive:

$$\begin{aligned} P(A) &= P(A \cap B_1) + P(A \cap B_2) + P(A \cap B_3) \\ &= P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + P(A|B_3)P(B_3) \\ &= \frac{3}{5} \cdot \frac{1}{3} + \frac{2}{7} \cdot \frac{1}{3} + 1 \cdot \frac{1}{3} \end{aligned}$$

The argument above can be extended to an arbitrary number of events B_k . This leads us to the following definition.

Definition 15.7 We say that the observable events $B_1, B_2, \dots \in \mathcal{A}$ form a *partition* of the sample space, if none of them has probability zero, and further

1. they are mutually exclusive, i.e. $B_i \cap B_j = \emptyset$ if $i \neq j$,
2. one of them occurs, i.e. $B_1 \cup B_2 \cup \dots = \Omega$.

Following the analogous argument of Example 15.6 for an arbitrary number of events B_k , we come up with the following theorem.

Theorem 15.8 (Theorem of Total Probability) *Let us suppose that in the probability space (Ω, \mathcal{A}, P) the events B_1, B_2, \dots form a partition of the sample space. Then for any event $A \in \mathcal{A}$ we have*

$$P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots$$

Proof. Indeed, if the events B_k form a partition of the sample space, then

$$A = (A \cap B_1) \cup (A \cap B_2) \cup (A \cap B_3) \cup \dots$$

where the terms of the union are mutually exclusive. Thus:

$$P(A) = P(A \cap B_1) + P(A \cap B_2) + P(A \cap B_3) + \dots$$

By the very definition of the conditional probability, for every index k

$$P(A \cap B_k) = P(A|B_k) \cdot P(B_k)$$

and the theorem ensues. \square

Example 15.9 If the probability that the number of incoming calls to a call center is n on a given day is given by $0 < q_n < 1$, and every call is a wrong number with probability $0 < p < 1$ (independently of each other), find the probability that the number of wrong calls is exactly k on that day.

Introduce the following notations. Let A be the event that the center receives k wrong calls, and B_n is the event that the total number of incoming calls is n . In this case the events B_n form a partition of the sample space, hence by the theorem of total probability

$$P(A) = \sum_{n=1}^{\infty} P(A|B_n) \cdot P(B_n) = \sum_{n=k}^{\infty} q_n \binom{n}{k} p^k (1-p)^{n-k}$$

In fact, for $n \geq k$ the number of wrong calls can be regarded as the outcome of a Bernoulli experiment: how many wrong calls do we have out of n incoming calls. Keep in mind that we have $P(A|B_n) = 0$, for $n < k$.

15.4 Bayes' Rule

Let us return to Example 15.6. Assume that someone has performed the draw (we did not see it) and tells us that the draw is a 2000 Ft bill. What is the probability that the bill was taken from the first envelope?

Using our former notations, we need to find the conditional probability $P(B_1|A)$.

$$P(B_1|A) = \frac{P(A \cap B_1)}{P(A)} = \frac{P(A|B_1)P(B_1)}{P(A)}$$

The denominator of the fraction on the right-hand side can be evaluated by the theorem of total probability:

$$\begin{aligned} P(B_1|A) &= \frac{P(A|B_1)P(B_1)}{P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + P(A|B_3)P(B_3)} \\ &= \frac{\frac{3}{5} \cdot \frac{1}{3}}{\frac{3}{5} \cdot \frac{1}{3} + \frac{2}{7} \cdot \frac{1}{3} + 1 \cdot \frac{1}{3}} \end{aligned}$$

This argument can be extended to any partition of the sample space.

Theorem 15.10 (Bayes' Rule) *Let us suppose that in the probability space (Ω, \mathcal{A}, P) the events B_1, B_2, \dots form a partition of the sample space. Then for any event $A \in \mathcal{A}$, $P(A) \neq 0$ and any index i we have*

$$P(B_i|A) = \frac{P(A|B_i)P(B_i)}{P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots}$$

Proof. Indeed, by the definition of the conditional probability

$$P(B_i|A) = \frac{P(A \cap B_i)}{P(A)} = \frac{P(A|B_i)P(B_i)}{P(A)},$$

and our statement is proven by applying the theorem of total probability. \square

Example 15.11 For instance, in our call center Example 15.9 the probability that the number of incoming calls on a given day is i provided that exactly k wrong calls have been registered is

$$P(B_i|A) = \frac{q_i \binom{i}{k} p^k (1-p)^{i-k}}{\sum_{n=k}^{\infty} q_n \binom{n}{k} p^k (1-p)^{n-k}}$$

for $i \geq k$, while this probability is 0, for $i < k$.

Recitation and Exercises

1. Reading: Textbook-2, Sections 2.6 and 2.7
2. Homework: Textbook-2, Exercises 2.80, 2.81, 2.87, 2.95, 2.97, 2.100, 2.109, 2.118
3. Review: Highschool Combinatorics and Binomial Theorem (Textbook-2, Section 2.3), and "Probability Exercises"

Chapter 16

Random variables and distributions

16.1 Random variables

Definition 16.1 Consider a probability space (Ω, \mathcal{A}, P) . The function

$$X : \Omega \rightarrow \mathbb{R}$$

is called *random variable*, if for any $x \in \mathbb{R}$

$$\{X < x\} = \{\omega \in \Omega : X(\omega) < x\} \in \mathcal{A}$$

that is all level sets are observable (and hence possess a probability).

In the examples below specify the range R of the given random variables!

Example 16.2

1. Toss a pair of dice. Let X denote the sum of the numbers. Then $R = \{2, 3, \dots, 12\}$
2. Let X be the least winning number in Hungarian lottery. Then $R = \{1, 2, \dots, 86\}$
3. Keep tossing a die until 6 comes out for the first time. Denote by X the number of tosses. Then $R = \mathbb{N}$.
4. Pick a point arbitrarily on the unit disc (with center at the origin and radius 1). Let X denote the distance of the point from the origin. Then $R = [0, 1]$.

Definition 16.3 We say that a random variable is *discrete*, if its range is a countable set (finite or infinite). That is the elements of the range can be arranged in a finite or infinite sequence.

In our examples the first three random variables are discrete, but the fourth is not.

16.2 Distribution of discrete variables

Definition 16.4 Let X be a discrete random variable, whose range is $R = \{x_1, x_2, \dots\}$. The sequence

$$p_k = P(X = x_k), \quad k = 1, 2, \dots$$

is called the *distribution* of X .

Example 16.5 Consider our introductory examples for random variables

1. If X means the sum of the numbers when a pair of dice tossed, then the distribution can be given by the following *chart*:

x_k	2	3	4	...	12
p_k	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$...	$\frac{1}{36}$

2. If X means the least winning number in lottery, then the distribution can be given by the following *formula*:

$$p_k = \frac{\binom{90-k}{4}}{\binom{90}{5}} \quad k = 1, 2, \dots, 86$$

3. If X means the number of tosses needed to get the first 6, the distribution of X is:

$$p_k = \left(\frac{5}{6}\right)^{k-1} \cdot \frac{1}{6} \quad k = 1, 2, \dots$$

Unlike in the previous two examples, this distribution is an infinite sequence.

The most important properties of distributions are summed up in the following theorem.

Theorem 16.6 Consider a discrete random variable X with range $R = \{x_1, x_2, \dots\}$ and distribution $p_k = P(X = x_k)$, $k = 1, 2, \dots$. Then

- $0 \leq p_k \leq 1$ for all indices $k = 1, 2, \dots$
- $p_1 + p_2 + \dots = 1$.
- If $a < b$ any real numbers, then

$$P(a < X < b) = \sum_{a < x_k < b} p_k$$

where the sum is taken for all indices k such that the inequality $a < x_k < b$ holds true. The last statement remains true if instead of the strict inequalities, the signs \leq are inserted simultaneously on both sides.

16.3 The cumulative distribution function

Definition 16.7 Consider a probability space (Ω, \mathcal{A}, P) , and a random variable $X : \Omega \rightarrow \mathbb{R}$. For every $x \in \mathbb{R}$ set

$$F(x) = P(X < x).$$

The function $F : \mathbb{R} \rightarrow [0, 1]$ is called the *cumulative distribution function* of X . (Or sometimes briefly *distribution function*.)

Example 16.8 It is easy to see that the distribution function of the random variable X defined in the introductory example 4, is

$$F(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ x^2 & \text{if } 0 < x \leq 1 \\ 1 & \text{if } x > 1 \end{cases} \quad (16.1)$$

In fact we mean that the probability that the randomly picked point belongs to a given subset of the unit disc is proportional to the area of the subset. In particular, for instance $P(0 \leq X < 1/2) = 1/4$.

In several problems in probability and statistics, and their applications we need to find a probability of the form $P(a \leq X < b)$. This probability can be expressed in term of the distribution function. The basic properties of the distribution function are summarized in the theorem below.

Theorem 16.9 Let X be a random variable and consider its distribution function F .

- For every $x \in \mathbb{R}$ we have $0 \leq F(x) \leq 1$.

- F is monotone increasing and at every point continuous from the left.
- $\lim_{x \rightarrow -\infty} F(x) = 0$ and $\lim_{x \rightarrow +\infty} F(x) = 1$.
- For any real numbers $a < b$ we have

$$P(a \leq X < b) = F(b) - F(a).$$

If the range of a discrete random variable X is given by $R = \{x_1, x_2, \dots\}$, where $x_1 < x_2 < \dots$, and X takes these values with the probabilities p_1, p_2, \dots respectively, then the distribution function of X has the form:

$$F(x) = \begin{cases} 0 & \text{if } x \leq x_1 \\ p_1 + \dots + p_k & \text{if } x_k < x \leq x_{k+1} \end{cases}$$

for each $k = 1, 2, \dots$. Sketch the graph!

This tells us that in this case the distribution function is piecewise constant. Instead of using the formula $P(a \leq X < b) = F(b) - F(a)$, it is reasonable to collect all elements of the range of X that are in the open interval (a, b) . In particular, if $P(X = x_k) = p_k$ for every k , then

$$P(a \leq X < b) = \sum_{a \leq x_k < b} p_k$$

On the right-hand side only the probabilities $P(X = x_k)$ appear, therefore, it is more convenient to rely on the distribution X .

16.4 The density function

Definition 16.10 We say that X is *continuously distributed*, if there exists an integrable function f on the real line with

$$F(x) = \int_{-\infty}^x f(t) dt$$

for every $x \in \mathbb{R}$. In this case the function f is called the *density function* of X .

For instance, in the example (16.1) we can easily verify that

$$f(t) = \begin{cases} 2t & \text{if } 0 < t < 1 \\ 0 & \text{elsewhere} \end{cases}$$

If the random variable X is continuously distributed, then the distribution function F is continuous. Moreover, at every point x where the density function f is continuous, the distribution function F is differentiable, namely

$$F'(x) = f(x)$$

Theorem 16.11 *If X is continuously distributed and f is its density function, then for any real numbers $a < b$*

$$P(a \leq X < b) = \int_a^b f(t) dt$$

What can we say about the probability that the random variable X takes a single point? Let $a \in \mathbb{R}$ be any real number, then we conclude that

$$\begin{aligned} P(X = a) &= P\left(\bigcap_{n=1}^{\infty} \left\{a \leq X < a + \frac{1}{n}\right\}\right) = \lim_{n \rightarrow \infty} P\left(a \leq X < a + \frac{1}{n}\right) \\ &= \lim_{n \rightarrow \infty} \left(F\left(a + \frac{1}{n}\right) - F(a)\right) = \lim_{x \rightarrow a^+} F(x) - F(a) \end{aligned}$$

Consequently $P(X = a)$ equals the "jump" of F at the point a . **ATTENTION:** Why can we pass to the limit in the first line of the array formula?

A simple consequence of the previous argument is that $P(X = a) = 0$ if and only if F is continuous at the point a . In particular, if X is continuously distributed, then F is continuous on the whole real line, hence for any real numbers $a < b$ we get

$$P(a < X < b) = P(a \leq X \leq b)$$

We sum up the basic properties of density functions.

Theorem 16.12 *If f is the density function of the random variable X , then*

1. $f(x) \geq 0$ for every $x \in \mathbb{R}$,
- 2.

$$\int_{-\infty}^{+\infty} f(x) dx = 1,$$

3. if $a < b$ are any real numbers, then

$$P(a < X < b) = P(a \leq X \leq b) = \int_a^b f(x) dx.$$

Example 16.13 Let us suppose that the density function of X is given by

$$f(x) = \begin{cases} x & \text{if } 0 < x \leq 1 \\ 2 - x & \text{if } 1 < x < 2 \\ 0 & \text{elsewhere} \end{cases}$$

ATTENTION! Verify that f fulfills all conditions of the previous theorem, so it is in fact a density function.

Then, for instance

$$\begin{aligned} P(0 \leq X \leq 3/2) &= P(0 < X < 3/2) = \int_0^{3/2} f(x) dx \\ &= \int_0^1 x dx + \int_1^{3/2} (2-x) dx \\ &= 1 - \int_{3/2}^2 (2-x) dx = \frac{7}{8} \end{aligned}$$

Recitation and Exercises

1. Reading: Textbook-2, Sections 3.1, 3.2 and 3.3
2. Homework: Textbook-2, Exercises 3.7, 3.9, 3.11, 3.14, 3.21, 3.22, 3.25, 3.26, 3.32 and 3.36
3. Review: Calculus, integration and infinite series and "Probability Exercises"

Chapter 17

Mean and variance

In everyday language by the mean (or expected value) of a random variable we think of the weighted average, by the standard deviation we think of the average deviation from the mean. Precise definitions will follow below.

17.1 Mean of discrete distributions

Definition 17.1 Consider a discrete random variable X whose distribution is given by

$$P(X = x_k) = p_k \quad k = 1, 2, \dots$$

We say that X has a mean (or expected value) if the series $\sum_{k=1}^{\infty} |x_k| \cdot p_k$ is convergent. In this case the sum

$$E(X) = \sum_{k=1}^{\infty} x_k \cdot p_k$$

is called the *mean* (or *expected value*) of X .

Remark that the convergence of the series $\sum_{k=1}^{\infty} |x_k| \cdot p_k$ is an important condition, because otherwise the sum $E(X)$ might depend on the rearrangement of the terms.

Example 17.2 Toss a pair of playing dice. Find the expected value of the sum of the two numbers.

Let X denote the sum of the two numbers, then the distribution of X is given in Example 16.5. Therefore, the mean of the sum is:

$$E(X) = \sum_{k=2}^{12} kp_k = 2 \cdot \frac{1}{36} + 3 \cdot \frac{2}{36} + \dots + 12 \cdot \frac{1}{36} = 7.$$

Example 17.3 Take a sample of 5 cards from a deck of 52 playing cards at random. Find the expected number of diamonds in the sample.

Denote by X the number of diamonds in the sample. By using sampling without replacement, the distribution of X is given by:

$$P(X = k) = \frac{\binom{13}{k} \cdot \binom{39}{5-k}}{\binom{52}{5}} \quad k = 0, \dots, 5$$

Hence, the expected value is:

$$\begin{aligned} E(X) &= \sum_{k=0}^5 k P(X = k) = \sum_{k=0}^5 k \frac{\binom{13}{k} \cdot \binom{39}{5-k}}{\binom{52}{5}} \\ &= \frac{13}{\binom{52}{5}} \sum_{k=1}^5 \binom{12}{k-1} \binom{39}{4-(k-1)} = \frac{13}{\binom{52}{5}} \cdot \binom{51}{4} = \frac{5}{4}. \end{aligned}$$

Example 17.4 Consider the Bernoulli experiment that we discussed in Section 14.4. and determine the expected number of occurrences of the event A out of n trials.

Let X denote the number of times A occurs, then the distribution of X is:

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k} \quad k = 0, 1, \dots, n$$

By virtue of the binomial theorem, the mean of X is:

$$\begin{aligned} E(X) &= \sum_{k=0}^n k \binom{n}{k} p^k (1-p)^{n-k} \\ &= np \sum_{k=1}^n \frac{(n-1)!}{(k-1)!(n-k)!} p^{k-1} (1-p)^{n-k} \\ &= np \sum_{k=1}^n \binom{n-1}{k-1} p^{k-1} (1-p)^{n-k} = np \end{aligned}$$

17.2 Mean of infinite distributions

In this section we investigate discrete random variables with infinite range.

Example 17.5 We keep tossing a die until 6 comes out for the first time. What is the expected number of tosses?

If X means the number of tosses, then the distribution of X is given by

$$P(X = k) = \left(\frac{5}{6}\right)^{k-1} \cdot \frac{1}{6} \quad k = 1, 2, \dots$$

Thus the expected value is

$$E(X) = \sum_{k=1}^{\infty} k \cdot \left(\frac{5}{6}\right)^{k-1} \cdot \frac{1}{6} = \frac{1}{6} \cdot \frac{1}{(1 - 5/6)^2} = 6$$

Example 17.6 Let λ be a given positive number, and consider a random variable X with the following distribution

$$P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda} \quad k = 0, 1, 2, \dots$$

In view of the power series of the exponential function, the mean of X is:

$$\begin{aligned} E(X) &= \sum_{k=0}^{\infty} k \frac{\lambda^k}{k!} e^{-\lambda} = \lambda \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} e^{-\lambda} \\ &= \lambda e^{-\lambda} \sum_{i=0}^{\infty} \frac{\lambda^i}{i!} = \lambda e^{-\lambda} e^{\lambda} = \lambda \end{aligned}$$

Example 17.7 In a box there is a black and a white ball. We take one ball at random. If it is black, we put it back, and add another black ball. We continue this process until the white ball is selected. Find the expected number of draws.

If X stands for the number of draws, then the distribution of X can be given like $P(X = 1) = 1/2$, and:

$$P(X = k) = \frac{1}{2} \cdot \frac{2}{3} \cdot \frac{3}{4} \cdots \frac{k-1}{k} \cdot \frac{1}{k+1} = \frac{1}{k(k+1)}, \quad k = 2, 3, \dots$$

Therefore, for the mean of X we obtain the following infinite series:

$$E(X) = \sum_{k=1}^{\infty} k P(X = k) = \sum_{k=1}^{\infty} k \frac{1}{k(k+1)} = \sum_{k=1}^{\infty} \frac{1}{k+1}$$

Apart from the first term, this series exactly coincides with the harmonic series, which is divergent. Consequently, this random variable does not have a mean.

17.3 Mean of continuous distributions

Definition 17.8 Let X be a continuously distributed random variable with density function f . We say that X has a mean if the improper integral $\int_{-\infty}^{\infty} |x| \cdot f(x) dx$ is convergent. In this case the integral

$$E(X) = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

is called the *mean* (or expected value) of X .

Example 17.9 Verify that the function f below defines a density function

$$f(x) = \frac{1}{\pi} \cdot \frac{1}{1+x^2} \quad -\infty < x < \infty$$

(this is the so-called Cauchy distribution), but it has no mean, since the improper integral

$$\frac{1}{\pi} \int_{-\infty}^{\infty} \frac{x}{1+x^2} dx$$

is divergent. See Example 9.5 for the details.

Example 17.10 Consider an interval $[a, b]$ on the real line, and suppose the density function of the random variable X is given by

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{elsewhere} \end{cases}$$

Verify that f is really a density function! Then the mean of X is

$$E(X) = \int_a^b \frac{x}{b-a} dx = \frac{1}{b-a} \cdot \frac{b^2 - a^2}{2} = \frac{a+b}{2}$$

which is the midpoint of the interval $[a, b]$.

17.4 Basic properties of the mean

The mean $E(X^2)$ is called the second moment of the random variable X (if it exists). It can be shown that

$$E(X^2) = \begin{cases} \sum x_k^2 p_k & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} x^2 f(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

Below two fundamental properties of the mean are formulated.

Theorem 17.11

1. If X has a mean, then for any real numbers α and β $E(\alpha X + \beta) = \alpha E(X) + \beta$.
2. If $E(X)$, $E(X^2)$ exist, then $E(\alpha X^2 + \beta X + \gamma) = \alpha E(X^2) + \beta E(X) + \gamma$.

Example 17.12 Let λ be a positive number, and assume that the density function of X is given by

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x > 0 \\ 0 & \text{elsewhere.} \end{cases}$$

Based on Example 9.3 this is really a density function, since

$$\int_0^{\infty} f(x) dx = 1.$$

On the other hand, Example 9.8 shows that the mean is

$$E(X) = \int_0^{\infty} x f(x) dx = \frac{1}{\lambda}.$$

The second moment can be evaluated by integration by parts (see Example 9.9):

$$E(X^2) = \int_0^{\infty} x^2 f(x) dx = \frac{2}{\lambda^2}.$$

17.5 Variance and standard deviation

The variance of a random variable is the average squared deviation from the mean.

Definition 17.13 The *variance* of a random variable of X (if it exists) is defined by

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

Then the *standard deviation* of X is $D(X) = \sqrt{\text{Var}(X)}$.

Sometimes the notation $D^2(X)$ is also used for the variance (for obvious reason).

The variance can be evaluated in the following simplified way:

$$\begin{aligned} \text{Var}(X) &= E((X - E(X))^2) = E(X^2 - 2E(X)X + E(X)^2) \\ &= E(X^2) - 2E(X)^2 + E(X)^2 = E(X^2) - E(X)^2 \end{aligned}$$

Basic properties of the variance:

$$\text{Var}(\alpha X + \beta) = \alpha^2 \text{Var}(X), \quad D(\alpha X + \beta) = |\alpha| \cdot D(X)$$

Verify these two directly, based on the definition!

Example 17.14 Find the variance and standard deviation of the continuously distributed random variable X in Example 17.12 (where $\lambda > 0$ is a given constant).

$$\text{Var}(X) = E(X^2) - E(X)^2 = \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{1}{\lambda^2},$$

and in particular

$$D(X) = \frac{1}{\lambda}$$

Example 17.15 Consider now the continuously distributed random variable X examined in Example 17.10. We can calculate the second moment this way:

$$\begin{aligned} E(X^2) &= \int_{-\infty}^{\infty} x^2 f(x) dx = \int_a^b \frac{x^2}{b-a} dx = \frac{1}{b-a} \left[\frac{x^3}{3} \right]_a^b \\ &= \frac{b^3 - a^3}{3(b-a)} = \frac{b^2 + ab + a^2}{3} \end{aligned}$$

Therefore, the variance is:

$$\text{Var}(X) = E(X^2) - E(X)^2 = \frac{b^2 + ab + a^2}{3} - \frac{a^2 + 2ab + b^2}{4} = \frac{(b-a)^2}{12}$$

moreover, the standard deviation of X is the square root of the variance:

$$D(X) = \frac{b-a}{2\sqrt{3}}.$$

Recitation and Exercises

1. Reading: Textbook-2, Sections 4.1 and 4.2.
2. Homework: Textbook-2, Exercises 4.1, 4.2, 4.4, 4.8, 4.12, 4.13, 4.14, 4.34, 4.37, 4.38, 4.43 and 4.50
3. Review: Calculus, integration, improper integrals and infinite series, and "Probability Exercises"

Chapter 18

Special discrete distributions

This chapter gives a summary of the most widely applied discrete distributions.

18.1 Characteristic distribution

Let (Ω, \mathcal{A}, P) be a probability space and consider an event $A \in \mathcal{A}$ with $P(A) = p$, and $0 < p < 1$. Then the random variable

$$X = \begin{cases} 1 & \text{if } A \text{ occurs} \\ 0 & \text{if } A \text{ does not occur} \end{cases}$$

possesses the distribution

$$P(X = 0) = 1 - p \quad P(X = 1) = p$$

This is called the *characteristic distribution* associated with the event A .

Theorem 18.1

- *The parameter of the distribution is: $0 < p < 1$.*
- *The mean of this distribution: $E(X) = p$*
- *The variance of this distribution: $\text{Var}(X) = p(1 - p)$.*

Proof. We only need to verify the variance. Since the second moment is $E(X^2) = p$, the statement ensues. \square

18.2 Binomial distribution

Let (Ω, \mathcal{A}, P) be a probability space, and consider the Bernoulli experiment, where we carry out n independent experiments in a row, and every time we observe if a given event A occurs. Suppose that $P(A) = p$, $0 < p < 1$ is given. Let X denote how many times A comes out. By the Bernoulli experiment the distribution of X is given by

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k} \quad k = 0, 1, 2, \dots, n$$

This distribution is called the *binomial distribution*.

Theorem 18.2

- The parameters of the distribution: $n \in \mathbb{N}$ and $0 < p < 1$.
- The mean of the distribution: $E(X) = np$
- The variance of the distribution: $\text{Var}(X) = np(1-p)$.

Proof. In view of Example 17.4 we only need to check the variance. First find the second moment.

$$\begin{aligned} E(X^2) &= \sum_{k=1}^n k^2 \binom{n}{k} p^k (1-p)^{n-k} = \\ &= \sum_{k=2}^n k(k-1) \binom{n}{k} p^k (1-p)^{n-k} + \sum_{k=1}^n k \binom{n}{k} p^k (1-p)^{n-k} \\ &= n(n-1)p^2 \sum_{k=2}^n \binom{n-2}{k-2} p^{k-2} (1-p)^{n-k} + np = (n^2 - n)p^2 + np. \end{aligned}$$

Therefore, the variance is

$$\text{Var}(X) = E(X^2) - E(X)^2 = n(n-1)p^2 + np - n^2p^2 = np(1-p)$$

where we observed that the second sum in the second line is precisely the mean. \square

18.3 Hypergeometric distribution

Examine the following sampling without replacement problem. Consider a set of N objects in which m of them are defective. Select a sample of n objects

without replacement from the whole set ($n \leq m$). Let X denote the number of defective objects in the sample. Then the distribution of X is:

$$P(X = k) = \frac{\binom{m}{k} \cdot \binom{N-m}{n-k}}{\binom{N}{n}} \quad k = 0, 1, 2, \dots, n$$

This distribution is called the *hypergeometric distribution*.

Theorem 18.3

- The parameters of the distribution: $N, m, n \in \mathbb{N}$.
- The mean of the distribution:

$$E(X) = n \cdot \frac{m}{N}$$

- The variance of the distribution:

$$\text{Var}(X) = \frac{N-n}{N-1} \cdot n \cdot \frac{m}{N} \left(1 - \frac{m}{N}\right).$$

Proof. By applying the argument of Example 17.3, we again only have to calculate the variance. First find the second moment.

$$\begin{aligned} E(X^2) &= \sum_{k=1}^n k^2 \frac{\binom{m}{k} \binom{N-m}{n-k}}{\binom{N}{n}} = \sum_{k=2}^n k(k-1) \frac{\binom{m}{k} \binom{N-m}{n-k}}{\binom{N}{n}} + \sum_{k=1}^n k \frac{\binom{m}{k} \binom{N-m}{n-k}}{\binom{N}{n}} \\ &= \frac{m(m-1)n(n-1)}{N(N-1)} \sum_{k=2}^n \frac{\binom{m-2}{k-2} \cdot \binom{N-m}{n-k+2}}{\binom{N-2}{n-2}} + n \frac{m}{N} \\ &= \frac{m(m-1)n(n-1)}{N(N-1)} + n \frac{m}{N}. \end{aligned}$$

Then we conclude

$$\text{Var}(X) = \frac{m(m-1)n(n-1)}{N(N-1)} + n \frac{m}{N} - n^2 \frac{m^2}{N^2} = \frac{N-n}{N-1} \cdot n \cdot \frac{m}{N} \left(1 - \frac{m}{N}\right),$$

just as we stated. \square

18.4 Geometric distribution

Take a probability space (Ω, \mathcal{A}, P) , and consider an event A such that $P(A) = p$, wher $0 < p < 1$ is given. Keep performing the experiment until the event A occurs for the first time. Let X denote the number of trials. The distribution of X is given by:

$$P(X = k) = (1-p)^{k-1} p \quad k = 1, 2, \dots$$

This distribution is called the *geometric distribution*.

Theorem 18.4

- *The parameter of the distribution:* $0 < p < 1$.
- *The mean of the distribution:*

$$E(X) = \frac{1}{p}$$

- *The variance of the distribution:*

$$\text{Var}(X) = \frac{1-p}{p^2}.$$

Proof. The mean of this distribution is easily obtained by following the argument of Example 17.5, so we only need to find the variance. The second moment can be evaluated the following way. Using the second derivative of the power series at $|x| < 1$, we have

$$\sum_{k=2}^{\infty} k(k-1)x^{k-2} = \frac{2}{(1-x)^3}$$

If we employ this identity with $x = 1-p$ we receive

$$\begin{aligned} E(X^2) &= \sum_{k=1}^{\infty} k^2(1-p)^{k-1}p = \sum_{k=2}^{\infty} k(k-1)(1-p)^{k-1}p + \sum_{k=1}^{\infty} k(1-p)^{k-1}p \\ &= p(1-p) \sum_{k=1}^{\infty} k(k-1)(1-p)^{k-2} + \frac{1}{p} = \frac{2p(1-p)}{p^3} + \frac{1}{p}. \end{aligned}$$

Thus we get

$$\text{Var}(X) = E(X^2) - E(X)^2 = \frac{2p(1-p)}{p^3} + \frac{1}{p} - \frac{1}{p^2} = \frac{1-p}{p^2}$$

and this is what we needed. \square

18.5 Poisson distribution

Suppose that X is a random variable, whose range is $\{0\} \cup \mathbb{N}$ and its distribution is defined by

$$P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda} \quad k = 0, 1, 2, \dots$$

wher $\lambda > 0$ is a given number.

It is not hard to see that we really defined a distribution. Indeed,

$$\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} = e^{-\lambda} \cdot e^{\lambda} = 1$$

based on the power series of the natural exponential function. This infinite distribution is called the *Poisson distribution*.

Theorem 18.5

- *The parameter of the distribution: $\lambda > 0$.*
- *The mean of the distribution: $E(X) = \lambda$,*
- *The variance of the distribution: $\text{Var}(X) = \lambda$.*

Proof. In view of Example 17.6 we only have to calculate the variance. The second moment is obtained as follows.

$$\begin{aligned} E(X^2) &= \sum_{k=1}^{\infty} k^2 \frac{\lambda^k}{k!} e^{-\lambda} = \sum_{k=1}^{\infty} k(k-1) \frac{\lambda^k}{k!} e^{-\lambda} + \sum_{k=1}^{\infty} k \frac{\lambda^k}{k!} e^{-\lambda} \\ &= \lambda^2 \sum_{k=2}^{\infty} \frac{\lambda^{k-2}}{(k-2)!} e^{-\lambda} + \lambda. \end{aligned}$$

Hence, the variance is

$$\text{Var}(X) = E(X^2) - E(X)^2 = \lambda^2 + \lambda - \lambda^2 = \lambda$$

and that completes the proof. \square

Let us remark that the Poisson distribution can be regarded as the "limit distribution" of the binomial distribution as it is explained in the following.

Theorem 18.6 *If $\lambda > 0$ is fixed and $0 < p_n < 1$ is a sequence with $np_n = \lambda$, then*

$$\lim_{n \rightarrow \infty} \binom{n}{k} p_n^k (1-p_n)^{n-k} = \frac{\lambda^k}{k!} e^{-\lambda}$$

for every $k = 0, 1, 2, \dots$

Proof. Indeed, for each fixed index k we have

$$\begin{aligned} \binom{n}{k} p_n^k (1-p_n)^{n-k} &= \binom{n}{k} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ &= \frac{n(n-1)\dots(n-k+1)}{n^k} \cdot \frac{\lambda^k}{k!} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-k} \end{aligned}$$

Here examine the limits of the four factors separately. It is easy to see that they are 1, $\lambda^k/k!$, $e^{-\lambda}$ and 1 respectively. That proves our theorem. \square

Practically, this theorem means that for large values of n and for small values of p the binomial distribution can be approximated by the Poisson distribution, i.e.

$$\binom{n}{k} p_n^k (1 - p_n)^{n-k} \approx \frac{\lambda^k}{k!} e^{-\lambda}$$

for every index $0 \leq k \leq n$.

Example 18.7 Let us suppose that in a brand new Suzuki Vitara the probability that the airbag is defective, is 0.002 independently from each other. The factory announces withdrawal if at least 10 malfunctions are reported for the 2000 cars that are manufactured in a month. Find the probability that no withdrawal has to be announced.

let X denote the number of defective cars in a given month. Since this is a Bernoulli-experiment (the probability of malfunction is 0.002 independently from each other), it follows that X has binomial distribution with parameters $n = 2000$ and $p = 0.002$. Therefore the exact value of the probability is

$$P(X \leq 9) = \sum_{k=0}^9 \binom{2000}{k} 0.002^k 0.998^{2000-k}$$

which not easy to handle. Based on our theorem, we can give an approximation of this probability by using the Poisson distribution (we say that " n is sufficiently large and p is sufficiently small"), moreover $\lambda = np = 4$, so

$$\sum_{k=0}^9 \binom{2000}{k} 0.002^k 0.998^{2000-k} \approx \sum_{k=0}^9 \frac{4^k}{k!} e^{-4} \approx 0.9919$$

This latter value can be determined by looking up in the Poisson tables that can be found on page 732 in our Textbook.

Recitation and Exercises

1. Reading: Textbook-2, Sections 5.1, 5.2, 5.3 and 5.5
2. Homework: Textbook-2, Exercises 5.5, 5.9, 5.10, 5.15, 5.27, 5.33, 5.47, 5.56, 5.60, 5.66, 5.70 and 5.72
3. Review: Calculus, integration, improper integrals and infinite series, and "Probability Exercises"

Chapter 19

Special continuous distributions

19.1 Uniform distribution

Let $[a, b]$ be a given finite interval. Consider a random variable X with the following density function:

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{elsewhere} \end{cases}$$

This random variable X is said to have *uniform distribution* on the interval $[a, b]$. The name comes from the fact that the probability that X is in a subinterval of $[a, b]$ is proportional to the length of the subinterval.

Theorem 19.1

- *The parameters of the distribution: a and b , $a < b$.*
- *The mean of the distribution:*

$$E(X) = \frac{a+b}{2}$$

- *The variance of the distribution:*

$$Var(X) = \frac{(b-a)^2}{12}.$$

Proof. These statements are immediate consequences of the results in Examples 17.10 and 17.15. \square

Example 19.2 Let X be a uniformly distributed random variable with $E(X) = 5$ and $Var(X) = 3$. Find the probability $P(4 < X < 10)$.

The unknown endpoints of the interval a and b satisfy the following equations:

$$\begin{aligned}\frac{a+b}{2} &= 5 \\ \frac{(b-a)^2}{12} &= 3\end{aligned}$$

whose solutions are $a = 2$ and $b = 8$. Therefore,

$$P(4 < X < 10) = P(4 < X < 8) = \frac{2}{3}$$

since the subinterval beyond $[4, 8]$ comes with 0 probability.

19.2 Exponential distribution

Let $\lambda > 0$ be a fixed number. Consider the random variable X with the following density function

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x > 0 \\ 0 & \text{elsewhere} \end{cases}$$

In this case we say that X has *exponential distribution*. nevezük.

ATTENTION: Verify that f really defines a density function! Sketch the graph of the function!

Theorem 19.3

- *The parameter of the distribution:* $\lambda > 0$.
- *The mean of the distribution:* $E(X) = 1/\lambda$,
- *The variance of the distribution:* $\text{Var}(X) = 1/\lambda^2$.

Proof. Our theorem is an immediate consequence of the equalities in Examples 17.12 and 17.14. \square

Example 19.4 Consider an exponentially distributed random variable X with a given parameter $\lambda > 0$. Find the probability $P(X > E(X))$.

Our theorem claims that $E(X) = 1/\lambda$, thus

$$P(X > E(X)) = P\left(X > \frac{1}{\lambda}\right) = \int_{1/\lambda}^{\infty} \lambda e^{-\lambda x} dx = [-e^{-\lambda x}]_{1/\lambda}^{\infty} = \frac{1}{e}.$$

We say that the exponential distribution is memoryless in the following sense. If X is exponentially distributed with a given parameter $\lambda > 0$, and $t, s > 0$ are given positive numbers, then

$$P(X > t + s | X > t) = P(X > s).$$

Indeed, the event $\{X > t + s\}$ implies the event $\{X > t\}$, therefore, the conditional probability on the left-hand side can be written like

$$\begin{aligned} P(X > t + s | X > t) &= \frac{P(X > t + s)}{P(X > t)} = \frac{1 - \int_0^{t+s} \lambda e^{-\lambda x} dx}{1 - \int_0^t \lambda e^{-\lambda x} dx} \\ &= \frac{e^{-\lambda(t+s)}}{e^{-\lambda t}} = e^{-\lambda s} = 1 - \int_0^s \lambda e^{-\lambda x} dx = P(X > s). \end{aligned}$$

If for instance X denotes the waiting time between two occurrences (i.e. two telephone calls, or two customers, etc.), then the lack of memory means that the further waiting time does not depend on how much we have been waiting.

Conversely, it can also be proven that if a continuous distribution is memoryless, then it is necessarily the exponential distribution.

There is an interesting relationship between the Poisson distribution and the exponential distribution. In particular, if the waiting times between successive occurrences are independent, exponentially distributed random variables with identical parameter $\lambda > 0$, then the number of occurrences in a unit time interval has Poisson distribution with the same parameter. These features will be discussed in later chapters.

19.3 The standard normal distribution

Because of the central role of the standard normal distribution we use a distinguished notation for its density function and cumulative distribution function.

Definition 19.5 We say that the random variable Z has *standard normal distribution*, if its density function is given by

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad -\infty < x < \infty$$

In view of formula (9.2), we see that φ really defines a density function. As an exercise analyze the function φ , and show that it possesses the following properties.

$$\lim_{x \rightarrow -\infty} \varphi(x) = \lim_{x \rightarrow +\infty} \varphi(x) = 0$$

moreover φ is strictly monotone increasing on the interval $(-\infty, 0)$, strictly monotone decreasing on the interval $(0, \infty)$, and reaches its global maximum at $x = 0$.

By analyzing the second derivative, we can see that φ is convex on the intervals $(-\infty, 1)$ and $(1, +\infty)$, while it is concave on the interval $(-1, 1)$, and consequently has points of inflection at $x = -1$ and $x = 1$ respectively.

EXERCISE: CREATE THE GRAPH OF THE FUNCTION!

Theorem 19.6

- *The parameter of the distribution: no parameter.*
- *The mean of the distribution: $E(Z) = 0$.*
- *The variance of the distribution: $Var(Z) = 1$.*

Proof. Example 9.6 shows that $E(Z) = 0$, and equality (9.3) tells us that $E(Z^2) = 1$. Therefore

$$Var(Z) = E(Z^2) - E(Z)^2 = 1.$$

as we stated. □

Let Φ denote the standard normal cumulative distribution function, i.e.

$$\Phi(x) = \int_{-\infty}^x \varphi(t) dt.$$

This function has the properties of cumulative distribution functions, but its interesting feature is that it cannot be expressed explicitly in terms of elementary functions or their finite combinations.

Observe however that φ is an even function, in other words it is symmetric with respect to the y -axis. This implies that $\Phi(0) = 1/2$, and further

$$\Phi(-x) = 1 - \Phi(x) \tag{19.1}$$

for every real number x .

Example 19.7 Because of its central role in Statistics and other applications we can find tables for the values of the Φ function in most probability textbooks and spreadsheet programs like the Microsoft Windows Office Excel application. See the tables on pages 735–736 of our Textbook!

If for example Z is a standard normally distributed random variable, the find the probability

$$P(-2 < Z < 2)$$

Using the table on page 736 of our Textbook, we get

$$\begin{aligned} P(-2 < Z < 2) &= \Phi(2) - \Phi(-2) = \Phi(2) - (1 - \Phi(2)) = 2\Phi(2) - 1 = \\ &= 2 \cdot 0.9772 - 1 = 0.9544 \end{aligned}$$

where we exploited the symmetry property (19.1).

19.4 Normal distribution

Definition 19.8 Let m and σ be given real numbers where $\sigma > 0$. Let Z be a standard normally distributed random variable, then the random variable

$$X = \sigma Z + m$$

is said to have *normal distribution with (m, σ) -parameters* (or briefly (m, σ) -normal distribution).

Making use of the properties of the standard normal distribution, and the properties of the mean and the variance (refer to Theorem 17.11) we get the following theorem for (m, σ) -normal distributions.

Theorem 19.9

- *The parameters of the distribution: $m \in \mathbb{R}, \sigma > 0$.*
- *The mean of the distribution: $E(X) = m$,*
- *The variance of the distribution: $\text{Var}(X) = \sigma^2$.*

How can we find the cumulative distribution function and the density function of this random variable X ? Let F denote the cumulative distribution function of X , and take a real number x arbitrarily. Then

$$F(x) = P(X < x) = P(\sigma Z + m < x) = P\left(Z < \frac{x - m}{\sigma}\right) = \Phi\left(\frac{x - m}{\sigma}\right)$$

IMPORTANT! It is vital that $\sigma > 0$, so when we divide by σ the inequality will not change!

We get the density function of X by differentiating F : for every $x \in \mathbb{R}$ we have

$$f(x) = F'(x) = \frac{1}{\sigma} \varphi\left(\frac{x - m}{\sigma}\right) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-m)^2}{2\sigma^2}}$$

by the Chain-Rule. This function has a global maximum at $x = m$, furthermore it has points of inflection at $x = m - \sigma$ and $x = m + \sigma$ respectively. CREATE A PICTURE!

Example 19.10 For an (m, σ) -normally distributed random variable X the probability of being in an interval can always be expressed in terms of the standard normal cumulative distribution function Φ .

Indeed, if $a < b$ are arbitrarily taken real numbers, then

$$P(a < X < b) = F(b) - F(a) = \Phi\left(\frac{b - m}{\sigma}\right) - \Phi\left(\frac{a - m}{\sigma}\right).$$

For example, for a normally distributed random variable X with parameters $m = 10$ and $\sigma = 2$ we have

$$\begin{aligned} P(7 < X < 13) &= F(13) - F(7) = \Phi(1.5) - \Phi(-1.5) = 2\Phi(1.5) - 1 = \\ &= 2 \cdot 0.9332 - 1 = 0.8664 \end{aligned}$$

where we used the symmetry of Φ , and the tables on page 736 in the Textbook.

Recitation and Exercises

1. Reading: Textbook, Sections 6.1, 6.2, 6.3, 6.4, 6.6
2. Homework: Textbook, Exercises 6.2, 6.3, 6.4, 6.6, 6.7, 6.9, 6.11, 6.15, 6.17, 6.18, 6.45 and 6.46 5.66, 5.70 and 5.72
3. Review: Calculus, integration, improper integrals and infinite series and "Probability Exercises"

Chapter 20

Joint distributions

20.1 Joint cumulative distribution function

Definition 20.1 Let X and Y be random variables (not necessarily on the same sample space). For any real numbers x and y the function

$$F(x, y) = P(X < x, Y < y)$$

is called the *joint cumulative distribution function* of X and Y .

The following statement comes directly from the definition.

Proposition 20.2 *If F is a joint cumulative distribution function, then*

$$\lim_{x \rightarrow -\infty} F(x, y) = \lim_{y \rightarrow -\infty} F(x, y) = 0$$

for any fixed real y and x respectively, moreover

$$\lim_{x, y \rightarrow +\infty} F(x, y) = 1$$

Similarly to the one dimensional case, we separately discuss discrete and continuous distributions.

20.2 Discrete joint distributions

Definition 20.3 Assume that the range of the variable X is $\{x_1, x_2, \dots\}$, and the range of the variable Y is $\{y_1, y_2, \dots\}$. Then the joint distribution of X and Y is given by

$$p_{ik} = P(X = x_i, Y = y_k) \quad i = 1, 2, \dots \quad k = 1, 2, \dots$$

These values can be arranged in a chart:

$y \setminus x$	x_1	x_2	x_3	\cdots
y_1	p_{11}	p_{21}	p_{31}	\cdots
y_2	p_{12}	p_{22}	p_{32}	\cdots
y_3	p_{13}	p_{23}	p_{33}	\cdots
\vdots	\vdots	\vdots	\vdots	\vdots

Obviously for all indices $p_{ik} \geq 0$ and $\sum_i \sum_k p_{ik} = 1$.

Let A be a subset of the plane. By using the joint distribution, how can we evaluate the probability $P((X, Y) \in A)$? Collect all values x_i and y_k for which $(x_i, y_k) \in A$, then

$$P((X, Y) \in A) = \sum_{(x_i, y_k) \in A} p_{ik}$$

Example 20.4 For instance, if we consider the following joint distribution

$y \setminus x$	0	1	2	3
0	0.1	0.08	0.13	0.04
1	0.04	0.2	0.08	0
2	0.03	0	0.05	0.25

then for the subset $A = \{(x, y) \in \mathbb{R}^2 : x + y \geq 3\}$ we have:

$$P(X + Y \geq 3) = 0.04 + 0.08 + 0.05 + 0.25 = 0.42$$

A natural question to ask is that based on the joint distribution, how can we determine the distributions of X and Y alone? As we conclude from the definition

$$p_i = P(X = x_i) = \sum_k p_{ik} = \sum_k P(X = x_i, Y = y_k) \quad i = 1, 2, \dots$$

Namely, the probability $p_i = P(X = x_i)$ can be obtained by taking the sum of the elements in the i -th column. Therefore, the sums of columns provide the distribution of X .

In an analogous way,

$$q_k = P(Y = y_k) = \sum_i p_{ik} = \sum_i P(X = x_i, Y = y_k) \quad k = 1, 2, \dots$$

which means that the distribution of Y is obtained by taking the sums of rows.

Definition 20.5 The distributions of X and Y are called the *marginal distributions* of the joint distribution.

20.3 Continuous joint distributions

Definition 20.6 We say that X and Y are continuously distributed, if there exists a non-negative integrable function f on the plane such that for all real numbers x and y we have

$$F(x, y) = \int_{-\infty}^x \int_{-\infty}^y f(t, s) ds dt$$

where F is the joint cumulative distribution function of the random variables X and Y . This function f is called the *joint density function* of X and Y .

Clearly, if f is a joint density function, then

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1$$

Example 20.7 Let A be a subset of the plane. How can we find the probability $P((X, Y) \in A)$? If f is the joint density function of X and Y , then

$$P((X, Y) \in A) = \iint_A f(x, y) dy dx$$

For example if we consider the joint density function

$$f(x, y) = \begin{cases} \frac{2}{3}(x + 2y) & \text{if } 0 < x < 1, 0 < y < 1 \\ 0 & \text{elsewhere} \end{cases} \quad (20.1)$$

then for the set $A = \{(x, y) \in \mathbb{R}^2 : x < 1/2, y < 1/2\}$ we have

$$P(X < 1/2, Y < 1/2) = \frac{2}{3} \int_0^{1/2} \int_0^{1/2} (x + 2y) dy dx = \frac{1}{8}$$

If the joint density function is given, how can we find the density of X or Y alone? It can be shown that if f_X denotes the density of X , then for every point x

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy$$

and analogously

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx$$

for every point y .

Definition 20.8 The functions f_X and f_Y are called the *marginal densities* of the joint density function.

Example 20.9 For instance in the case of the joint density in the previous example

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy = \begin{cases} \int_0^1 \frac{2}{3}(x+2y) dy & \text{if } 0 < x < 1 \\ 0 & \text{elsewhere} \end{cases}$$

The marginal density of Y in a similar way

$$f_Y(y) = \begin{cases} \frac{1}{3}(4y+1) & \text{if } 0 < y < 1 \\ 0 & \text{elsewhere} \end{cases}$$

20.4 Independence

Definition 20.10 Let X and Y be random variables with joint cumulative distribution function F . Denote by F_X and F_Y the marginal cumulative distribution functions of X and Y respectively. We say that X and Y are *independent*, if

$$F(x, y) = F_X(x) \cdot F_Y(y)$$

for all real numbers x, y .

In other words we may say that X and Y are independent, if

$$P(X < x, Y < y) = P(X < x) \cdot P(Y < y)$$

for all real numbers x, y . Now we reformulate this definition for the discrete and for the continuous case.

Let X and Y be discrete random variables with joint distribution

$$P(X = x_i, Y = y_k) = p_{ik} \quad i = 1, 2, \dots \quad k = 1, 2, \dots$$

Consider the marginal distributions of X and Y :

$$P(X = x_i) = p_i \quad i = 1, 2, \dots \quad P(Y = y_k) = q_k \quad k = 1, 2, \dots$$

Theorem 20.11 X and Y are independent if and only if

$$p_{ik} = p_i \cdot q_k$$

for all indices i and k .

Our theorem states that the random variables are independent if and only if their joint distribution can be expressed as the product of the marginal distributions. For instance in Example 20.4 the variables are not independent, since for the very first element

$$0.17 \cdot 0.35 = p_1 \cdot q_1 \neq p_{11} = 0.1,$$

VERIFY!

Now let X and Y be continuously distributed random variables with joint density function f . Denote by f_X and f_Y the marginal densities of X and Y respectively.

Theorem 20.12 *X and Y are independent if and only if*

$$f(x, y) = f_X(x) \cdot f_Y(y)$$

for every real x and y .

Proof. Easily follows from the equality $F(x, y) = F_X(x) \cdot F_Y(y)$. \square

Example 20.13 In Example 20.1 the random variables are not independent, since

$$f_X(x) \cdot f_Y(y) \neq f(x, y),$$

i.e. the joint density cannot be expressed as the product of the marginal densities.

However, if the joint density of X and Y is given by

$$f(x, y) = \begin{cases} 4xy & \text{if } 0 < x < 1, 0 < y < 1 \\ 0 & \text{elsewhere} \end{cases}$$

then X and Y are independent. Indeed

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy = \int_0^1 4xy dy = \begin{cases} 2x & \text{if } 0 < x < 1 \\ 0 & \text{elsewhere} \end{cases}.$$

and by the symmetry of f the marginal density f_Y has the same form with respect to y . Thus

$$f(x, y) = f_X(x) \cdot f_Y(y)$$

for all real numbers x and y .

20.5 Conditional distributions

Consider the discrete random variables X és Y with joint distribution $P(X = x_i, Y = y_k) = p_{ik}$, where $i = 1, 2, \dots$ and $k = 1, 2, \dots$

Definition 20.14 Suppose that for a specific index k we have $P(Y = y_k) > 0$. Then by the *conditional distribution* of X under the condition $Y = y_k$ we mean the distribution

$$P(X = x_i | Y = y_k) = \frac{P(X = x_i, Y = y_k)}{P(Y = y_k)}, \quad i = 1, 2, \dots$$

ATTENTION! Verify directly that we really have defined a distribution!

Definition 20.15 By the *conditional expected value* of X under the condition $Y = y_k$ we mean the sum

$$E(X | Y = y_k) = \sum_{i=1} x_i \cdot P(X = x_i | Y = y_k)$$

that may consist of finitely many or infinitely many terms depending on the range of X (this is why we do not indicate the upper bound of the summation).

Example 20.16 Let us examine again the joint distribution in Example 20.4. Then $P(Y = 1) = 0.32$, and the conditional expected value of X under the condition $Y = 1$

$$E(X | Y = 1) = 0 \cdot 0.04 + 1 \cdot 0.2 + 2 \cdot 0.08 + 3 \cdot 0 = 0.36$$

Verify this calculation!

Recitation and Exercises

1. Reading: Textbook-2, Section 3.4
2. Homework: Textbook-2, Exercises 3.39, 3.40, 3.41, 3.42, 3.43, 3.45, 3.47, 3.49, 3.50, 3.51, 3.52 and 3.53
3. Review: Calculus, integration, improper integrals and infinite series, and "Probability Exercises"

Chapter 21

Covariance and correlation

21.1 Mean of a sum

Tétel 21.1 *If the random variables X and Y both have a mean, then so does $X + Y$ and*

$$E(X + Y) = E(X) + E(Y)$$

Proof. We give an outline of the proof in the discrete case, the continuous case is analogous.

$$\begin{aligned} E(X + Y) &= \sum_i \sum_k (x_i + y_k) P(X = x_i, Y = y_k) \\ &= \sum_i x_i \sum_k p_{ik} + \sum_k y_k \sum_i p_{ik} \\ &= \sum_i x_i P(X = x_i) + \sum_k y_k P(Y = y_k) = E(X) + E(Y) \quad \square \end{aligned}$$

This theorem remains true for a sum with a finite number of terms (use induction!).

Example 21.2 Suppose that on n pieces of cards we wrote the integers $1, \dots, n$, and then placed them in a hat. We choose m pieces of cards from the hat at random, with replacement. Let X denote the sum of the integers. Find $E(X)$.

The distribution of X in that problem is hard to find. Give it a try!

Denote by X_1, \dots, X_m the numbers selected. In view of the selection with replacement, each X_k is identically distributed, namely:

$$P(X_k = i) = \frac{1}{n} \quad i = 1, \dots, n$$

This means that for every k

$$E(X_k) = \sum_{i=1}^n i \cdot \frac{1}{n} = \frac{n(n+1)}{2} \cdot \frac{1}{n} = \frac{n+1}{2}.$$

On the other hand, clearly $X = X_1 + \dots + X_m$, and therefore

$$E(X) = E(X_1) + \dots + E(X_m) = m \cdot \frac{n+1}{2}$$

Thus $E(X)$ can be found without even knowing the distribution of X !

21.2 Mean of a product

If the discrete random variables X and Y , then

$$E(XY) = \sum_i \sum_k x_i y_k \cdot p_{ik}$$

where the range of X is $\{x_1, x_2, \dots\}$, and the range of Y is $\{y_1, y_2, \dots\}$ respectively, and p_{ik} denotes their joint distribution.

In a completely similar way, if X and Y continuously distributed, both have a mean, and their joint density function is f , then

$$E(XY) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy \cdot f(x, y) dx dy$$

Theorem 21.3 *If X and Y are independent, then*

$$E(XY) = E(X) \cdot E(Y)$$

Proof. We just focus on the continuous case, the discrete case can be treated similarly.

$$\begin{aligned} E(XY) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy \cdot f(x, y) dx dy = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy \cdot f_X(x) \cdot f_Y(y) dx dy \\ &= \int_{-\infty}^{\infty} x f_X(x) dx \cdot \int_{-\infty}^{\infty} y f_Y(y) dy = E(X) \cdot E(Y) \end{aligned}$$

since the independence implies that the joint density is the product of the marginal densities, i.e. $f(x, y) = f_X(x) \cdot f_Y(y)$. \square

21.3 Variance of a sum

Theorem 21.4 *Assume that X and Y are independent, and they both have a variance. Then*

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$$

The statement can be extended to any finite number of terms.

Proof. Exploit our theorem about the mean of the product, then we get:

$$\begin{aligned} \text{Var}(X + Y) &= E((X + Y - E(X + Y))^2) \\ &= E((X - E(X))^2) + E((Y - E(Y))^2) \\ &\quad + 2E((X - E(X))(Y - E(Y))) \\ &= \text{Var}(X) + \text{Var}(Y) + 2(E(XY) - E(X)E(Y)) \\ &= \text{Var}(X) + \text{Var}(Y). \quad \square \end{aligned}$$

Example 21.5 Why do we think that by repeatedly performing an experiment and taking the average of the results we can expect a more accurate result?

Let us suppose that for determining an unknown quantity m we perform n observations, and the results are the random variables X_1, \dots, X_n . We assume that the variables are independent and identically distributed with

$$E(X_k) = m, \quad D(X_k) = \sigma, \quad k = 1, 2, \dots, n.$$

The assumption that all variables have the same distribution means that the observations (measurements) are carried out in the same circumstances. Then σ is interpreted as the expected error. Take the arithmetic average of our results, i.e. introduce the random variable

$$Y_n = \frac{X_1 + \dots + X_n}{n}$$

Then clearly $E(Y_n) = m$, moreover, according to our theorem above

$$\text{Var}(Y_n) = \text{Var}\left(\frac{1}{n}(X_1 + \dots + X_n)\right) = \frac{1}{n^2}n \cdot \sigma^2 = \frac{\sigma^2}{n}.$$

as a consequence of independence. Thus, for the standard deviation of Y_n we obtain:

$$D(Y_n) = \frac{\sigma}{\sqrt{n}}$$

for which $D(Y_n) \rightarrow 0$ as $n \rightarrow \infty$. Hence, the expected error tends to zero, when n approaches infinity.

21.4 Covariance and correlation

The following concepts are used for measuring the degree of dependence of random variables.

Definition 21.6 The *covariance* of random variables X and Y is defined by

$$\text{Cov}(X, Y) = E((X - E(X)) \cdot (Y - E(Y)))$$

and their *correlation coefficient* is given by

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{D(X) \cdot D(Y)}$$

As it is easy to see

$$\text{Cov}(X, Y) = E(XY - E(X)Y - E(Y)X + E(X)E(Y)) = E(XY) - E(X)E(Y),$$

and most of the time, this simpler expression is used to evaluate the covariance.

The covariance is NOT an absolute measurement of the independence, since for any $\alpha \neq 0$ we have

$$\text{Cov}(\alpha X, Y) = \alpha \text{Cov}(X, Y)$$

so it depends on the dimensions. Just think of the case when X and Y are costs given in Euro, but if we convert them to Forint, then their covariance will change to approximately 340^2 times higher. However, the correlation coefficient is independent of the dimension, since for any real numbers $\alpha \neq 0$ and β we have:

$$\text{Corr}(\alpha X + \beta, \alpha Y + \beta) = \text{Corr}(X, Y)$$

which means that the correlation is independent of linear transformations. ATTENTION! Verify this equality directly by the definition!

Theorem 21.7

1. $-1 \leq \text{Corr}(X, Y) \leq 1$
2. If X and Y are independent, then $\text{Cov}(X, Y) = 0$

Proof. For proving the first statement, take a real number $t \in \mathbb{R}$ arbitrarily, and consider the random variable

$$W = [X - E(X) + t(Y - E(Y))]^2$$

Since W is nonnegative, so is its mean. This means that

$$E(W) = E((X - E(X))^2) + 2tCov(X, Y) + t^2E((Y - E(Y))^2) \geq 0$$

for every real number t . This expression is quadratic with respect to t , and therefore it can only be nonnegative, if its discriminant is nonpositive, that is:

$$4Cov(X, Y)^2 - 4E((X - E(X))^2)E((Y - E(Y))^2) \leq 0.$$

Rearranging the terms, and taking the square root of both sides, we get:

$$|Cov(X, Y)| \leq D(X)D(Y)$$

The second statement is an immediate consequence of Theorem 21.3. \square

Example 21.8 ATTENTION! The example below shows that the converse of the second statement of our theorem is not true! Toss a coin twice in a row, and introduce the random variables:

$$X_k = \begin{cases} 0 & \text{if toss } k \text{ is a Head} \\ 1 & \text{if toss } k \text{ is a Tail} \end{cases}$$

($k = 1, 2$). Consider the variables $Y_1 = X_1 + X_2$ and $Y_2 = X_1 - X_2$. Then their joint distribution is:

$Y_2 \setminus Y_1$	0	1	2
-1	0	0.25	0
0	0.25	0	0.25
1	0	0.25	0

By examining the joint distribution, we see that Y_1 and Y_2 are not independent, but we can easily calculate that $Cov(Y_1, Y_2) = 0$

21.5 Theorem of Total Expectation

Consider the discrete random variables X and Y that have a joint distribution $P(X = x_i, Y = y_k) = p_{ik}$, and $P(Y = y_k) > 0$ for all indices $i = 1, 2, \dots$ and $k = 1, 2, \dots$

Definition 21.9 Create the conditional expected values of X under the conditions $Y = y_k$ that is:

$$m_k = E(X|Y = y_k) = \sum_{i=1} x_i P(X = x_i|Y = y_k)$$

for every $k = 1, 2, \dots$. This sequence is called the *conditional expectation* of X with respect to the variable Y . Its notation is $E(X|Y)$.

Observe that this way we have defined a random variable, namely

$$E(X|Y) = m_k, \quad \text{ha } Y = y_k, k = 1, 2, \dots$$

Below we determine the mean of this random variable. This result can be regarded as the generalization of Theorem of Total Probability.

Tétel 21.10 (Theorem of Total Expectation) $E(E(X|Y)) = E(X)$.

Proof. Indeed,

$$\begin{aligned} E(E(X|Y)) &= \sum_{k=1} m_k P(Y = y_k) = \sum_{k=1} \sum_{i=1} x_i P(X = x_i | Y = y_k) P(Y = y_k) \\ &= \sum_{i=1} x_i \sum_{k=1} P(X = x_i, Y = y_k) = \sum_{i=1} x_i P(X = x_i) = E(X) \end{aligned}$$

since, in the second line, we obtain precisely the marginal distribution of X . \square

ATTENTION! Why can we interchange the sums in the second line?

Example 21.11 In some situations it is easier to find $E(X)$ by our theorem than by the direct approach. The number of calls received by a call center on a given day has Poisson distribution with a parameter $\lambda > 0$. Every call is a wrong number with a given probability $p > 0$, independently from each other. Find the expected value of the wrong number calls on that day.

Let X denote the number of wrong calls, and Y the total number of calls. It is clear that for any fixed $n \in \mathbb{N}$ under the condition $Y = n$ we face the Bernoulli-experiment. Therefore,

$$P(X = k | Y = n) = \binom{n}{k} p^k (1-p)^{n-k} \quad \text{if } n \geq k$$

while $P(X = k | Y = n) = 0$, if $n < k$. Hence, the conditional expected value is given by

$$m_n = E(X | Y = n) = np, \quad n = 1, 2, \dots$$

Making use of the Theorem of Total Expectation, we obtain

$$E(X) = E(E(X|Y)) = \sum_{n=1}^{\infty} np \frac{\lambda^n}{n!} e^{-\lambda} = \lambda p$$

ATTENTION! Find $E(X)$ directly by using the distribution of X as well!

Recitation and Exercises

1. Reading: Textbook-2, Sections 4.1, 4.2 and 4.3.
2. Homework: Textbook-2, Exercises 4.23, 4.24, 4.52, 4.59, 4.60, 4.64, 4.70, 4.98.
3. Review: "Probability Exercises"

Chapter 22

Sums of random variables

22.1 Sums of discrete variables

Assume that X and Y are independent variables, and both have Poisson-distribution, with parameters $\lambda > 0$ and $\mu > 0$ respectively. Find the distribution of $X + Y$. Then for any fixed integer k

$$\begin{aligned}P(X + Y = k) &= \sum_{i=0}^k P(X = i, Y = k - i) = \sum_{i=0}^k P(X = i) \cdot P(Y = k - i) \\&= \sum_{i=0}^k \frac{\lambda^i}{i!} e^{-\lambda} \cdot \frac{\mu^{k-i}}{(k-i)!} e^{-\mu} = \frac{e^{-(\lambda+\mu)}}{k!} \sum_{i=0}^k \binom{k}{i} \lambda^i \mu^{k-i} \\&= \frac{(\lambda + \mu)^k}{k!} e^{-(\lambda+\mu)}\end{aligned}$$

by the independence. Thus, $X + Y$ has Poisson-distribution with the parameter $\lambda + \mu$.

Using induction, this result can be extended to any finite number of terms.

Tételet 22.1 *Assume that X_1, \dots, X_n are independent variables, and have Poisson-distribution with parameters $\lambda_1, \dots, \lambda_n$ respectively. Then the random variable*

$$Y_n = X_1 + \dots + X_n$$

has Poisson-distribution with parameter $\lambda_1 + \dots + \lambda_n$.

22.2 Sums of continuous variables

Let X and Y be independent, continuously distributed random variables with density functions f and g respectively. Denote by F and G their cumulative

distribution functions. Let H denote the cumulative distribution function of $X + Y$. To find H pick a real number $x \in \mathbb{R}$. Then (sketch a picture!):

$$\begin{aligned} H(x) &= \int \int_{t+s < x} f(s)g(t) \, ds \, dt = \int_{-\infty}^{\infty} \int_{-\infty}^{x-s} f(s)g(t) \, dt \, ds \\ &= \int_{-\infty}^{\infty} f(s) \left(\int_{-\infty}^{x-s} g(t) \, dt \right) \, ds = \int_{-\infty}^{\infty} f(s)G(x-s) \, ds. \end{aligned}$$

By taking the derivative of H , we get the density function h of $X + Y$

$$h(x) = \int_{-\infty}^{\infty} f(s)g(x-s) \, ds$$

This formula is called the *convolution integral* of f and g .

ATTENTION! Differentiating the integral is not straightforward! Examine this rule in some simple cases!

Example 22.2 Suppose now that X and Y are independent random variables that are uniformly distributed on the interval $[0, 1]$. Then (X, Y) is uniformly distributed on the unit square of the plane. By sketching a picture, show that if h stands for the density function of $X + Y$, then

$$h(x) = \begin{cases} x & \text{if } 0 < x < 1 \\ 2 - x & \text{if } 1 < x < 2 \\ 0 & \text{elsewhere.} \end{cases}$$

Example 22.3 Let X and Y be independent, exponentially distributed random variables, both with parameter $\lambda > 0$. Let h denote the density function of $X + Y$. If f denotes the density function of the exponential distribution with parameter λ , then the convolution integral is:

$$h(x) = \int_{-\infty}^{\infty} f(s)f(x-s) \, ds$$

Behind the integral sign f is zero on the negative part of the real line. Therefore, the integrand is not zero if and only if $s > 0$ and $x - s > 0$, that is $0 < s < x$. Thus,

$$h(x) = \int_0^x \lambda^2 e^{-\lambda s} e^{-\lambda(x-s)} \, ds = \lambda^2 \int_0^x e^{-\lambda x} \, ds = \lambda^2 x e^{-\lambda x}$$

for any given $x > 0$, since the last integrand does not depend on s .

By using induction, we can extend the above result to any finite number of terms.

Theorem 22.4 Assume that X_1, \dots, X_n are independent, exponentially distributed random variables with the same parameter $\lambda > 0$. Let h_n denote the

density function of the random variable

$$Y_n = X_1 + \dots + X_n$$

Then

$$h_n(x) = \frac{\lambda^n}{(n-1)!} x^{n-1} e^{-\lambda x}$$

if $x > 0$, and $h_n(x) = 0$, if $x \leq 0$.

22.3 The Poisson process

In this section we describe a deeper relationship between the exponential and the Poisson distributions.

Consider the random variables T_1, T_2, \dots which mean waiting times between consecutive "occurrences".

We can think of times between successive vehicles on a highway, times between incoming claims received by an insurance company, waiting times between consecutive clients at a customer service desk, time intervals between incoming calls to a call center, etc.

Assume that T_1, T_2, \dots independent, exponentially distributed random variables with identical parameter $\lambda > 0$. The smaller the value of λ , the longer are the expected waiting times (check the expectation!). The memoryless property of the exponential distribution means that the waiting time is independent on how long we have been waiting before.

Set $S_0 = 0$ denote by

$$S_n = T_1 + \dots + T_n$$

the total waiting time until the n -th occurrence. For a given $t > 0$ the event

$$\{S_n \leq t\}$$

means that the n -th occurrence arrives before t . This means that the number of occurrences in the time interval $[0, t]$ is at least n .

Denote by $N(t)$ the number of occurrences in the time interval $[0, t]$, then the events

$$\{N(t) \geq n\} = \{S_n \leq t\}$$

coincide. For every $t > 0$ we defined a random variable $N(t)$, this correspondence is called the *Poisson process*.

How can we find the distribution of $N(t)$ for a fixed $t > 0$? The event that there are exactly n occurrences in the time interval $[0, t]$ is given by

$$\{N(t) = n\} = \{S_n \leq t\} \cap \overline{\{S_{n+1} \leq t\}} = \{S_n \leq t < S_{n+1}\}.$$

Clearly $\{S_{n+1} \leq t\} \subset \{S_n \leq t\}$, and this implies

$$P(N(t) = n) = P(S_n \leq t) - P(S_{n+1} \leq t).$$

Let h_n be the density of S_n , and h_{n+1} be the density of S_{n+1} . Since T_1, T_2, \dots are independent, exponentially distributed random variables with the same parameter λ , then in view of Theorem 22.4 of the previous section we get

$$h_n(x) = \frac{\lambda^n}{(n-1)!} x^{n-1} e^{-\lambda x} \quad \text{and} \quad h_{n+1}(x) = \frac{\lambda^{n+1}}{n!} x^n e^{-\lambda x}$$

for every $x > 0$. Therefore

$$P(N(t) = n) = P(S_n \leq t) - P(S_{n+1} \leq t) = \int_0^t h_n(x) dx - \int_0^t h_{n+1}(x) dx.$$

Evaluate the first integral on the right-hand side by integration by parts:

$$\begin{aligned} \int_0^t h_n(x) dx &= \frac{\lambda^n}{(n-1)!} \int_0^t x^{n-1} e^{-\lambda x} dx \\ &= \frac{\lambda^n}{(n-1)!} \left[\frac{x^n}{n} e^{-\lambda x} \right]_0^t + \frac{\lambda^n}{(n-1)!} \int_0^t \frac{x^n}{n} \lambda e^{-\lambda x} dx \\ &= \frac{(\lambda t)^n}{n!} e^{-\lambda t} + \frac{\lambda^{n+1}}{n!} \int_0^t x^n e^{-\lambda x} dx. \end{aligned}$$

We can recognize that in the last integral we precisely have h_{n+1} . Hence,

$$P(N(t) = n) = \frac{(\lambda t)^n}{n!} e^{-\lambda t}$$

Theorem 22.5 *In the Poisson process the number of occurrences in the time interval $[0, t]$ is a Poisson random variable with parameter λt .*

22.4 Sum of standard normal distributions

Let Z_1 and Z_2 be independent, standard normally distributed random variables, and find the distribution of their sum:

$$Y = Z_1 + Z_2$$

Now, the convolution integral is

$$h(x) = \int_{-\infty}^{\infty} \varphi(s) \varphi(x-s) ds$$

where h is the density function of Y . Then

$$\begin{aligned} h(x) &= \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-s^2/2} e^{-(x-s)^2/2} ds = \frac{1}{2\pi} e^{-\frac{x^2}{2}} \int_{-\infty}^{\infty} e^{xs-s^2} ds \\ &= \frac{1}{2\pi} e^{-\frac{x^2}{2}} \int_{-\infty}^{\infty} e^{-(s-x/2)^2} e^{x^2/4} ds = \frac{1}{2\pi} e^{-\frac{x^2}{4}} \int_{-\infty}^{\infty} e^{-(s-x/2)^2} ds \end{aligned}$$

The last integral is precisely the Gauss integral, whose value is $\sqrt{\pi}$, thus

$$h(x) = \frac{1}{2\sqrt{\pi}} e^{-\frac{x^2}{4}} \quad -\infty < x < \infty$$

This is exactly the density function of the normal distribution with parameters $m = 0$ and $\sigma = \sqrt{2}$.

Using completely analogous arguments, we can formulate the following result.

Theorem 22.6 *Let Z_1, \dots, Z_n be independent, standard normally distributed random variables. Then $Y = Z_1 + \dots + Z_n$ is a normally distributed random variable with parameters $m = 0$ and $\sigma = \sqrt{n}$.*

22.5 Central Limit Theorem

Imagine the following experiment. To determine an unknown quantity m we carry out n independent observations (measurements). To approximate the unknown quantity we use the arithmetic mean (average) of the n outcomes.

Let us denote the outcomes by X_1, \dots, X_n and assume that they are independent and identically distributed random variables with

$$E(X_k) = m, \quad D(X_k) = \sigma, \quad k = 1, 2, \dots, n$$

(Identical distribution means that the observations are performed in identical circumstances.) For the standardized average let us introduce the following notation:

$$Y_n = \frac{\frac{1}{n}(X_1 + \dots + X_n) - m}{\sigma/\sqrt{n}}$$

Then Y_n has a mean of 0 and standard deviation 1.

It was the amazing discovery of the Russian mathematician Alexandr Lyapunov and the mathematics of his time (early 20-th century) that the distribution of this variable Y_n converges to the standard normal distribution.

Tétel 22.7 (Central Limit Theorem) *Under the above conditions let F_n denote the cumulative distribution function of Y_n . Then for every $x \in \mathbb{R}$ we have*

$$\lim_{n \rightarrow \infty} F_n(x) = \Phi(x).$$

Example 22.8 On a given day the number of visitors to a local convenience store is 100. Every visitor buys something with probability $p = 0.2$ (independently from each other). Find the probability that on that given day the the number of purchases will be between 15 and 25.

Let X be the number of purchases. Then X is binomially distributed (Bernoulli-experiment!) with parameters $n = 100$ and $p = 0.2$. For each visitor introduce the following notation:

$$X_k = \begin{cases} 0 & \text{if does not buy anything} \\ 1 & \text{if buys something} \end{cases}$$

then $X = X_1 + \dots + X_{100}$ and the terms are independent random variables. It is easy to see that for each k we have $E(X_k) = 0.2$ and $Var(X_k) = 0.16$, hence $D(X_k) = 0.4$. Therefore,

$$\begin{aligned} P(15 < X < 25) &= P\left(-\frac{5}{4} < \frac{X - 20}{4} < \frac{5}{4}\right) \\ &= P\left(-\frac{5}{4} < \frac{\frac{1}{100}(X_1 + \dots + X_{100}) - 0.2}{0.4/10} < \frac{5}{4}\right) \end{aligned}$$

Making use of the Central Limit Theorem

$$\begin{aligned} P(15 < X < 25) &\approx \Phi(1.25) - \Phi(-1.25) \\ &= 2\Phi(1.25) - 1 = 0.7888 \end{aligned}$$

by looking up the number in the table for the standard normal distribution, see Textbook-2, page 736 (Appendix A).

Recitation and Exercises

1. Reading: Textbook-2, Sections 6.5 and 6.6.
2. Homework: Textbook-2, Exercises 6.24, 6.26, 6.29, 6.34 and 6.38.
3. Review: "Probability Exercises"

Chapter 23

Law of Large Numbers

23.1 Chebyshev's Theorem

So far we have had to determine probabilities of the form

$$P(a < X < b)$$

This is easy to do if the distribution of the random variable X is known. In particular, in the case of a discrete variable we get

$$P(a < X < b) = \sum_{a < x_k < b} P(X = x_k)$$

while for a continuously distributed variable

$$P(a < X < b) = \int_a^b f(x) dx$$

where f is the density function of X . However, there are situations when this procedure cannot be completed. Namely, if

1. either the distribution of X is not known,
2. or the distribution of X is known, but too complicated to use.

In cases like these, we can be satisfied with an appropriate estimate on the given probability. This estimate is provided by Chebyshev's Theorem. Consider a random variable X that has a mean and a variance.

Theorem 23.1 (Chebyshev's Theorem) *The mean of X is $E(X) = m$ and its standard deviation is $D(X) = \sigma$. Then*

$$P(|X - m| < k \cdot \sigma) \geq 1 - \frac{1}{k^2}$$

for any $k > 0$.

Proof. We present the proof for a continuously distributed random variable. In the discrete case the proof can be carried out in a completely analogous way. Let f be the density function of X , then

$$\sigma^2 = \int_{-\infty}^{\infty} (x - m)^2 f(x) dx$$

If $k > 0$ is given, then the value of the integral on the right-hand side will not increase if we skip the interval $[m - k\sigma, m + k\sigma]$. In fact:

$$\sigma^2 \geq \int_{-\infty}^{m-k\sigma} (x - m)^2 f(x) dx + \int_{m+k\sigma}^{\infty} (x - m)^2 f(x) dx \quad (23.1)$$

since the integrand is nonnegative. On the other hand, at every point x of the interval $(-\infty, m - k\sigma]$ we have $(x - m)^2 \geq k^2\sigma^2$, and hence

$$\int_{-\infty}^{m-k\sigma} (x - m)^2 f(x) dx \geq \int_{-\infty}^{m-k\sigma} k^2\sigma^2 f(x) dx \geq k^2\sigma^2 P(X \leq m - k\sigma).$$

Completely similarly, at every point x of the interval $[m + k\sigma, \infty)$ we get $(x - m)^2 \geq k^2\sigma^2$, and consequently

$$\int_{m+k\sigma}^{\infty} (x - m)^2 f(x) dx \geq \int_{m+k\sigma}^{\infty} k^2\sigma^2 f(x) dx \geq k^2\sigma^2 P(X \geq m + k\sigma).$$

If we combine the latter two inequalities with the inequality (23.1), then we obtain

$$\sigma^2 \geq k^2\sigma^2 P(X \leq m - k\sigma) + k^2\sigma^2 P(X \geq m + k\sigma).$$

Dividing both sides with the positive expression $k^2\sigma^2$ we get

$$\frac{1}{k^2} \geq P(X \leq m - k\sigma) + P(X \geq m + k\sigma) = P(|X - m| \geq k\sigma).$$

By converting to the complement event, the proof is completed. \square

Note that the theorem gives an irrelevant result if $k \leq 1$, so we apply the inequality only for $k > 1$.

Example 23.2 For instance, if the distribution of the random variable X is not known, but its mean $E(X) = 8$ and its standard deviation $D(X) = 2$ are given, then

$$P(2 < X < 14) \geq 1 - \frac{1}{9} \approx 0.8889$$

since in this case $k = 3$.

23.2 Chebyshev's Theorem in equivalent form

Sometimes it is more convenient to use Chebyshev's Theorem in the following form:

$$P(|X - E(X)| < \varepsilon) \geq 1 - \frac{\text{Var}(X)}{\varepsilon^2}$$

where $\varepsilon > 0$. Indeed, this inequality is equivalent to our theorem by setting $k \cdot D(X) = \varepsilon > 0$, and then

$$\frac{1}{k^2} = \frac{\text{Var}(X)}{\varepsilon^2}$$

Let us formulate the theorem in the following equivalent form.

Theorem 23.3 *Consider a random variable X with a mean $E(X) = m$, and standard deviation $D(X) = \sigma$. Then for every fixed $\varepsilon > 0$ we have*

$$P(|X - m| < \varepsilon) \geq 1 - \frac{\sigma^2}{\varepsilon^2} \quad (23.2)$$

Example 23.4 On a given day a call center receives 2000 incoming calls. Every call is a wrong number with probability 0.002 (independently from each other). Find the probability that on that given day there are at most 8 wrong number calls.

Let X denote the number of wrong number calls. Clearly X is binomially distributed (Bernoulli experiment!), with parameters $n = 2000$ and $p = 0.002$. The solution to our problem is:

$$P(X \leq 8) = \sum_{k=0}^8 \binom{2000}{k} 0.002^k \cdot 0.998^{2000-k}$$

which is not easy to evaluate (although the distribution is known).

However, we can give a reasonable estimate by using Chebyshev's Theorem. Now $m = 4$ and $\sigma^2 = 4 \cdot 0.998 \approx 4$, and therefore

$$P(X \leq 8) = P(|X - 4| < 5) \geq 1 - \frac{4}{25} = 0.84$$

23.3 Poisson approximation

Example 23.5 In a large hospital with 2000 beds, the probability that a patient needs intensive care is 0.002 on any given day (independently from each other). The director wants to establish a new emergency ward so that if a patient needs intensive care, must get a bed with probability of at least 0.99. What should be the size of the emergency ward with minimal cost (smallest number of beds)?

Let N denote the number of beds in the emergency ward, and X be the number of patients who need intensive care on a given day. Then X is clearly binomially distributed (Bernoulli experiment!) with a mean of $m = 4$ and variance $\sigma^2 = 4 \cdot 0.998 \approx 4$. Then the inequality

$$P(X \leq N) = \sum_{k=0}^N \binom{2000}{k} 0.002^k 0.998^{2000-k} \geq 0.99$$

has to be solved for the smallest N (which means the lowest cost).

This is the situation when the distribution of X is known, but too complicated to use. Apply Chebyshev's Theorem instead:

$$P(|X - 4| < \varepsilon) \geq 1 - \frac{4}{\varepsilon^2} = 0.99$$

The lowest solution is $\varepsilon = 20$ and therefore $N = 23$ is obtained for the optimal smallest number of beds in the new emergency ward.

Chebyshev's Theorem is true for any distribution, so we cannot expect a very sharp estimate. We can get a much more accurate solution if we apply the Poisson approximation. The theorem on how to approximate the binomial distribution by the Poisson distribution is discussed in Section 18.5. In particular, in the present example:

$$\sum_{k=0}^N \binom{2000}{k} 0.002^k 0.998^{2000-k} \approx \sum_{k=0}^N \frac{4^k}{k!} e^{-4}$$

since " $n = 2000$ is large enough, and $p = 0.002$ sufficiently small", moreover $np = 4$. When we look at the Poisson tables (see Textbook-2, page 732, Appendix A) we can see that the sum on the right-hand side exceeds 0.99 at $N = 9$. Based on this approximation we claim that even an emergency ward of size $N = 9$ fulfills the criteria. (Examining how sharp this approximation is, goes beyond the scope of this book.)

23.4 Law of Large Numbers

We carry out an experiment n times in a row (independently from each other) and each time we observe whether or not a given event A occurs (Bernoulli experiment).

Suppose that the probability of the event A is $P(A) = p$ (where $0 \leq p \leq 1$) and let X_n be the number of experiments in which A occurs. The quotient X_n/n means the relative frequency of the event A .

We want to examine whether the relative frequency converges to the real value of the probability when the number of experiments is increased that is $n \rightarrow \infty$?

From theoretical point of view, this question is of fundamental importance. If the answer is affirmative, it justifies our axiomatic approach to probability. Indeed, within the framework of our theory that we have developed from the axioms, we are able to derive a theorem that can directly be verified in reality. In other words, our axioms are set properly, and their consequences reflect real phenomena.

As is well known, X_n is binomially distributed with parameters n and p paraméterekkel. Pick a number $\varepsilon > 0$ and apply Chebyshev's Theorem:

$$P\left(\left|\frac{X_n}{n} - p\right| \geq \varepsilon\right) = P(|X_n - np| \geq n\varepsilon)$$

Since $E(X_n) = np$ and $Var(X_n) = np(1-p)$, we get

$$P(|X_n - np| \geq n\varepsilon) \leq \frac{np(1-p)}{n^2\varepsilon^2}$$

We have $p(1-p) \leq 1/4$ for any real number p , so from here

$$P\left(\left|\frac{X_n}{n} - p\right| \geq \varepsilon\right) \leq \frac{1}{4n\varepsilon^2} \rightarrow 0$$

if $n \rightarrow \infty$. We formulate this result in the theorem below.

Theorem 23.6 (Bernoulli's Law of Large Numbers)

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{X_n}{n} - p\right| < \varepsilon\right) = 1$$

for every $\varepsilon > 0$.

This theorem is sometimes called "Bernoulli's Weak Law of Large Numbers" to distinguish it from more advanced and complicated "Strong Law" results.

Example 23.7 A consulting agency makes a forecast of the support of a political party before the upcoming parliamentary election. They interview potential voters about their preferences. The agency wants to be 90% sure that their prediction should be within the 1% margin (i.e. the difference between the predicted ratio and real ratio is less than 1%). How many people have to be interviewed?

Let $0 < p < 1$ denote the unknown real ratio (the real support of the party), this will be estimated by the relative frequency. Assume that the size of sample (number of interviews) is n (yet to be determined) and X_n is the number of voters who support the party. Then the anticipated support ratio is X_n/n .

This is a Bernoulli experiment, therefore X_n is binomially distributed with $E(X_n) = np$ and $Var(X_n) = np(1 - p)$. Then the following inequality holds:

$$P\left(\left|\frac{X_n}{n} - p\right| \leq 0.01\right) \geq 1 - \frac{1}{4n \cdot 10^{-4}}$$

If the agency wants to guarantee this accuracy with at least 90% certainty, then

$$1 - \frac{1}{4n \cdot 10^{-4}} = 0.90$$

from which we have $n = 25\,000$.

In reality, using advanced statistical methods, even a smaller sample might be sufficient. However, in most situations it is hard to guarantee that the set of interviewed voters is homogeneous and representative (in the sense that the sample ratio reflects the ratio for the whole voting society).

Under the conditions of Theorem 23.6 the following stronger statement can also be proven.

Theorem 23.8 *Under the conditions of Theorem 23.6 we have*

$$P\left(\lim_{n \rightarrow \infty} \frac{X_n}{n} = p\right) = 1$$

Intuitively, Theorem 23.6 claims that very likely the relative frequency gets close to the probability p as n increases. However, it does not exclude that large differences can occur beyond any arbitrarily large index n . It just says that such large differences are unlikely. Theorem 23.8 tells us however, that such large differences come with probability zero. (The proof is due to Lyapunov and to Kolmogorov in a more general form in the 30's of the last century.)

Recitation and Exercises

1. Reading: Textbook-2, Section 4.4.
2. Homework: Textbook-2, Exercises 4.75, 4.76, 4.77, 4.78 and 4.91.
3. Review: "Probability Exercises"

Part III

Third Semester: Linear algebra

Chapter 25

Vector spaces and subspaces

25.1 The vector space \mathbb{R}^n

Let n be a given integer. The set \mathbb{R}^n is defined as the set of all n -tuples of real numbers that is:

$$\mathbb{R}^n = \left\{ x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} : x_1, \dots, x_n \in \mathbb{R} \right\}$$

The elements of this set are called vectors, their components are called coordinates. In a geometric interpretation this set for $n = 2$ means the plane, for $n = 3$ it means the three dimensional space.

In the sequel, vectors are denoted by lower case latin letters, real numbers (or scalars) are denoted by lower case greek letters.

For the vectors of the space \mathbb{R}^n we introduce the following operations:

Sums of vectors

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \text{ and } y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \text{ then } x + y = \begin{bmatrix} x_1 + y_1 \\ \vdots \\ x_n + y_n \end{bmatrix} \in \mathbb{R}^n$$

Vector multiplied by a scalar

$$\alpha \in \mathbb{R} \text{ and } x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \text{ then } \alpha x = \begin{bmatrix} \alpha x_1 \\ \vdots \\ \alpha x_n \end{bmatrix} \in \mathbb{R}^n$$

The set \mathbb{R}^n equipped with these operations is called a *vector space*.

Definition 25.1 Consider the vectors a_1, \dots, a_k in the vector space, and let $\alpha_1, \dots, \alpha_k$ be arbitrary real numbers (scalars). The vector

$$\alpha_1 a_1 + \dots + \alpha_k a_k$$

is called a *linear combination* of the vectors a_1, \dots, a_k .

Example 25.2 For instance, if

$$a_1 = \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix} \text{ and } a_2 = \begin{bmatrix} 3 \\ 0 \\ 4 \end{bmatrix} \text{ further } \alpha_1 = 3 \text{ and } \alpha_2 = -2$$

then

$$\alpha_1 a_1 + \alpha_2 a_2 = \begin{bmatrix} 0 \\ -3 \\ 1 \end{bmatrix}$$

25.2 Subspaces

Definition 25.3 A subset M of the vector space \mathbb{R}^n is called a *subspace*, if

- for every $x, y \in M$ we have $x + y \in M$, and
- for every $x \in M$ and $\alpha \in \mathbb{R}$ we have $\alpha x \in M$.

It is clear from the definition that a subspace always contains the zero vector 0 . The smallest subspace is $\{0\}$, the largest subspace is the whole vector space.

Theorem 25.4 If M is a subspace, then for all vectors $a_1, \dots, a_k \in M$ and all scalars $\alpha_1, \dots, \alpha_k \in \mathbb{R}$ we have

$$\alpha_1 a_1 + \dots + \alpha_k a_k \in M$$

In other words: a subspace is closed for linear combinations. It is easy to see that in the vector space \mathbb{R}^3 the straight lines and planes that pass through the origin are all subspaces.

Example 25.5 Consider the following subset in the vector space \mathbb{R}^n

$$M = \left\{ x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n : x_1 + x_2 = 0 \right\}$$

Verify that M is a subspace.

Indeed, if $x, y \in M$, then $x_1 + x_2 = 0$ and $y_1 + y_2 = 0$, and hence, for the first two coordinates of $x + y$ we have $(x_1 + y_1) + (x_2 + y_2) = 0$. This implies $x + y \in M$.

Similarly, if $x \in M$ and $\alpha \in \mathbb{R}$, then equality $x_1 + x_2 = 0$ implies that $\alpha(x_1 + x_2) = 0$, therefore, $\alpha x \in M$.

It is easily visible that for $n = 3$ the subspace M above is a plane that is perpendicular to the xy -plane and their intersection is the straight line with the angle of -45° degree. CREATE A PICTURE!

On the other hand, if in the definition of M the sum $x_1 + x_2$ were set to be any number different zero, then M would not be a subspace. In that situation the addition and the scalar multiplication may go out of M .

Theorem 25.6 *The intersection of subspaces is again a subspace.*

Proof. It is enough to prove the statement for two subspaces. The proof for any number of subspaces can be carried out analogously.

Let L and M be subspaces. If $x, y \in L \cap M$, then $x + y \in L$ and $x + y \in M$, because both are subspaces. Thus, $x + y \in L \cap M$.

Similarly, if $x \in L \cap M$ and $\alpha \in \mathbb{R}$, then $\alpha x \in L$ and $\alpha x \in M$, because both are subspaces. Consequently, $\alpha x \in L \cap M$. \square

25.3 Generated subspace

In view of Theorem 25.6 we can speak about the smallest subspace containing given vectors. This is formulated in the following definition.

Definition 25.7 The smallest subspace spanned by the vectors a_1, \dots, a_k is denoted by

$$\text{lin}\{a_1, \dots, a_k\}$$

and it is defined as the intersection of all subspaces containing these vectors. It is called the *generated (or spanned) subspace*.

A subspace contains all linear combinations of its vectors. Since all linear combinations already form a subspace, we come to the following theorem.

Theorem 25.8 *The subspace generated by the vectors a_1, \dots, a_k is the set of all linear combinations*

$$\alpha_1 a_1 + \dots + \alpha_k a_k$$

where $\alpha_1, \dots, \alpha_k \in \mathbb{R}$.

Example 25.9 For instance, in the vector space \mathbb{R}^3 consider the vectors

$$a_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \text{ and } a_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

then

$$\text{lin}\{a_1, a_2\} = \left\{ x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \in \mathbb{R}^3 : x_3 = 0 \right\}$$

i.e. the set of vectors whose third coordinate is zero. Please verify that this set is really a subspace!

25.4 Linear independence

Definition 25.10 In the vector space \mathbb{R}^n the vectors a_1, \dots, a_k are said to be *linearly independent*, if the equality

$$\alpha_1 a_1 + \dots + \alpha_k a_k = 0$$

implies $\alpha_1 = \dots = \alpha_k = 0$.

In the opposite situation the vectors are called *linearly dependent*.

Linear independence is one of the most profound concept of algebra, it formulates that a linear combination is zero **ONLY** if all coefficients are zero.

ATTENTION! The definition does not say that if all coefficients are zero, then the linear combination is also zero. This is obvious! The implication is the opposite.

In a collection of linearly independent vectors none of them can be expressed as the linear combination of the others. This stated in the following theorem.

Theorem 25.11 *The vectors a_1, \dots, a_k are linearly dependent if and only if one of them can be expressed as the linear combination of the others.*

Proof. If one vector, say a_1 can be expressed as the linear combination of the others, then

$$a_1 = \alpha_2 a_2 + \dots + \alpha_k a_k .$$

Rearrange the equality, then we get

$$-a_1 + \alpha_2 a_2 + \dots + \alpha_k a_k = 0 .$$

This shows that their linear combination is zero, although the first coefficient is not zero. Hence, the vectors cannot be linearly independent.

Conversely, assume that the vectors are linearly dependent. Then there exists a linear combination

$$\alpha_1 a_1 + \dots + \alpha_k a_k = 0 ,$$

where not all coefficients are zero, say $\alpha_1 \neq 0$. Then

$$a_1 = -\frac{\alpha_2}{\alpha_1} a_2 - \dots - \frac{\alpha_k}{\alpha_1} a_k ,$$

that means a_1 can be expressed as the linear combination of the others. \square

Example 25.12 Consider the following vectors in \mathbb{R}^3

$$a_1 = \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix} \quad a_2 = \begin{bmatrix} 3 \\ 0 \\ 4 \end{bmatrix} \quad a_3 = \begin{bmatrix} 1 \\ -2 \\ 2 \end{bmatrix}$$

and find out if they are independent.

Easy calculation shows that $a_3 = 2a_1 - a_2$, so the vectors are dependent.

The statements below follow easily from the definition. Verify them!

Theorem 25.13 Consider the vectors a_1, \dots, a_k in the vector space \mathbb{R}^n .

- If the vectors are linearly independent, then so is any subset of them.
- If the zero vector is an element of the collection, then they are linearly dependent.
- If there are two identical vectors in the collection, then they are dependent.
- If the vectors are linearly dependent, then any extension is dependent.

Proof. Just hints are given, detailed proof is a homework.

- Consider any linear combination of a subset, and insert the missing vectors with zero coefficients.

- Consider the linear combination in which the zero vector comes with the coefficient 1, and all other vectors with 0.
- Consider the linear combination in which the identical vectors come with the coefficients +1 and -1 respectively, and all other vectors with 0.
- Consider a linear combination which is zero, but not all coefficients are zero, and insert the vectors in the extension with zero coefficients.

Example 25.14 Suppose the vectors a , b and c are linearly independent. Is it true that the vectors $a + b$, $b + c$, $c + a$ are linearly independent as well?

Take a linear combination, and make it equal zero:

$$\alpha_1(a + b) + \alpha_2(b + c) + \alpha_3(c + a) = 0.$$

Rearrange the equality this way:

$$(\alpha_1 + \alpha_3)a + (\alpha_1 + \alpha_2)b + (\alpha_2 + \alpha_3)c = 0$$

Independence of a_1 , a_2 , a_3 implies that

$$\alpha_1 + \alpha_3 = 0 \quad \alpha_1 + \alpha_2 = 0 \quad \alpha_2 + \alpha_3 = 0$$

The only solution of this system is $\alpha_1 = \alpha_2 = \alpha_3 = 0$. Thus, the vectors $a + b$, $b + c$ és $c + a$ are linearly independent.

Recitation and Exercises

1. Reading: Textbook-1: Sections 12.1, 12.2, 12.3 and 14.1.
2. Homework: Textbook-1, Section 14, Exercises 1, 2, 3, 4, 5 and 7.
3. Review: "Linear Algebra Exercises

Chapter 26

Linear independence and basis

26.1 Generating system

Definition 26.1 In the vector space \mathbb{R}^n a collection of vectors a_1, \dots, a_k is said to be a *generating system*, if

$$\text{lin}\{a_1, \dots, a_k\} = \mathbb{R}^n,$$

i.e. all vectors of the space can be expressed as the linear combination of the given vectors.

Example 26.2 Consider the following vectors in the vector space \mathbb{R}^3

$$a_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad a_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad a_3 = \begin{bmatrix} 3 \\ -2 \\ 0 \end{bmatrix}$$

and decide whether or not they form a generating system.

We can easily see that the whole space is not spanned by these vectors, since no vector with a nonzero third coordinate belongs to the span. These vectors are not independent either, because $a_3 = 3a_1 - 2a_2$.

On the other hand, the following set of vectors

$$e_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad e_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad e_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

forms a generating system, since every vector x can be expressed this way:

$$x = x_1e_1 + x_2e_2 + x_3e_3,$$

where x_1, x_2, x_3 are the coordinates of x . It is easy to check that these vectors are linearly independent as well.

Completely analogously, we can define the generating system of a subspace M in the vector space \mathbb{R}^n .

Definition 26.3 The vectors a_1, \dots, a_k form a generating system of the subspace M (or they span the subspace M), if every vector in M can be given as the linear combination of the vectors a_1, \dots, a_k .

26.2 Basis

Definition 26.4 A collection of vectors a_1, \dots, a_k is called a *basis* of the vector space \mathbb{R}^n , if

- they are linearly independent,
- they form a generating system.

Quite analogously, we can define the basis of a subspace.

Example 26.5 As we have seen above, in the vector space \mathbb{R}^3 the vectors

$$e_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad e_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad e_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

form a basis, as they are linearly independent and they span the whole space. Similarly, the vectors

$$a_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad a_2 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \quad a_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

form a basis as well. (VERIFY!) However, the vectors

$$a_1 = \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix} \quad a_2 = \begin{bmatrix} 0 \\ -1 \\ 0 \end{bmatrix}$$

do not form a basis, for they do not span the whole space (although linearly independent).

The following properties of a basis can be verified directly by the definition.

- A basis is a maximal linearly independent system.
- A basis is a minimal generating system.
- In a vector space every basis has the same number of elements. (ATTENTION! NOT OBVIOUS.)
- In a vector space every vector can uniquely be expressed as the linear combination of the basis.

Definition 26.6 The *standard basis* in the vector space \mathbb{R}^n is given by

$$e_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad e_2 = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \dots \quad e_n = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

and we always use the notation e_k for these vectors ($k = 1, \dots, n$).

Verify that they really form a basis! In a certain sense this is the "simplest" basis, since for every vector x we have

$$x = x_1 e_1 + \dots + x_n e_n.$$

where x_1, \dots, x_n are the coordinates of x .

26.3 Dimension

Based on the properties of a basis, we can introduce the following definition.

Definition 26.7 The dimension of a vector space or subspace M is defined as the number of elements in a maximal linearly independent system (i.e. a basis). Its notation is

$$\dim M$$

Example 26.8 In view of Example 26.6 for every integer n we have

$$\dim \mathbb{R}^n = n$$

since e_1, \dots, e_n is a maximal linearly independent system, i.e. a basis.

Example 26.9 Consider now the subspace M spanned by the vectors

$$a_1 = \begin{bmatrix} 2 \\ 2 \\ -3 \end{bmatrix} \quad a_2 = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \quad a_3 = \begin{bmatrix} 5 \\ 3 \\ -5 \end{bmatrix}$$

and find the dimension of M .

We can easily check that a_1 and a_2 are linearly independent, but a_1, a_2, a_3 are not, since

$$2a_1 + a_2 = a_3.$$

Therefore, the dimension of the generated subspace is

$$\dim M = \dim \text{lin}\{a_1, a_2, a_3\} = 2.$$

Of course, it is also true that $\dim \text{lin}\{a_1, a_2\} = 2$.

Definition 26.10 The *rank* of a collection of vectors a_1, \dots, a_k is defined as the dimension of the subspace spanned by the given vectors. Notation:

$$\text{rank}\{a_1, \dots, a_k\} = \dim \text{lin}\{a_1, \dots, a_k\}.$$

26.4 Gauss-Jordan-elimination

In this section we exhibit a very simple but powerful procedure to check quickly if a collection of vectors is linearly independent.

Consider a vector a in the vector space \mathbb{R}^n that can be given in the form

$$a = \alpha_1 e_1 + \dots + \alpha_n e_n \tag{26.1}$$

in the standard basis. Consider another vector b with

$$b = \beta_1 e_1 + \dots + \beta_n e_n \quad \text{where } \beta_1 \neq 0. \tag{26.2}$$

QUESTION: What linear combination will express the vector a , if we use the basis b, e_2, \dots, e_n instead of the standard basis, i.e. the vector e_1 is replaced by the vector b ?

REMARK: As we see, the condition $\beta_1 \neq 0$ implies that the collection b, e_2, \dots, e_n is a basis. Indeed, on the one hand, it has n elements, on the other hand b is independent of the others (i.e. to express b we need the vector e_1 as well).

Isolate the vector e_1 from the equality (26.2) the we get:

$$e_1 = \frac{1}{\beta_1}b - \frac{\beta_2}{\beta_1}e_2 - \dots - \frac{\beta_n}{\beta_1}e_n$$

and replace e_1 by this expression in equality (26.1). After rearranging we have that

$$a = \frac{\alpha_1}{\beta_1}b + \left(\alpha_2 - \frac{\alpha_1}{\beta_1}\beta_2\right)e_2 + \dots + \left(\alpha_n - \frac{\alpha_1}{\beta_1}\beta_n\right)e_n. \quad (26.3)$$

This procedure, called *Gauss-Jordan-elimination*, provides the expression of vector a with respect to the new basis b, e_2, \dots, e_n .

Example 26.11 Using Gauss-Jordan-elimination decide if the vectors

$$a = \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix} \quad b = \begin{bmatrix} 2 \\ 1 \\ -1 \end{bmatrix} \quad c = \begin{bmatrix} 3 \\ 3 \\ -4 \end{bmatrix}$$

are linearly independent. This process is as follows:

$$\left| \begin{array}{ccc|cc} \boxed{1} & 2 & 3 & 2 & 3 \\ -1 & 1 & 3 & \boxed{3} & 6 \\ 2 & -1 & -4 & -5 & -10 \end{array} \right| \begin{array}{c} -1 \\ 2 \\ 0 \end{array}$$

The calculation shows that the vectors are not independent. In particular, we obtain that the vector c can be expressed in terms of a and b , namely $c = -a + 2b$. Thus, the rank of the collection a, b, c is 2, in other words

$$\dim \text{lin}\{a, b, c\} = 2.$$

Example 26.12 Consider the following vectors in the vector space \mathbb{R}^4

$$a_1 = \begin{bmatrix} 1 \\ 0 \\ -1 \\ 2 \end{bmatrix} \quad a_2 = \begin{bmatrix} 2 \\ 1 \\ 0 \\ -1 \end{bmatrix} \quad a_3 = \begin{bmatrix} -1 \\ -2 \\ -3 \\ 8 \end{bmatrix} \quad a_4 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 3 \end{bmatrix}$$

and find the dimension of the subspace M spanned by the vectors a_1, a_2, a_3, a_4

$$M = \text{lin}\{a_1, a_2, a_3, a_4\}$$

Carry out the Gauss-Jordan-elimination process for the given vectors, then we get

$$\left| \begin{array}{cccc|ccc} \boxed{1} & 2 & -1 & 1 & 2 & -1 & 1 & 3 & -1 \\ 0 & 1 & -2 & 1 & \boxed{1} & -2 & 1 & -2 & 1 \\ -1 & 0 & -3 & 1 & 2 & -4 & 2 & 0 & 0 \\ 2 & -1 & 8 & -3 & -5 & 10 & -5 & 0 & 0 \end{array} \right|$$

The result tells us that the vectors a_3 and a_4 linearly depend on the vectors a_1 and a_2 , specifically

$$a_3 = 3a_1 - 2a_2$$

and

$$a_4 = -a_1 + a_2$$

Consequently, the maximum number of linearly independent vectors in the system of vectors a_1, a_2, a_3, a_4 is 2, namely a_1 and a_2 . Therefore, we deduce that

$$\dim M = \dim \text{lin} \{a_1, a_2, a_3, a_4\} = 2.$$

Recitation and Exercises

1. Reading: Textbook-1, Sections 12.1, 12.2, 12.3 and 14.1.
2. Homework: Textbook-1, Section 14, Exercises 1, 2, 3, 4, 5 and 7.
3. Review: "Linear Algebra Exercises"

Chapter 27

Linear mappings and matrices

27.1 Linear mappings

Let n and m be integers and consider the vector spaces \mathbb{R}^n and \mathbb{R}^m (with dimensions n and m respectively).

Definition 27.1 The map $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called a *linear map* if

- $A(x + y) = Ax + Ay$
- $A(\alpha x) = \alpha Ax$

for all vectors $x, y \in \mathbb{R}^n$ and every scalar $\alpha \in \mathbb{R}$.

If $n = m$, i.e. A maps the vector space into itself, then A is called a *linear transformation*.

We can easily check the following properties of a linear map:

- $A(\alpha x + \beta y) = \alpha Ax + \beta Ay$ for all vectors x, y and scalars α, β .
- $A0 = 0$, that is the image of the vector 0 is always the vector 0 .

Example 27.2 Below we define some mappings that map the plane \mathbb{R}^2 into itself.

1. Let A be the map that associates with every vector x its λ -multiple, i.e. $Ax = \lambda x$ ($\lambda \in \mathbb{R}$).
2. Let A be the map that associates with every x reflection with respect to the horizontal axis.

3. Let A be the map that associates with every vector its projection onto the straight line $y = x$ (the 45° line bisecting the right angle).
4. Let A be the map that associates with every vector its rotation around the origin by the angle φ (in positive direction).

In each of the above examples show that A defines a linear transformation of the plane \mathbb{R}^2 (i.e. fulfills both equalities). CREATE PICTURES!

27.2 Matrix of a linear map

In this section we take advantage of the simple fact: given a linear map, if we know the images of the basis vectors, then we can calculate the images of all vectors.

Indeed, let a linear map $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be given, and consider the standard basis e_1, \dots, e_n in the vector space \mathbb{R}^n . If we pick any vector $x \in \mathbb{R}^n$, then x is given in the standard basis:

$$x = x_1e_1 + \dots + x_ne_n.$$

Apply the map A on both sides, then by the linearity

$$Ax = x_1Ae_1 + \dots + x_nAe_n,$$

that means for Ax we only need the images Ae_1, \dots, Ae_n .

Let us denote by f_1, \dots, f_m the standard basis in the vector space \mathbb{R}^m , then we can express the vectors Ae_1, \dots, Ae_n this way:

$$\begin{aligned} Ae_1 &= a_{11}f_1 + a_{21}f_2 + \dots + a_{m1}f_m \\ Ae_2 &= a_{12}f_1 + a_{22}f_2 + \dots + a_{m2}f_m \\ &\vdots \\ Ae_n &= a_{1n}f_1 + a_{2n}f_2 + \dots + a_{mn}f_m \end{aligned}$$

If we now collect the coefficients in these equalities in chart, we obtain a *matrix* of size $m \times n$. That is called the *matrix* of the linear map A :

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ & & \vdots & \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$

that consists of m rows and n columns. The j -th column of the matrix is the vector Ae_j in the standard basis f_1, \dots, f_m of the space \mathbb{R}^m . Thus, we can obtain

the image Ax of the vector x by multiplying the matrix A by the coordinates of x this way:

$$Ax = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + \cdots + a_{1n}x_n \\ a_{21}x_1 + \cdots + a_{2n}x_n \\ \vdots \\ a_{m1}x_1 + \cdots + a_{mn}x_n \end{bmatrix} \quad (27.1)$$

i.e. the product Ax is a vector with m coordinates in the space \mathbb{R}^m .

Clearly, every linear map has a matrix representation in given bases. Conversely, it is also evident that the equality (27.1) defines a linear map on the vector space. We conclude that there is a one-to-one correspondence between linear maps and matrices.

ATTENTION! In the rest of this book we are not going to make a difference between linear maps and their matrices.

Example 27.3 Consider the linear transformations introduced in Example 27.2. Their matrices in the standard bases are (in this order):

$$\begin{aligned} 1. A &= \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} & 2. A &= \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} & 3. A &= \begin{bmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{bmatrix} \\ 4. A &= \begin{bmatrix} \cos \varphi & -\sin \varphi \\ \sin \varphi & \cos \varphi \end{bmatrix} \end{aligned}$$

Use pictures to verify these results!

27.3 Rank and degree of freedom of a matrix

Definition 27.4 Consider a linear map $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ (i.e. an $m \times n$ matrix). The range of A

$$\text{im } A = \{y \in \mathbb{R}^m : \text{van olyan } x \in \mathbb{R}^n, \text{ hogy } y = Ax\}$$

is called the *image* of A , and the set

$$\ker A = \{x \in \mathbb{R}^n : Ax = 0\}$$

is called the *kernel* of A . It is easy to see that both $\ker A$ and $\text{im } A$ are subspaces in the vector spaces \mathbb{R}^n and \mathbb{R}^m respectively.

In view of the definition the subspace $\text{im } A$ is the subspace of vectors that can be expressed as the linear combinations of the columns of A . In other words, if a_1, \dots, a_n denote the columns of A , then

$$\text{im } A = \text{lin } \{a_1, \dots, a_n\}.$$

Quite similarly, $\ker A$ is the subspace of those vectors x for which

$$Ax = x_1a_1 + \dots + x_na_n = 0,$$

where x_1, \dots, x_n denote the coordinates of x .

Example 27.5 Consider the following matrix A

$$A = \begin{bmatrix} 1 & 3 & 1 \\ -2 & 1 & 5 \\ 2 & 2 & -2 \end{bmatrix}$$

and denote the columns by a_1 , a_2 and a_3 . Using Gauss-Jordan-elimination we see that the columns are not independent, since $a_3 = -2a_1 + a_2$. Therefore, the image of A (i.e. the subspace spanned by the columns) is:

$$\operatorname{im} A = \operatorname{lin} \{a_1, a_2\}.$$

On the other hand, if we rearrange the equality above, we get

$$-2a_1 + a_2 - a_3 = 0.$$

This tells us that

$$\ker A = \operatorname{lin} \left\{ \left[\begin{array}{c} -2 \\ 1 \\ -1 \end{array} \right] \right\} = \left\{ t \cdot \left[\begin{array}{c} -2 \\ 1 \\ -1 \end{array} \right] \in \mathbb{R}^3 : t \in \mathbb{R} \right\},$$

that is, all these vectors multiplied by A result in zero vector.

Definición 27.6 The *rank* of an $m \times n$ matrix A is defined as the dimension of its image, that is:

$$\operatorname{rank} A = \dim \operatorname{im} A$$

which is equal to the maximum number of linearly independent columns of A . The *degree of freedom* of A is defined by

$$\operatorname{deg} A = \dim \ker A.$$

For instance, in the case of the matrix A in Example 27.5, we have $\operatorname{rank} A = 2$ and $\operatorname{deg} A = 1$.

Theorem 27.7 For any linear map A on \mathbb{R}^n we have $\text{rank } A + \text{deg } A = n$

Proof. Let a_1, \dots, a_k be the basis vectors of the subspace $\ker A$, and let Ab_1, \dots, Ab_m be the basis vectors of the sub space $\text{im } A$. We show that vectors

$$a_1, \dots, a_k, b_1, \dots, b_m$$

combined form a basis of the vector space \mathbb{R}^n . Indeed, on the one hand, they are linearly independent, because the equality

$$\alpha_1 a_1 + \dots + \alpha_k a_k + \beta_1 b_1 + \dots + \beta_m b_m = 0$$

multiplied by A gives us

$$\beta_1 Ab_1 + \dots + \beta_m Ab_m = 0$$

and this yields $\beta_1 = \dots = \beta_m = 0$. It follows that $\alpha_1 = \dots = \alpha_k = 0$.

On the other hand, they form a generating system: if $x \in \mathbb{R}^n$ is taken arbitrarily, then $Ax \in \text{im } A$, and hence, it can be expressed in terms of the vectors Ab_1, \dots, Ab_m

$$Ax = \beta_1 Ab_1 + \dots + \beta_m Ab_m$$

This means that $x - (\beta_1 b_1 + \dots + \beta_m b_m) \in \ker A$, so it is the linear combination of the basis vectors a_1, \dots, a_k

$$x - (\beta_1 b_1 + \dots + \beta_m b_m) = \alpha_1 a_1 + \dots + \alpha_k a_k.$$

Consequently, $k + m = n$. □

27.4 Multiplication of matrices

Suppose we are given two linear maps $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $B : \mathbb{R}^m \rightarrow \mathbb{R}^k$. Then we can consider the composition mapping $B \circ A : \mathbb{R}^n \rightarrow \mathbb{R}^k$ that we denote as the product of the two maps:

$$BA = B \circ A.$$

We can easily verify that $BA : \mathbb{R}^n \rightarrow \mathbb{R}^k$ is a linear map as well, therefore its matrix in the standard basis is of the size $k \times n$. How can we compute this matrix?

For a given index j consider the image of the vector Ae_j with the map B (i.e. the vector $B(Ae_j)$). The i -th coordinate of this image vector is:

$$b_{i1}a_{1j} + \dots + b_{im}a_{mj},$$

which is precisely the entry in the product matrix BA of the i -th row and the j -th column. Conclusion: we carry out the multiplication of the matrices so that we multiply all rows of B by the columns of A according to the above rule.

ATTENTION! The order is important! The product BA does not coincide with AB (except some special situations). It may happen that the other is not even defined.

Example 27.8 In view of the rule above, verify the following multiplication of matrices:

$$\begin{bmatrix} 2 & 3 & 1 & 0 \\ -1 & 0 & 2 & 1 \\ 1 & 1 & -1 & 1 \end{bmatrix} \cdot \begin{bmatrix} -1 & 2 & -4 \\ 1 & -1 & 1 \\ 0 & 3 & -2 \\ 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 4 & -7 \\ 2 & 5 & -1 \\ 1 & -1 & -1 \end{bmatrix}.$$

so the product matrix is of size 3×3 .

Regarding the multiplication of matrices as the composition of mappings, we the following associative property:

$$C(BA) = (CB)A \quad (27.2)$$

in all cases when the multiplication is well defined.

Example 27.9 Consider the matrix A defined with a parameter α

$$A = \begin{bmatrix} 1 & 2 & 5 \\ 2 & -1 & 0 \\ -3 & 1 & \alpha \end{bmatrix}$$

and find its rank and degree of freedom. Using Gauss-Jordan-elimination we conclude:

$$\left[\begin{array}{ccc|cc} \boxed{1} & 2 & 5 & 2 & 5 \\ 2 & -1 & 0 & \boxed{-5} & -10 \\ -3 & 1 & \alpha & 7 & \alpha + 15 \end{array} \middle| \begin{array}{c} 1 \\ 2 \\ \alpha + 1 \end{array} \right]$$

This tells us that

$$\text{rank } A = \begin{cases} 3 & \text{if } \alpha \neq -1 \\ 2 & \text{if } \alpha = -1 \end{cases}$$

and making use of Theorem 27.7 we have

$$\text{deg } A = \begin{cases} 0 & \text{if } \alpha \neq -1 \\ 1 & \text{if } \alpha = -1. \end{cases}$$

Recitation and Exercises

1. Reading: Textbook-1, Sections 12.6, 12.7, 12.8 and 14.2
2. Homework: Textbook-1, Section 12, Exercises 3, 4, 5, 6, 7, 8, and Section 14, Exercises 1, 2 and 3.
3. Review: "Linear Algebra Exercises"

Chapter 28

Linear systems

28.1 Homogeneous systems

By a homogeneous system we mean the following system of linear equations:

$$\begin{aligned} a_{11}x_1 + \dots + a_{1n}x_n &= 0 \\ &\vdots \\ a_{m1}x_1 + \dots + a_{mn}x_n &= 0 \end{aligned}$$

where the coefficients a_{ij} are given real numbers. Solve the system for the unknowns x_1, \dots, x_n .

If we compile the matrix A of coefficients and the vector x of unknowns this way:

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ & \vdots & \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

then the homogeneous system can be rewritten in the form:

$$Ax = 0 \tag{28.1}$$

whose solution set is exactly the subspace $\ker A$. The vector $x = 0$ is always a solution, if this is not the only one, we have infinitely many solutions.

Example 28.1 Find all solutions of the homogeneous system

$$\begin{aligned} x_1 + 3x_2 - 3x_3 &= 0 \\ 2x_1 + x_2 + 4x_3 &= 0 \\ x_1 + 2x_2 - x_3 &= 0 \end{aligned}$$

Now the matrix A of coefficients is

$$A = \begin{bmatrix} 1 & 3 & -3 \\ 2 & 1 & 4 \\ 1 & 2 & -1 \end{bmatrix}$$

and consider the system $Ax = 0$. Use Gauss-Jordan-elimination:

$$\left| \begin{array}{ccc|cc} \boxed{1} & 3 & -3 & 3 & -3 & 3 \\ 2 & 1 & 4 & \boxed{-5} & 10 & -2 \\ 1 & 2 & -1 & -1 & 2 & 0 \end{array} \right|$$

It shows that the columns of A are dependent. Denote the columns by a_1, a_2, a_3 , then we get

$$a_3 = 3a_1 - 2a_2, \quad \text{that is} \quad 3a_1 - 2a_2 - a_3 = 0.$$

Therefore, $x_1 = 3$, $x_2 = -2$ and $x_3 = -1$ are solutions. The whole solution set is given by: megoldás pedig az

$$x = t \cdot \begin{bmatrix} 3 \\ -2 \\ -1 \end{bmatrix}, \quad t \in \mathbb{R}$$

Clearly, in this situation we have

$$\deg A = 1 \quad \text{and} \quad \text{rank } A = 2.$$

28.2 Inhomogeneous systems

Consider the $m \times n$ matrix A , and let $b \in \mathbb{R}^m$ be a nonzero vector. The system

$$Ax = b \tag{28.2}$$

is called an *inhomogeneous system* of linear equations.

Theorem 28.2 *The system (28.2) has a solution if and only if $b \in \text{im } A$, i.e. the vector b can be expressed as a linear combinations of the columns of A .*

The solution is unique only if the columns of A are independent, i.e. its degree of freedom is zero.

Proof. We only need to prove the unicity of solution. The proof is by contradiction: if there were two solutions, x and y , then

$$A(x - y) = Ax - Ay = b - b = 0,$$

but this means that the columns of A are dependent. \square

Assume now that we know a particular solution \bar{x} of the inhomogeneous system. Then all solutions can be given by using the solution set of the homogeneous system. This is formulated in the following theorem.

Theorem 28.3 *Let \bar{x} be a given solution to the inhomogeneous system. Then all solutions can be given in the form*

$$x = \bar{x} + x_0$$

where x_0 is a solution of the homogeneous system. Conversely, if x_0 is any solution of the homogeneous system, then $\bar{x} + x_0$ solves the inhomogeneous system.

Proof. Indeed, if x is any solution of the inhomogeneous system, then consider the vector $x_0 = x - \bar{x}$. This vector solves the homogeneous system, because

$$Ax_0 = A(x - \bar{x}) = Ax - A\bar{x} = b - b = 0,$$

and we see that $x = \bar{x} + x_0$.

Conversely, take a solution x_0 of the homogeneous system arbitrarily, and set $x = \bar{x} + x_0$. Then

$$Ax = A(\bar{x} + x_0) = A\bar{x} + Ax_0 = b + 0 = b,$$

that means x solves the inhomogeneous system. \square

Example 28.4 Find all solutions of the inhomogeneous system below:

$$\begin{aligned} x_1 - x_2 + 3x_3 + 3x_4 &= 1 \\ x_1 - 2x_2 + x_3 - x_4 &= 4 \\ 2x_1 + x_2 - x_3 + 5x_4 &= 6. \end{aligned}$$

Using the notations above, we have

$$A = \begin{bmatrix} 1 & -1 & 3 & 3 \\ 1 & -2 & 1 & -1 \\ 2 & 1 & -1 & 5 \end{bmatrix} \quad \text{and} \quad b = \begin{bmatrix} 1 \\ 4 \\ 6 \end{bmatrix}.$$

The matrix A has 3 rows and 4 columns, therefore its degree of freedom is at least 1, so the solution is certainly not going to be unique (if any). Gauss-Jordan-elimination shows:

$$\left| \begin{array}{cccc|c} \boxed{1} & -1 & 3 & 3 & 1 \\ 1 & -2 & 1 & -1 & 4 \\ 2 & 1 & -1 & 5 & 6 \end{array} \right| \left| \begin{array}{ccc|c} -1 & 3 & 3 & 1 \\ \boxed{-1} & -2 & -4 & 3 \\ 3 & -7 & -1 & 4 \end{array} \right| \left| \begin{array}{cc|c} 5 & 7 & -2 \\ 2 & 4 & -3 \\ \boxed{-13} & -13 & 13 \end{array} \right| \left| \begin{array}{c|c} 2 & 3 \\ 2 & -1 \\ 1 & -1 \end{array} \right|$$

Consequently, the rank of A is 3, and its degree of freedom is 1. If we denote the columns of A by a_1, a_2, a_3, a_4 , then the above calculation tells us (look up the last two columns) that

$$a_4 = 2a_1 + 2a_2 + a_3 \quad \text{moreover} \quad b = 3a_1 - a_2 - a_3.$$

This way we can give a solution to the inhomogeneous and homogeneous system as well:

$$\bar{x} = \begin{bmatrix} 3 \\ -1 \\ -1 \\ 0 \end{bmatrix} \quad \text{and} \quad x_0 = \begin{bmatrix} 2 \\ 2 \\ 1 \\ -1 \end{bmatrix}.$$

According to our theorem above, the solution set of the inhomogeneous system is given by:

$$x = \bar{x} + t \cdot x_0 = \begin{bmatrix} 3 \\ -1 \\ -1 \\ 0 \end{bmatrix} + t \cdot \begin{bmatrix} 2 \\ 2 \\ 1 \\ -1 \end{bmatrix} \quad t \in \mathbb{R}.$$

Example 28.5 For what values of the unknown parameters α and β does the inhomogeneous system $Ax = b$ have solutions?

$$A = \begin{bmatrix} 1 & -1 & 0 \\ 2 & 1 & 3 \\ -1 & 1 & \alpha \end{bmatrix} \quad \text{and} \quad b = \begin{bmatrix} 1 \\ 2 \\ \beta \end{bmatrix}.$$

Perform the Gauss-Jordan-elimination:

$$\left[\begin{array}{ccc|ccc} \boxed{1} & -1 & 0 & 1 & -1 & 0 \\ 2 & 1 & 3 & 2 & 3 & 3 \\ -1 & 1 & \alpha & -1 & 1 & \alpha \end{array} \right] \left[\begin{array}{ccc|ccc} -1 & 0 & 1 & 1 & 1 & 1 \\ \boxed{3} & 3 & 0 & 1 & 1 & 0 \\ 0 & \alpha & \beta + 1 & \alpha & \beta + 1 & \alpha \end{array} \right]$$

Then we conclude:

- there is exactly one solution if $\alpha \neq 0$ and β is any real number,
- there are infinitely many solutions if $\alpha = 0$ and $\beta = -1$,
- there is no solution if $\alpha = 0$ and $\beta \neq -1$

The rank and degree of freedom of A depend on the parameters this way:

$$\text{rank } A = \begin{cases} 3 & \text{if } \alpha \neq 0 \\ 2 & \text{if } \alpha = 0 \end{cases}$$

and

$$\text{deg } A = \begin{cases} 0 & \text{if } \alpha \neq 0 \\ 1 & \text{if } \alpha = 0. \end{cases}$$

28.3 Inverse of a matrix

Let E denote the $n \times n$ matrix where all diagonal elements are 1, and all elements outside the diagonal are 0. This matrix is called the $n \times n$ *unit matrix*. If we regard E as linear transformation, it is the *identity*: for every vector x we have $Ex = x$.

Definition 28.6 Consider an $n \times n$ square matrix A . We say that A is *invertible*, if there exists an $n \times n$ square matrix A^{-1} such that

$$A \cdot A^{-1} = E.$$

The matrix A^{-1} is called the *inverse matrix* of A .

Clearly, A^{-1} is the inverse linear map, that is $AA^{-1}x = x$ for every $x \in \mathbb{R}^n$. It is easy to see that we also have $A^{-1}A = E$ in this case, and $(A^{-1})^{-1} = A$.

The necessary and sufficient condition for the existence of the inverse is that the map A is one-to-one. The next theorem is based on this observation.

Theorem 28.7 For an $n \times n$ matrix A the following statements are equivalent:

1. A is invertible.
2. The columns of A are independent.
3. $\ker A = \{0\}$
4. $\operatorname{im} A = \mathbb{R}^n$
5. $\operatorname{rank} A = n$
6. $\operatorname{deg} A = 0$.

Instead of checking the equivalences pairwise (there are 15 of them!), it is enough to prove the following array of implications:

$$1 \Rightarrow 2 \Rightarrow 3 \Rightarrow 4 \Rightarrow 5 \Rightarrow 6 \Rightarrow 1$$

They all come simply from the definition, and from Theorem 27.7.

ATTENTION! Please make sure you understand that the array of implications above really substitutes the pairwise equivalences. In mathematics this is a commonly used (and quick) method for proving equivalences of several statements.

In the case of invertible matrices the solution of a linear system can be given explicitly.

Theorem 28.8 *Let A be an $n \times n$ invertible matrix. Then the inhomogeneous system*

$$Ax = b$$

has a unique solution for every $b \in \mathbb{R}^n$, and the solution is given by

$$x = A^{-1}b$$

It is worth mentioning though, that in most cases solving the system by Gauss-Jordan-elimination is much faster than finding the inverse. Finding the inverse is mainly advantageous when the system has to be solved multiple times with different vectors b on the right-hand side.

28.4 Finding the inverse

Consider an $n \times n$ invertible matrix A . Finding the inverse basically means that we are looking for an $n \times n$ matrix X so that $AX = E$. If the columns of this unknown matrix X are denoted by x_1, \dots, x_n , then this process means solving n copies of inhomogeneous systems of the form

$$Ax_k = e_k$$

($k = 1, \dots, n$). This process is illustrated in the following example, where we solve all copies simultaneously with Gauss-Jordan-elimination.

Example 28.9 Is the matrix A below invertible? If yes, find the inverse.

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix}.$$

We solve the systems $Ax = e_1$, $Ax = e_2$, $Ax = e_3$ simultaneously:

$$\left[\begin{array}{ccc|ccc} \boxed{1} & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \end{array} \right] \left\| \begin{array}{ccc|ccc} 1 & 1 & & 1 & 0 & 0 \\ \boxed{1} & 1 & & 0 & 1 & 0 \\ -1 & 0 & & -1 & 0 & 1 \end{array} \right\| \left[\begin{array}{ccc|ccc} 0 & & & 1 & -1 & 0 \\ 1 & & & 0 & 1 & 0 \\ \boxed{1} & & & -1 & 1 & 1 \end{array} \right] \left\| \begin{array}{ccc|ccc} 1 & -1 & 0 & 1 & -1 & 0 \\ 1 & 0 & -1 & 1 & 0 & -1 \\ -1 & 1 & 1 & -1 & 1 & 1. \end{array} \right.$$

This shows us that

$$A^{-1} = \begin{bmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \\ -1 & 1 & 1 \end{bmatrix}.$$

We can directly verify this result by carrying out the multiplication $AA^{-1} = E$.

Example 28.10 Let A and B be $n \times n$ invertible matrices. Regarding A and B as mappings, it is obvious that the product AB is also invertible, since AB is one-to-one as well.

How can we find the inverse matrix $(AB)^{-1}$? We show that

$$(AB)^{-1} = B^{-1}A^{-1}.$$

Indeed, the matrix on the right-hand side is the inverse of the product AB , since

$$(AB)B^{-1}A^{-1} = A(BB^{-1})A^{-1} = AEA^{-1} = E$$

in view of the associative rule (27.2).

Recitation and Exercises

1. Reading: Textbook-1, Sections 13.6, 13.7 and 14.3.
2. Homework: Textbook-1, Section 13.6, Exercises 2, 3, 4, 5, 6, 8, 10, 12, Section 13.7, Exercises 2 and 4, Section 14.3, Exercises 1, 2, 3, 5 and 6.
3. Review: "Linear Algebra Exercises"

Chapter 29

Eigenvalue, eigenvector

29.1 Eigenvalue, eigenvector

Definition 29.1 Consider a linear transformation $A : \mathbb{R}^n \rightarrow \mathbb{R}^n$ that is, an $n \times n$ matrix. We say that a real number λ is an *eigenvalue* of A if there exists a vector $v \neq 0$ for which

$$Av = \lambda v.$$

In this case v is called an *eigenvector* of A associated with the eigenvalue λ .

Example 29.2 Consider the linear transformation A of the vector space \mathbb{R}^3 and the vector v , where

$$A = \begin{bmatrix} 2 & 1 & -1 \\ 0 & 1 & 1 \\ 2 & 0 & -2 \end{bmatrix} \quad \text{and} \quad v = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}.$$

It is easy to verify that

$$Av = 2v$$

which means that $\lambda = 2$ is an eigenvalue of A , and v is an associated eigenvector. On the other hand, for the vector

$$u = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$$

we have

$$Au = 0$$

thus, $\lambda = 0$ is an eigenvalue as well, and u is an associated eigenvector. Verify that $\lambda = -1$ is also an eigenvalue of A and try to find an associated eigenvector!

ATTENTION! An eigenvector that belongs to (or associated with) a given eigenvalue λ is never unique. Just think of the fact that if v is an eigenvector then so is any $\alpha \neq 0$ scalar multiple. Indeed,

$$A(\alpha v) = \alpha Av = \alpha \cdot \lambda v = \lambda(\alpha v).$$

Example 29.3 Examine the planar transformations in Example 27.2 and find their eigenvalues.

1. If A is the α -multiple of all vectors, then α is the only eigenvalue of A , and every nonzero vector of the plane is an eigenvector.
2. If A is the reflection with respect to the horizontal axis, then $\lambda_1 = 1$ is an eigenvalue, and e_1 is an associated eigenvector, further $\lambda_2 = -1$ is also an eigenvalue, and e_2 is an associated eigenvector.
3. If A is the projection onto the 45° degree line, then the eigenvalues and the corresponding eigenvectors are as follows:

$$\lambda_1 = 1 \quad \text{and} \quad v_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \text{moreover} \quad \lambda_2 = 0 \quad \text{and} \quad v_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

4. If A is the rotation around the origin by the angle $0 \leq \varphi < 2\pi$ in positive direction, then for $\varphi = 0$ we have $\lambda = 1$, while for $\varphi = \pi$ we have $\lambda = -1$ and they are the only eigenvalues. In both cases all nonzero vectors of the plane are eigenvectors.

For other angles the rotation A has no real eigenvalues.

29.2 Eigensubspace

Definition 29.4 Take a linear transformation A of the vector space \mathbb{R}^n and suppose the λ is an eigenvalue of A . All vectors v with $Av = \lambda v$ form a subspace the is called the *eigensubspace* of A associated with λ . Notation:

$$S_A(\lambda) = \{v \in \mathbb{R}^n : Av = \lambda v\}$$

Example 29.5 Consider the following matrix

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

Simple calculation shows that $\lambda = 2$ is an eigenvalue of A and e_2, e_3 are eigenvectors (linearly independent).

More independent eigenvectors associated with $\lambda = 2$ cannot be found, thus

$$\dim S_A(2) = 2$$

and e_2 and e_3 form the basis of this eigensubspace.

29.3 Finding eigenvectors

In this section we are given an $n \times n$ matrix A and suppose that λ is an eigenvalue. Question: how can we find all associated eigenvectors?

Let E denote the $n \times n$ identity matrix, and assume that v is an eigenvector that belongs to λ . Then

$$Av = \lambda v = \lambda E v$$

or by moving all terms to the left-hand side:

$$(A - \lambda E)v = 0$$

Consequently, v is a solution of a homogeneous system. This observation is formulated in the next theorem.

Theorem 29.6 *Let A be an $n \times n$ matrix, and λ is an eigenvalue of A . Then*

$$S_A(\lambda) = \ker(A - \lambda E)$$

that is the eigensubspace is given by all solutions of a homogeneous system.

The case $\lambda = 0$ is particularly interesting. In fact, if 0 is an eigenvalue, then the homogeneous system $Av = 0$ possesses a nonzero solution, and hence the rank of A cannot be n . We emphasize this fact in a separate theorem.

Theorem 29.7 *A is invertible if and only if $\lambda = 0$ is not an eigenvalue.*

In most cases it is a lot more difficult problem to find the eigenvalues of a linear transformation A .

29.4 Independent eigenvectors

Consider a linear transformation A on the vector space \mathbb{R}^n , and suppose that $\lambda_1, \dots, \lambda_k$ are all different eigenvalues of A . Take the nonzero corresponding eigenvectors v_1, \dots, v_k . We show that these vectors are linearly independent.

Theorem 29.8 *The eigenvectors that belong to different eigenvalues are linearly independent.*

Proof. We prove by induction. The statement is trivial $k = 1$. Let us suppose that the statement is true up to $k - 1$. Now we prove by contradiction: assume that v_1, \dots, v_k are dependent. This means that the vectors have a linear combination

$$\alpha_1 v_1 + \dots + \alpha_k v_k = 0 \quad (29.1)$$

where not all coefficients are zero, for simplicity, say $\alpha_1 \neq 0$. Multiply both sides by the matrix A , then we get

$$\alpha_1 A v_1 + \dots + \alpha_k A v_k = \alpha_1 \lambda_1 v_1 + \dots + \alpha_k \lambda_k v_k = 0$$

If we subtract the λ_k -multiple of equality (29.1) from this latter equality, then the last term will be cancelled, and what remains is:

$$\alpha_1 (\lambda_1 - \lambda_k) v_1 + \dots + \alpha_{k-1} (\lambda_{k-1} - \lambda_k) v_{k-1} = 0.$$

Since all eigenvalues are different, and $\alpha_1 \neq 0$, we see that the coefficient of v_1 is not zero, which means the vectors v_1, \dots, v_{k-1} are dependent. However, this contradicts to the assumption in the induction. \square

Example 29.9 Consider again the matrix A examined in Example 29.2. Simple calculation shows that $\lambda_1 = 2$, $\lambda_2 = 0$ and $\lambda_3 = -1$ are all eigenvalues of A and the associated eigenvectors are

$$v_1 = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}, \quad v_2 = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \quad \text{and} \quad v_3 = \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix}$$

respectively. Using Gauss-Jordan-elimination, verify that vectors v_1, v_2, v_3 are really linearly independent.

29.5 Diagonal form of transformations

A matrix is called diagonal if all entries outside the diagonal are zero (for example like in the identity matrix). Diagonal matrices are easy to work with (multiplication, power, etc.) that is why they are useful in linear algebra and

its applications. This justifies the important question: if a transformation is given, does there exist a basis in which its matrix becomes diagonal? As we will see in this section, such a basis consists of eigenvectors (if it exists).

Let us suppose that the eigenvalues of an $n \times n$ matrix A are $\lambda_1, \dots, \lambda_n$ (not necessarily all different), and the corresponding eigenvectors are v_1, \dots, v_n , respectively. Assume that the eigenvectors are independent. Since their number is n , they form a basis of the space \mathbb{R}^n .

Find the matrix of the transformation A in the basis of the eigenvectors! Let \hat{A} denote the matrix in this new basis. Since for all eigenvectors we have

$$Av_k = \lambda v_k \quad k = 1, \dots, n$$

this shows that \hat{A} will look like:

$$\hat{A} = \begin{bmatrix} \lambda_1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & \lambda_n \end{bmatrix}$$

Let us examine the relationship between the matrices of the transformation A with respect to the bases e_1, \dots, e_n resp. v_1, \dots, v_n . Denote by S the $n \times n$ matrix, whose columns are the the eigenvectors v_1, \dots, v_n , i.e.

$$Se_k = v_k$$

Clearly, S invertible because its columns are linearly independent. Moreover

$$Av_k = S\hat{A}e_k$$

for every index $k = 1, \dots, n$. Multiply both sides by the matrix S^{-1} , then we get

$$S^{-1}ASe_k = \hat{A}e_k$$

for every index k , therefore

$$\hat{A} = S^{-1}AS$$

We summarize these results in the following theorem.

Theorem 29.10 *Suppose that A is an $n \times n$ matrix whose eigenvalues are $\lambda_1, \dots, \lambda_n$, and the corresponding eigenvectors v_1, \dots, v_n form a basis of the space. Then the matrix of A with respect to the basis of the eigenvectors is a diagonal matrix \hat{A} with the eigenvalues in the diagonal. More specifically:*

$$\hat{A} = S^{-1}AS$$

where the columns of S are the eigenvectors v_1, \dots, v_n .

Példa 29.11 Consider again the linear transformation A in Example 29.2. For the eigenvalues $\lambda_1 = 2$, $\lambda_2 = 0$ and $\lambda_3 = -1$ the corresponding eigenvectors

are

$$v_1 = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}, \quad v_2 = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \quad \text{and} \quad v_3 = \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix}$$

These vectors are independent, therefore they form a basis of the space \mathbb{R}^3 . With respect to this basis the matrix of A will be diagonal. More specifically:

$$S = \begin{bmatrix} 2 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & -1 & 2 \end{bmatrix} \quad \text{and} \quad \hat{A} = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

and with these notations

$$\hat{A} = S^{-1}AS$$

Please verify this identity directly by determining the inverse matrix S^{-1} and by carrying out the indicated multiplications!

Example 29.12 Unfortunately, not every transformation has a diagonal form. The reason is that the eigenvectors may not form the basis of the space. For instance, in the two dimensional situation, the matrix

$$A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

has one eigenvalue: $\lambda = 1$, and there is only one independent eigenvector (for instance e_1). VERIFY!

Recitation and Exercises

1. Reading: Textbook-1, Sections 14.4 and 14.5.
2. Homework: Textbook-1, Section 14.4, Exercises 1, 2, 3, 4, 5, 6, 7, and Section 14.5, Exercises 1, 2 and 3.
3. Review: "Linear Algebra Exercises"

Chapter 30

Determinant

30.1 Permutations

Consider the set $H = \{1, \dots, n\}$ of the first n integers.

Definition 30.1 A one-to-one map $p : H \rightarrow H$ is called a *permutation* of the set H .

Intuitively, a permutation is an arrangement of the elements in H . The number of all permutations of H is $n!$ (n factorial).

Definition 30.2 Consider a permutation p of the set H , for which

$$p(1) = i_1, \quad \dots \quad p(n) = i_n$$

that is the arrangement $\{i_1, \dots, i_n\}$. We say that the elements i_j and i_k form an *inversion* if $j < k$ and $i_j > i_k$.

Example 30.3 For instance, in case of $n = 5$ the permutation

$$\{1, 3, 2, 4, 5\}$$

contains a single inversion, while in the permutation

$$\{2, 3, 1, 5, 4\}$$

we find three inversions.

Definition 30.4 A permutation p of the set H is called *odd* if the number of inversions is odd, otherwise we say that the permutation is *even*.

30.2 The determinant

Consider the following $n \times n$ matrix A

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & & & \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}$$

Definition 30.5 By the *determinant* of the matrix A we mean the following expression:

$$\det A = \sum (-1)^{\mathcal{P}} a_{1i_1} a_{2i_2} \cdots a_{ni_n}$$

where the summation is carried out for all permutations $\{i_1, \dots, i_n\}$ of the set $H = \{1, \dots, n\}$ thus, the sum contains $n!$ terms. The exponent of (-1) is odd or even, if the permutation $\{i_1, \dots, i_n\}$ is odd or even, respectively.

Some other usual widely used notations (each one is used throughout this book):

$$\det A = |A| = \begin{vmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & & & \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{vmatrix}$$

As we can see from the definition, the products behind the sum sign are compiled so that they contain precisely one factor from each row and from each column of the matrix.

Example 30.6 Verify directly by the definition that for the matrices

$$A = \begin{bmatrix} 2 & -1 \\ 3 & 4 \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} 2 & 1 & -1 \\ 0 & 1 & 3 \\ 1 & 0 & -2 \end{bmatrix}$$

we have $\det A = 11$ and $\det B = 0$.

30.3 Properties of the determinant

Theorem 30.7 $\det A = \det A^T$.

Indeed, if the rows and the columns of the matrix are interchanged, then the parity of the inversions will not change.

Theorem 30.8 *If all elements of one row of the matrix are zero, then the determinant of the matrix is zero*

Indeed, in this case all terms behind the sum sign contain a zero factor.

Theorem 30.9 *If we interchange two rows of the matrix, then the determinant of the matrix changes the sign.*

Indeed, in this case the parity of inversions in every term will change.

Theorem 30.10 *If two rows in the matrix are identical, then the determinant of the matrix is zero.*

Indeed, if we interchange the two identical rows, then on the one hand the determinant changes the sign, on the other hand it remains unchanged, thus $\det A = -\det A$. Hence, $\det A = 0$.

Theorem 30.11 *If a row of a matrix is multiplied by λ , then its determinant is multiplied by λ as well.*

Indeed, in this case every term behind the sum sign is multiplied by λ , since every term contains precisely one factor from each row.

Theorem 30.12 *If in a matrix one row is a λ -multiple of another row, then the determinant of the matrix is zero.*

Indeed, if λ is factored out from the matrix, then we obtain a matrix with two identical rows.

Theorem 30.13 *If in a matrix the i -th row is given in the form of a sum like*

$$a_{ij} = b_{ij} + c_{ij} \quad j = 1, \dots, n$$

then its determinant is the sum of the two determinants with i -th rows of elements b_{ij} and c_{ij} respectively.

Indeed, behind the sum sign every product is the sum of such terms.

Theorem 30.14 *If in a matrix one row is the linear combination of the other rows, then its determinant is zero.*

Indeed, divide the determinant into a sum as in the preceding theorem. If we now factor out the scalar coefficients from each term, we obtain two identical rows in each determinant. Therefore, each of them is zero.

Theorem 30.15 *If in a matrix we add a λ scalar multiple of a row to another row, then the determinant remains unchanged.*

Indeed, in this case the determinant can be divided into a sum of two determinants, where the second term is zero.

Theorem 30.16 *If A and B are $n \times n$ matrices, then $\det(AB) = \det A \cdot \det B$.*

This statement can be verified by carrying out the matrix multiplication step-by-step, and by exploiting our previous theorems.

Finally, as a consequence, we can formulate the following fundamental property.

Theorem 30.17 *The columns of a square matrix A are linearly dependent if and only if $\det A = 0$.*

The necessity is an immediate corollary of Theorem 30.14. The sufficiency comes from the fact that if the columns of A are linearly independent, then A is invertible, and hence $AA^{-1} = E$. Thus,

$$\det(AA^{-1}) = \det A \cdot \det A^{-1} = \det E = 1$$

and consequently, $\det A \neq 0$.

30.4 Evaluating the determinant

We conclude directly from the definition that the determinant of a 2×2 matrix is given by

$$\det A = a_{11}a_{22} - a_{12}a_{21}$$

Completely analogously, the determinant of a 3×3 matrix can be evaluated like

$$\det A = a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} - a_{12} \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} + a_{13} \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}$$

If we apply this observation inductively, we come to the following result.

Theorem 30.18 *Consider the $n \times n$ matrix A , and denote by A_{1j} the $(n-1) \times (n-1)$ matrix that we obtain by discarding the first row and the j -th column of A . Then*

$$\det A = \sum_{j=1}^n (-1)^{j+1} a_{1j} \det A_{1j}$$

This procedure is called the division into subdeterminants.

Example 30.19 Apply the division into subdeterminants procedure for the matrix

$$A = \begin{bmatrix} 3 & 6 & 0 & 2 \\ 0 & 1 & 2 & 4 \\ 4 & 8 & 3 & 5 \\ 1 & 2 & 0 & 0 \end{bmatrix}$$

Step-by-step, first taking the 3×3 subdeterminants, then the 2×2 subdeterminants, verify that we finally get $\det A = -6$.

30.5 Finding the eigenvalues

Consider an $n \times n$ matrix A , and suppose that λ is an eigenvalue with a corresponding eigenvector $v \neq 0$, i.e. $Av = \lambda v$. This can be rewritten like

$$Av - \lambda v = (A - \lambda E)v = 0.$$

This equality means that the columns of the matrix $A - \lambda E$ are linearly dependent, since the homogeneous system possesses a nonzero solution. Making use of Theorem 30.17 we come to the following conclusion.

Theorem 30.20 *The scalar λ is an eigenvalue of the matrix A if and only if $\det(A - \lambda E) = 0$.*

This necessary and sufficient condition ultimately means finding the roots of the n -th degree polynomial $\det(A - \lambda E)$. This polynomial is called the *characteristic polynomial* of the matrix A .

Example 30.21 Consider the matrix

$$A = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

and find the eigenvalues. Expand the determinant of the matrix $A - \lambda E$

$$\det(A - \lambda E) = (\lambda + 1)^2(\lambda - 2)$$

The roots of this cubic polynomial are $\lambda_1 = -1$ (with multiplicity 2), and $\lambda_2 = 2$, and they are the eigenvalues of A . As a routine calculation, find the corresponding eigenvectors as well. The linearly independent solutions of the homogeneous system

$$(A - \lambda E)v = 0$$

are as follows. For $\lambda = -1$ we have

$$v_1 = \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix} \quad \text{and} \quad v_2 = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$

as well as for $\lambda = 2$ we get

$$v_3 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

The vectors v_1, v_2, v_3 form a basis of \mathbb{R}^3 , and the matrix A admits the diagonal form

$$\hat{A} = \begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

in this basis.

Recitation and Exercises

1. Reading: Textbook-1, Sections 13.1, 13.2, 13.3, 13.4 and 13.5.
2. Homework: Textbook-1, Section 13.2, Exercises 2, 6, Section 13.4, Exercises 4, 6, 8, Section 13.5, Exercises 2, 3, Section 14.4, Exercises 1, 2, 3 and 4.
3. Review: "Linear Algebra Exercises"

Chapter 31

Scalar product

31.1 Scalar product

Definition 31.1 The *scalar product* of the vectors x and y in the vector space \mathbb{R}^n is defined by

$$\langle x, y \rangle = x_1y_1 + \dots + x_ny_n$$

where on the right-hand side we have the coordinates of the vectors.

Please observe that in the case of $n = 2$ this concept coincides with the one studied in highschool.

Definition 31.2 The *norm* or absolute value of a vector $x \in \mathbb{R}^n$ is defined by

$$\|x\| = (\langle x, x \rangle)^{1/2} = \sqrt{x_1^2 + \dots + x_n^2}$$

that we also call the length of the vector.

Clearly, in view of the Pythagorean theorem, this concept complies with our geometric intuition. It is also easy to see that $\|x\| = 0$ if and only if $x = 0$.

Definition 31.3 The *distance* of the vectors x and y is defined by $\|x - y\|$.

Example 31.4 For instance, if we consider the vectors

$$x = \begin{bmatrix} 2 \\ -2 \\ 1 \end{bmatrix} \quad \text{and} \quad y = \begin{bmatrix} 4 \\ -3 \\ 0 \end{bmatrix}$$

then their scalar product is $\langle x, y \rangle = 14$, and their norms are $\|x\| = 3$ and $\|y\| = 5$. The distance of the two vectors is $\|x - y\| = \sqrt{6}$. VERIFY!

31.2 Angle of vectors, perpendicularity

Theorem 31.5 Cauchy-Schwarz-inequality For all vectors $x, y \in \mathbb{R}^n$ we have

$$|\langle x, y \rangle| \leq \|x\| \cdot \|y\|$$

Proof. Let t be an arbitrarily chosen real number, and consider the following quadratic polynomial:

$$g(t) = \langle x + ty, x + ty \rangle$$

On the one hand, by the definition of the scalar product we get:

$$g(t) = \|x\|^2 + 2t\langle x, y \rangle + t^2\|y\|^2$$

on the other hand this is the square of the norm of the vector $x + ty$, therefore it is nonnegative, i.e.

$$g(t) \geq 0 \quad \text{for every } t.$$

If a quadratic polynomial is nonnegative, then its discriminant is nonpositive that is:

$$4(\langle x, y \rangle)^2 \leq 4\|x\|^2 \cdot \|y\|^2$$

Dividing by 4, and taking the square root of both sides we have

$$|\langle x, y \rangle| \leq \|x\| \cdot \|y\|$$

and this is exactly that we had to prove. \square

Theorem 31.6 (Triangle-inequality) $\|x + y\| \leq \|x\| + \|y\|$.

Proof. Indeed, in view of the Cauchy-Schwarz-inequality we get

$$\begin{aligned} \|x + y\|^2 &= \langle x + y, x + y \rangle = \|x\|^2 + 2\langle x, y \rangle + \|y\|^2 \\ &\leq \|x\|^2 + 2\|x\| \cdot \|y\| + \|y\|^2 = (\|x\| + \|y\|)^2 \end{aligned}$$

and by taking the square root of both sides, the statement ensues. \square

Definition 31.7 The *angle* of the nonzero vectors x and y is defined as the angle $0 \leq \varphi \leq \pi$ for which

$$\cos \varphi = \frac{\langle x, y \rangle}{\|x\| \cdot \|y\|}$$

Moreover, we say that the vectors x and y are *orthogonal* (or perpendicular), notation: $x \perp y$, if

$$\langle x, y \rangle = 0$$

Obviously, in this case $\cos \varphi = 0$, i.e. $\varphi = \pi/2$. The vector 0 is orthogonal to any other vector.

ATTENTION! Please observe that the angle of vectors is well defined. Indeed, in view of the Cauchy-Schwarz-inequality we see that $-1 \leq \cos \varphi \leq 1$.

31.3 Orthogonal systems

Definition 31.8 We say that the vectors a_1, \dots, a_k form an *orthogonal system* in the vector space \mathbb{R}^n , if none of them is the zero vector and they are pairwise orthogonal, i.e.

$$\langle a_i, a_j \rangle = 0$$

for all indices $i \neq j$.

Theorem 31.9 *Every orthogonal system is linearly independent.*

Proof. Consider the orthogonal system a_1, \dots, a_k , and suppose that

$$\alpha_1 a_1 + \dots + \alpha_k a_k = 0.$$

Take scalar product of both sides by the vector a_i vektorral. By the pairwise orthogonality each product is zero except the i -th term, we get

$$\alpha_i \|a_i\|^2 = 0$$

Since $a_i \neq 0$, we conclude that $\alpha_i = 0$. This argument can be applied for all indices $i = 1, \dots, k$ therefore, we have that $\alpha_1 = \dots = \alpha_k = 0$. This exactly means that the vectors a_1, \dots, a_k are linearly independent. \square

The above result tells us that in the space \mathbb{R}^n the maximum number of elements of an orthogonal system is n . At the same time, an orthogonal system with n elements forms the basis of the whole space.

Definition 31.10 By an *orthogonal basis* of the space \mathbb{R}^n we mean an orthogonal system a_1, \dots, a_n . We say that this basis is *orthonormal*, if all basis vectors have unit length, i.e. $\|a_i\| = 1$ for every index $i = 1, \dots, n$.

The orthogonal or orthonormal basis of a subspace M is defined completely analogously.

Example 31.11 For instance, it is easy to check that the vectors

$$a_1 = \begin{bmatrix} \frac{\sqrt{3}}{2} \\ 0 \\ -\frac{1}{2} \end{bmatrix}, \quad a_2 = \begin{bmatrix} \frac{1}{2} \\ 0 \\ \frac{\sqrt{3}}{2} \end{bmatrix}, \quad a_3 = \begin{bmatrix} 0 \\ -1 \\ 0 \end{bmatrix}$$

form an orthonormal basis of the space \mathbb{R}^3 .

31.4 Gram-Schmidt-procedure

Vectors given in an orthonormal basis are easy to work with (think of scalar products!), so it is a natural question to ask if there exists an orthonormal basis in any subspace. An affirmative answer is given by the Gram-Schmidt-procedure, which even provides an algorithm showing how to create this basis.

Consider a subspace M in the vector space \mathbb{R}^n and suppose that the vectors a_1, \dots, a_k form a basis of M . Starting with this basis, we show how we can construct an orthonormal basis of M .

Put $b_1 = a_1$. Then set $b_2 = a_2 + \alpha_1 b_1$, where the unspecified scalar α_1 is chosen so that b_2 becomes orthogonal to the vector b_1 . This means

$$\langle b_2, b_1 \rangle = \langle a_2, b_1 \rangle + \alpha_1 \langle b_1, b_1 \rangle = 0$$

From this equality we get

$$\alpha_1 = -\frac{\langle b_1, a_2 \rangle}{\|b_1\|^2}$$

and therefore

$$b_2 = a_2 - \frac{\langle b_1, a_2 \rangle}{\|b_1\|^2} b_1.$$

Very similarly, we look for the vector b_3 in the form

$$b_3 = a_3 + \beta_1 b_1 + \beta_2 b_2$$

where the unspecified scalars should be chosen so that b_3 becomes orthogonal to both vectors b_1 and b_2 . Finding the scalars from the two conditions, we obtain Ekkor a két feltételből az együtthatókat meghatározva

$$b_3 = a_3 - \frac{\langle b_1, a_3 \rangle}{\|b_1\|^2} b_1 - \frac{\langle b_2, a_3 \rangle}{\|b_2\|^2} b_2.$$

By continuing this process, finally we come to an orthogonal basis of the subspace M . We formulate this result in the theorem below.

Theorem 31.12 (Gram-Schmidt) *In every subspace of \mathbb{R}^n there exists an orthonormal basis.*

Proof. In the construction above divide each vector b_i by the positive scalar $\|b_i\|$, then we obtain an orthonormal basis. \square

31.5 Orthogonal complement

Consider a subspace M in the vector space \mathbb{R}^n .

Definition 31.13 The set of all vectors that are orthogonal to every vector of M is denoted by

$$M^\perp = \{y \in \mathbb{R}^n : \langle y, x \rangle = 0 \text{ for every } x \in M\}$$

and is called the *orthogonal complement* of M .

We can easily check that M^\perp is a subspace in \mathbb{R}^n .

Példa 31.14 In the three dimensional space the orthogonal complement of a straight line that passes through the origin is the perpendicular plane passing through the origin. Conversely, the orthogonal complement of a plane is the perpendicular line cutting the plane in the origin.

For instance, if M is a subspace spanned by two vectors:

$$M = \text{lin} \left\{ \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \right\}$$

then the orthogonal complement is the subspace spanned by a single vector:

$$M^\perp = \text{lin} \left\{ \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \right\}$$

The converse statement is also true, this is claimed in the following theorem.

Theorem 31.15 For any subspace M we have $(M^\perp)^\perp = M$.

The statement follows easily from the definition.

Theorem 31.16 Pick a vector $a \in \mathbb{R}^n$ and consider a subspace M . Then there exists exactly one vector $u \in M$ for which

$$a - u \in M^\perp$$

Proof. Take an orthonormal basis b_1, \dots, b_k in the subspace M . We try to find the vector u of the subspace M in the form:

$$u = \alpha_1 b_1 + \dots + \alpha_k b_k$$

The unspecified scalar coefficients have to be chosen so that $a - u$ is orthogonal to each basis vector. This means the following equalities:

$$\langle b_i, a - u \rangle = \langle b_i, a \rangle - \alpha_i = 0$$

for all indices $i = 1, \dots, k$. The unknown scalars are uniquely determined by these equalities. \square

This vector u is called the *orthogonal projection* of a onto the subspace M .

Theorem 31.17 *Let M be a subspace in \mathbb{R}^n . Then every vector $a \in \mathbb{R}^n$ can uniquely be given in the form*

$$a = u + v$$

where $u \in M$ and $v \in M^\perp$.

Bizonyítás. Let $u \in M$ denote the orthogonal projection of a onto the subspace M . Then the vector $v = a - u$ is orthogonal to M (verify!), and consequently $v \in M^\perp$.

The unicity comes from the fact that if we have

$$a = u' + v'$$

for another two vectors, then by subtracting the second equality from the first, we get $u - u' = v - v'$. This implies $u - u' \in M$ and $u - u' \in M^\perp$. Thus, $u - u'$ is orthogonal to itself, which means

$$0 = \langle u - u', u - u' \rangle = \|u - u'\|^2$$

This is only possible if $u - u' = 0$ and similarly $v - v' = 0$. □

Theorem 31.18 *Let M be any subspace in \mathbb{R}^n . Then*

$$\dim M + \dim M^\perp = n.$$

Proof. Take an orthonormal basis u_1, \dots, u_k in the subspace M , and take an orthonormal basis v_1, \dots, v_m in the subspace M^\perp . Then, in view of our previous theorem, the collection of vectors

$$u_1, \dots, u_k, v_1, \dots, v_m$$

spans the whole vector space, i.e. it is a generating system, and hence, $k+m \geq n$. On the other hand the vectors of this collection are pairwise orthogonal, so they are linearly independent (see Theorem 31.9), therefore $k+m \leq n$. We conclude that $k+m = n$. □

Recitation and Exercises

1. Reading: Textbook-1, Section 12.4 and 12.5.
2. Homework: Textbook-1, Section 12.4, Exercises 3, 4, 5, 6, 7, 8, 9, Section 12.5, Exercises 2, 3 and 4.
3. Review: "Linear Algebra Exercises"

Chapter 32

The spectral theorem

32.1 Transpose of a matrix

Consider a linear transformation A of the space \mathbb{R}^n , i.e. an $n \times n$ matrix.

Definition 32.1 The *transpose* of A is the linear transformation A^T for which

$$\langle A^T y, x \rangle = \langle y, Ax \rangle$$

for every $x, y \in \mathbb{R}^n$.

What does the matrix A^T look like? Apply the definition specifically on the basis vectors, then

$$\langle A^T e_i, e_j \rangle = \langle e_i, Ae_j \rangle$$

for all indices i and j . On the right-hand side of the equality we have a_{ij} , which is the element in the i -th row and j -th column of A , while on the left-hand side we get the element in the j -th row and i -th column of A^T . Therefore the matrix A^T is obtained by interchanging the rows and the columns of A .

In other words, we may also say that the matrix A^T is created by reflecting the elements of A with respect to the diagonal. Clearly, $(A^T)^T = A$.

Theorem 32.2 For any square matrix A we have

$$\ker A = (\operatorname{im} A^T)^\perp$$

Proof. On the one hand, if a vector x is orthogonal to the subspace $\operatorname{im} A^T$, then

$$0 = \langle A^T y, x \rangle = \langle y, Ax \rangle$$

for every vector y . This is only possible if $Ax = 0$, and it means $x \in \ker A$.

Conversely, in view of the above equality, every vector in $\ker A$ is orthogonal to the subspace $\operatorname{im} A^T$ \square

This observation leads us to the Rank Theorem of matrices.

Theorem 32.3 (Rank theorem of matrices) *For any $n \times n$ matrix A we have*

$$\operatorname{rank} A = \operatorname{rank} A^T$$

Proof. Indeed, the previous theorem and Theorem 31.18 imply that

$$\dim \ker A + \dim \operatorname{im} A^T = n$$

and in view of Theorem 27.7 Tétel we get $\dim \operatorname{im} A = \dim \operatorname{im} A^T$, that proves our statement. \square

This last theorem can also be reformulated like: in any square matrix the number of linearly independent columns is equal to the number of independent rows.

32.2 Orthogonal matrices

Definition 32.4 A linear transformation S of the space \mathbb{R}^n is called *orthogonal*, if it is invertible and $S^{-1} = S^T$.

What does the matrix S look like? The equality $S^T S = E$ means that the scalar product of the i -th column of S by itself is 1, moreover, in case of $i \neq j$ the scalar product of the i -th and the j -th columns is zero. That tells us that each column has a unit norm, and the different columns are pairwise orthogonal. This is where the name comes from.

Example 32.5 Let S stand for the rotation of the vectors of the plane around the origin by the angle φ in positive direction. As we have seen, the matrix of this transformation is given by

$$S = \begin{bmatrix} \cos \varphi & -\sin \varphi \\ \sin \varphi & \cos \varphi \end{bmatrix}$$

This matrix is easily seen to be orthogonal, since both columns are of unit norm, and the scalar product of the two columns is zero. Therefore,

$$S^{-1} = \begin{bmatrix} \cos \varphi & \sin \varphi \\ -\sin \varphi & \cos \varphi \end{bmatrix}$$

which is precisely the matrix of the rotation by the angle $-\varphi$. Verify this directly by evaluating the product $S^T S$.

Theorem 32.6 *An orthogonal transformation keeps the length of the vectors.*

Proof. Indeed, if S is an orthogonal transformation, then

$$\|Sx\|^2 = \langle Sx, Sx \rangle = \langle S^T Sx, x \rangle = \langle x, x \rangle = \|x\|^2$$

for every vector x . □

Theorem 32.7 *The absolute value of any eigenvalue of an orthogonal matrix is 1.*

Proof. If λ is an eigenvalue of the orthogonal transformation S , and $v \neq 0$ is an associated eigenvector, then

$$|\lambda|^2 \cdot \|v\|^2 = \langle \lambda v, \lambda v \rangle = \langle Sv, Sv \rangle = \langle S^T Sv, v \rangle = \|v\|^2$$

that implies $|\lambda|^2 = 1$. □

32.3 Symmetric matrices

Definition 32.8 A linear transformation A of the space \mathbb{R}^n is called *symmetric* if $A = A^T$.

Obviously, in this case the matrix A is symmetric with respect to the diagonal (that explains the name), i.e. $a_{ij} = a_{ji}$ for all indices i and j . As a special case of Theorem 32.2 we have the following statement.

Theorem 32.9 *If A is symmetric, then*

$$\ker A = (\operatorname{im} A)^\perp$$

What can we say about the eigenvalues and eigenvectors of an $n \times n$ symmetric transformation A ? Set

$$P(\lambda) = \det(A - \lambda E)$$

which is the *characteristic polynomial* of A (of degree n).

Theorem 32.10 *The polynomial P has a real root. Consequently, the symmetric transformation A has a real eigenvalue and an associated eigenvector.*

The proof of this theorem is fairly complicated (technically goes far beyond this course), so we skip it. It can also be proven that (with multiplicity) P has exactly n real roots.

Theorem 32.11 *If A is a symmetric transformation, then the space \mathbb{R}^n possesses an orthonormal basis that consists of eigenvectors of A .*

Proof. The proof is carried out by induction on the indices $1 \leq k \leq n$. Our preceding theorem claims that if $k = 1$, then the transformation A has a real eigenvalue λ_1 and an associated eigenvector v_1 with $\|v_1\| = 1$.

Let us assume that we have found $k-1$ orthonormal eigenvectors v_1, \dots, v_{k-1} and let M denote their spanned subspace. Then A maps the subspace M^\perp into itself (invariant), because if $x \in M^\perp$ is any given vector, then

$$\langle v_i, Ax \rangle = \langle Av_i, x \rangle = \langle \lambda v_i, x \rangle = 0$$

for every $i = 1, \dots, k-1$, and hence $Ax \in M^\perp$. If we now consider the symmetric transformation A restricted onto the subspace M^\perp , then (in view of the previous theorem) we can again find a real eigenvalue λ_k and an associated eigenvector v_k with $\|v_k\| = 1$.

As a result of our construction, this v_k is orthogonal to the eigenvectors v_1, \dots, v_{k-1} , therefore the vectors v_1, \dots, v_k form an orthonormal system. \square

32.4 Spectral theorem of symmetric matrices

Now we exhibit how we can find the diagonal form of a symmetric matrix.

Take an $n \times n$ symmetric matrix A . Theorem 32.11 claims that we can find an orthonormal basis of the space \mathbb{R}^n that consists of eigenvectors of A . Let S denote the matrix whose columns are these eigenvectors.

Then obviously, S is an orthogonal matrix, which means $S^{-1} = S^T$. Summing up, we can reformulate Theorem 29.10 specifically for symmetric matrices in the following way.

Theorem 32.12 (Spectral theorem for symmetric matrices) *Let A be an $n \times n$ symmetric matrix, and consider an orthonormal basis of the space that consists of eigenvectors of A . Let S denote the matrix of the eigenvectors. Then S is orthogonal, and the matrix of A in the basis of the eigenvectors is*

$$\hat{A} = S^T A S$$

where \hat{A} is the following diagonal matrix:

$$\hat{A} = \begin{bmatrix} \lambda_1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & \lambda_n \end{bmatrix}$$

where the diagonal elements are the corresponding eigenvalues, respectively.

Take a look at how we can construct S for a given symmetric matrix A .

The simple situation is when the matrix A admits n different eigenvalues. Then the corresponding eigenvectors are automatically orthogonal. In this case we construct the matrix S by simply inserting the eigenvectors with unit norm into the columns of S .

Example 32.13 Consider the following symmetric matrix A :

$$A = \begin{bmatrix} 2 & 0 & 2 \\ 0 & -2 & 0 \\ 2 & 0 & 5 \end{bmatrix}$$

Find the eigenvalues first! The characteristic polynomial of A is given by:

$$P(\lambda) = (\lambda^2 - 7\lambda + 6)(-2 - \lambda)$$

whose roots are $\lambda_1 = 1$, $\lambda_2 = -2$ and $\lambda_3 = 6$. The corresponding eigenvectors for different eigenvalues λ can be obtained by finding nonzero solutions of the homogeneous system $(A - \lambda E)x = 0$. These solutions are, for instance

$$v_1 = \begin{bmatrix} 2 \\ 0 \\ -1 \end{bmatrix} \quad v_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad v_3 = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}$$

These eigenvectors are automatically orthogonal, and they should be converted into unit norm. Finally, the orthogonal matrix S and the diagonal matrix \hat{A} will look like:

$$S = \begin{bmatrix} \frac{2}{\sqrt{5}} & 0 & \frac{1}{\sqrt{5}} \\ 0 & 1 & 0 \\ -\frac{1}{\sqrt{5}} & 0 & \frac{2}{\sqrt{5}} \end{bmatrix} \quad \text{and} \quad \hat{A} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -2 & 0 \\ 0 & 0 & 6 \end{bmatrix}$$

Then we have

$$\hat{A} = S^T A S$$

and also verify this identity by performing the indicated multiplications.

The situation is slightly more complicated, when we have more than one linearly independent eigenvector associated with an eigenvalue. Since they are

not necessarily orthogonal, we want to use the Gram-Schmidt-procedure to make them orthogonal. This situation is illustrated in the following example.

Example 32.14 Modify the previous example this way:

$$A = \begin{bmatrix} 2 & 0 & 2 \\ 0 & 6 & 0 \\ 2 & 0 & 5 \end{bmatrix}$$

then the characteristic polynomial is

$$P(\lambda) = (6 - \lambda)(\lambda^2 - 7\lambda + 6)$$

whose roots are $\lambda_1 = 1$ and $\lambda_2 = 6$, and the latter with multiplicity 2. For the eigenvalue λ_1 it is convenient to choose the eigenvector v_1 in the preceding example. For the eigenvalue $\lambda_2 = 6$ the degree of freedom of the homogeneous system $(A - 6E)x = 0$ is 2, thus, we have 2 linearly independent solutions, for instance

$$v_2 = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix} \quad \text{and} \quad v_3 = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$$

However, these two vectors are not orthogonal, so we apply the Gram-Schmidt-procedure. As a result, instead of v_3 we obtain the following eigenvector u_3 :

$$u_3 = v_3 - \frac{\langle v_3, v_2 \rangle}{\|v_2\|^2} v_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Ultimately, we come to the matrices

$$S = \begin{bmatrix} \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} & 0 \\ 0 & 0 & 1 \\ -\frac{1}{\sqrt{5}} & \frac{2}{\sqrt{5}} & 0 \end{bmatrix} \quad \text{and} \quad \hat{A} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 6 \end{bmatrix}$$

so that

$$\hat{A} = S^T A S$$

that we should verify again directly by performing the indicated multiplications.

Recitation and Exercises

1. Reading: Textbook-1, Sections 14.5 and 14.6.
2. Homework: Textbook-1, Section 14.5, Exercises 1, 2, 3, Section 14.6, Exercises 1, 2 and 3.
3. Review: "Linear Algebra Exercises"

Chapter 33

Quadratic forms

33.1 Quadratic forms

A purely quadratic function defined on the space \mathbb{R}^n (i.e. no linear or constant terms) is introduced the following way.

Definició 33.1 A function $Q : \mathbb{R}^n \rightarrow \mathbb{R}$ is called a *quadratic form* if it is given by

$$Q(x) = Q(x_1, \dots, x_n) = a_{11}x_1^2 + a_{12}x_1x_2 + a_{22}x_2^2 + a_{13}x_1x_3 + \dots + a_{nn}x_n^2$$

that is a power function where all terms are purely quadratic.

For any quadratic form Q we can always find an $n \times n$ matrix A so that for every x

$$Q(x) = \langle x, Ax \rangle \tag{33.1}$$

thus, Q is given in terms of a scalar product. This is illustrated in the example below.

Example 33.2 On the space \mathbb{R}^3 consider the quadratic form

$$Q(x) = Q(x_1, x_2, x_3) = 3x_1^2 - 4x_1x_2 + 2x_1x_3 + x_2^2 + 6x_2x_3 - 5x_3^2$$

Collect the coefficients in matrix A the following way:

$$A = \begin{bmatrix} 3 & -4 & 2 \\ 0 & 1 & 6 \\ 0 & 0 & -5 \end{bmatrix}$$

Verify that for this matrix A we really have $Q(x) = \langle x, Ax \rangle$ for every x . However, this is not the only matrix with this property. If we introduce

$$B = \begin{bmatrix} 3 & -2 & 1 \\ -2 & 1 & 3 \\ 1 & 3 & -5 \end{bmatrix}$$

then we again have $Q(x) = \langle x, Bx \rangle$ for every $x \in \mathbb{R}^3$.

We can observe that there are infinitely many matrices A with the equality (33.1), but there exists only one symmetric matrix B that satisfies this property. If any quadratic form $Q(x) = \langle x, Ax \rangle$ is given, this symmetric matrix B is determined by the equality

$$B = \frac{1}{2}(A + A^T)$$

and then we have

$$Q(x) = \langle x, Ax \rangle = \langle x, Bx \rangle$$

for every $x \in \mathbb{R}^n$.

33.2 Symmetric matrix of a quadratic form

We can formulate the last observation of the previous section for symmetric matrices.

Theorem 33.3 *Consider a quadratic form $Q : \mathbb{R}^n \rightarrow \mathbb{R}$. Then there exists exactly one $n \times n$ symmetric matrix B for which*

$$Q(x) = \langle x, Bx \rangle$$

for every $x \in \mathbb{R}^n$. Conversely, for every symmetric matrix B the above scalar product defines a quadratic form.

This theorem basically tells us that there is a one-to-one correspondence between quadratic forms and symmetric matrices.

Example 33.4 Consider the quadratic form

$$Q(x_1, x_2, x_3) = 2x_1^2 - 2x_1x_2 + 4x_1x_3 - x_2^2 + 8x_2x_3 + 3x_3^2$$

and find the corresponding symmetric matrix B .

Absolutely analogously to the preceding example, we bisect the coefficients of the mixed products, and then we come to the following symmetric matrix:

$$B = \begin{bmatrix} 2 & -1 & 2 \\ -1 & -1 & 4 \\ 2 & 4 & 3 \end{bmatrix}$$

Check that we really have $Q(x) = \langle x, Bx \rangle$ for every $x \in \mathbb{R}^3$.

33.3 Definite quadratic forms

For finding minimum and maximum values of multivariate functions we will need to examine the sign of quadratic forms. For this purpose we introduce the following definition.

Definition 33.5 We say that a quadratic form Q is

- positive definite, if $Q(x) > 0$ for every $x \in \mathbb{R}^n$ and $x \neq 0$,
- positive semidefinite, if $Q(x) \geq 0$ for every $x \in \mathbb{R}^n$ and there exists an $x_0 \neq 0$, with $Q(x_0) = 0$,
- negative definite, if $Q(x) < 0$ for every $x \in \mathbb{R}^n$ and $x \neq 0$.
- negative semidefinite, if $Q(x) \leq 0$ for every $x \in \mathbb{R}^n$ and there exists an $x_0 \neq 0$, with $Q(x_0) = 0$,
- indefinite, if none of the above.

Example 33.6 For instance, the quadratic form with three variables

$$Q(x_1, x_2, x_3) = 2x_1^2 - 2x_1x_2 + x_2^2 + 3x_3^2$$

is positive definite, because it can be transformed into a sum of squares:

$$Q(x_1, x_2, x_3) = x_1^2 + (x_1 - x_2)^2 + 3x_3^2,$$

and this is positive for every vector $x \neq 0$.

Similarly, the quadratic form

$$Q(x_1, x_2, x_3) = x_1^2 - 2x_1x_2 + x_2^2 + 3x_3^2 = (x_1 - x_2)^2 + 3x_3^2$$

is positive semidefinite, since on the one hand $Q(x) \geq 0$ for every vector x , on the other hand $Q(1, 1, 0) = 0$, that is we can find a nonzero vector where Q take the value zero.

Further, we can see that the quadratic form

$$Q(x_1, x_2, x_3) = 2x_1^2 - 2x_1x_2 + x_2^2 - 3x_3^2$$

is indefinite, since it takes both positive and negative values. In particular, direct substitution shows that $Q(1, 1, 0) = 1 > 0$, and $Q(0, 0, 1) = -3 < 0$.

Throughout the rest of the book we will use these concepts completely analogously for symmetric matrices associated with quadratic forms.

33.4 Completing the square

The definite property of a quadratic form is very easy to decide if it consists purely of square terms (there are no mixed products).

Example 33.7 Consider the following quadratic form on the space \mathbb{R}^4 :

$$Q(x) = 5x_1^2 + 3x_2^2 + 9x_3^2 + 2x_4^2$$

This is clearly positive definite, since the sum of the squares is positive if $x \neq 0$.

On the other hand, the quadratic form

$$R(x) = 3x_1^2 + 5x_2^2 - 2x_4^2$$

is indefinite, because its value on the vector

$$x_1 = 1 \quad x_2 = 1 \quad x_3 = 0 \quad x_4 = 0$$

is positive, while the value of R on the vector

$$x_1 = 0 \quad x_2 = 0 \quad x_3 = 0 \quad x_4 = 1$$

is negative.

Very similarly, we can check that the quadratic form

$$P(x) = 3x_2^2 + x_3^2 + x_4^2$$

(ATTENTION, x_1 is missing!) positive semidefinite. Indeed, on the one hand the sum of squares is nonnegative, on the other hand there exists a vector $x \neq 0$, namely

$$x_1 = 1 \quad \text{and} \quad x_2 = x_3 = x_4 = 0$$

wher P takes the value zero. Thus, P cannot be positive definite.

We summarize the observations of this example in the following theorem.

Theorem 33.8 *Suppose that the quadratic form Q contains only purely square terms:*

$$Q(x) = b_1x_1^2 + b_2x_2^2 + \dots + b_nx_n^2$$

Based on the signs of the coefficients, we can distinguish the following cases.

- *If for every index k we have $b_k > 0$, then Q is positive definite.*
- *If every $b_k \geq 0$, and there exists an index j with $b_j = 0$, then Q is positive semidefinite.*
- *If we have both positive and negative coefficients, then Q is indefinite.*

We can formulate analogous statements for negative, resp. nonpositive coefficients. These observations lead us to examine how we can transform a quadratic form so that it contains purely square terms (completing the square in n dimension).

33.5 Definite property based on eigenvalues

Consider a quadratic form $Q : \mathbb{R}^n \rightarrow \mathbb{R}$ and let B denote the corresponding $n \times n$ symmetric matrix.

By the spectral theorem of symmetric matrices (previous chapter) the space \mathbb{R}^n has an orthonormal basis

$$v_1, \dots, v_n$$

that consists of eigenvectors of B that is

$$Bv_1 = \lambda_1 v_1 \quad \dots \quad Bv_n = \lambda_n v_n.$$

In this basis the matrix B takes the following diagonal form:

$$\hat{B} = \begin{bmatrix} \lambda_1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & \lambda_n \end{bmatrix}$$

where $\lambda_1, \dots, \lambda_n$ are the corresponding eigenvalues of B (not necessarily different). Using this diagonal matrix the quadratic form will contain purely square terms. Indeed, for any vector $y = y_1 v_1 + \dots + y_n v_n$ in the space

$$\langle y, \hat{B}y \rangle = \lambda_1 y_1^2 + \dots + \lambda_n y_n^2$$

This leads us to the following theorem.

Theorem 33.9 *Consider the quadratic form Q , and let B denote the corresponding symmetric matrix, i.e.*

$$Q(x) = \langle x, Bx \rangle$$

for every $x \in \mathbb{R}^n$. Examine the eigenvalues of B .

- If all eigenvalues are positive, then Q is positive definite.
- If all eigenvalues are nonnegative and at least one of them is zero, then Q is positive semidefinite.
- If all eigenvalues are negative, then Q is negative definite.
- If all eigenvalues are nonpositive and at least one of them is zero, then Q is negative semidefinite.
- If there are both positive and negative eigenvalues, then Q is indefinite.

Proof. We only have to prove that B and \hat{B} have the same definite property. By keeping the usual notations, if S denotes the matrix of the eigenvectors of B , then we have

$$\langle y, \hat{B}y \rangle = \langle y, S^T B S y \rangle = \langle S y, B S y \rangle = \langle x, B x \rangle$$

for every vector y in the space. Since S is invertible, the set of vectors $x = Sy$ is the whole space (i.e. the range of S is the whole space). \square

Example 33.10

Specify the definite property of the quadratic form

$$Q(x) = 2x_1^2 + 5x_1x_3 + 5x_2^2 - x_3x_1 - x_3^2.$$

The corresponding symmetric matrix B is given by:

$$B = \begin{bmatrix} 2 & 0 & 2 \\ 0 & 5 & 0 \\ 2 & 0 & -1 \end{bmatrix}$$

whose characteristic polynomial is

$$\det(B - \lambda E) = (5 - \lambda)(\lambda - 3)(\lambda + 2).$$

The roots (that is the eigenvalues of B) are

$$\lambda_1 = 5 \quad \lambda_2 = 3 \quad \text{and} \quad \lambda_3 = -2$$

We have both positive and negative eigenvalues, therefore Q is indefinite.

Recitation and Exercises

1. Reading: Textbook-1, Sections 15.8 and 15.9.
2. Homework: Textbook-1, Section 15.8, Exercises 1, 2, 3, Section 15.9, Exercises 1, 2, 3 and 4.
3. Review: "Linear Algebra Exercises"

Chapter 34

Functions with several variables

34.1 Partial derivatives

Consider a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. We can view it as $f(x) = f(x_1, \dots, x_n)$, that is the function of the n coordinates of the vector $x \in \mathbb{R}^n$.

Definition 34.1 We say that f is *partially differentiable* with respect to the k -th variable at the point x , if the function

$$F(t) = f(x + te_k)$$

is differentiable with respect to t at $t = 0$, where e_k is the k -th standard basis vector. Notation:

$$F'(0) = f'_k(x) = \frac{\partial f}{\partial x_k}(x)$$

is the *partial derivative* of f with respect to the k -th variable at the point x .

In other words, to determine the partial derivative with respect to the variable k , we regard all other variables constant and differentiate with respect to x_k only.

Example 34.2 Consider the function

$$f(x, y) = 5xe^{-2x+3y^2}$$

on the plane. Then

$$\frac{\partial f}{\partial x}(x, y) = 5e^{-2x+3y^2} - 10xe^{-2x+3y^2}$$

by the product differentiation rule, and similarly

$$\frac{\partial f}{\partial y}(x, y) = 30xye^{-2x+3y^2}$$

at every point (x, y) .

Example 34.3 When we want to find the partial derivative at a given point, it is sometimes much quicker to first substitute the fixed coordinates of the point, and then perform the differentiation. For instance, consider the function

$$f(x, y, z) = \sqrt{1 + x^2 + 3y^2 + 2z^2} \cdot (5 - x^2 - y^2) \cdot e^{-x-2y-2z}$$

on the three dimensional space, and find the partial derivative with respect to z at the point $P(2, 1, 2)$.

Of course, we could formally calculate the partial derivative function with respect to z , and substitute the coordinates of the given point P . This is immensely time consuming, and requires a lot of calculations.

A much quicker way is to first substitute $x = 2$, $y = 1$, then we get

$$f(2, 1, z) = 0$$

for every z . Thus

$$\frac{\partial f}{\partial z}(2, 1, 2) = 0.$$

Clearly, the partial derivative is zero at any other point $P(2, 1, z)$ as well.

34.2 The derivative

Definition 34.4 Assume that the partial derivatives of the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ exist at a point x (with respect to all variables). Then the vector

$$f'(x) = \left[\frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_n}(x) \right]$$

is called the *derivative* of f at x .

Sometimes this vector is also called the *gradient* of f at x .

Example 34.5 For example for the function f on the three dimensional space:

$$f(x, y, z) = 2xy\sqrt{x^2 + y^2 + z^2}$$

at the point $P(2, 1, 2)$ we have

$$f'(2, 1, 2) = \left[\frac{26}{3}, \frac{40}{3}, \frac{8}{3} \right]$$

Verify this by first calculating the vector $f'(x, y, z)$, and then substituting the coordinates of the point $P(2, 1, 2)$.

Example 34.6 Let B be an $n \times n$ symmetric matrix and consider the quadratic form:

$$Q(x) = \langle x, Bx \rangle = \sum_{i=1}^n \sum_{j=1}^n b_{ij} x_i x_j$$

where b_{ij} is the element of the i -th row and j -th column of B . Find the partial derivative of Q with respect to x_k . In this case we have

$$\frac{\partial Q}{\partial x_k}(x) = \sum_{i=1}^n b_{ik} x_i + \sum_{j=1}^n b_{kj} x_j$$

since all those terms have zero derivatives that do not contain x_k . Making use of the symmetry of B (namely $b_{ij} = b_{ji}$ for all indices), this can be rewritten like

$$\frac{\partial Q}{\partial x_k}(x) = 2 \sum_{j=1}^n b_{kj} x_j$$

for every $k = 1, \dots, n$. On the right-hand side we exactly have the k -th coordinate of the product vector $2Bx$. Therefore, the derivative of the quadratic form Q is given by

$$Q'(x) = 2Bx$$

for every x . (Please observe that this result completely complies with the derivative of a quadratic function of one variable.)

34.3 Chain-rule

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function with continuous partial derivatives with respect to all variables, and let g_1, \dots, g_n be differentiable functions on the real line. Consider the composition

$$F(t) = f(g_1(t), \dots, g_n(t))$$

for $t \in \mathbb{R}$. For simplicity, introduce the notation

$$g(t) = \begin{bmatrix} g_1(t) \\ \vdots \\ g_n(t) \end{bmatrix}$$

then we can write

$$F = f \circ g$$

on the real line. Quite analogously to the elementary chain-rule (see Theorem 4.7) we can prove the following statement.

Theorem 34.7 (Chain-rule) *Under the above conditions the composition function F is differentiable, in particular*

$$F'(t) = \langle f'(g(t)), g'(t) \rangle = \sum_{k=1}^n \frac{\partial f}{\partial x_k}(g(t)) g'_k(t)$$

at every point $t \in \mathbb{R}$.

Example 34.8 Consider the following function on the plane:

$$f(x, y) = x^2 - xy - 2y^2$$

and put

$$x = g_1(t) = \cos t \quad \text{moreover} \quad y = g_2(t) = \sin t$$

By applying the Chain-rule, the derivative of the composition $F = f \circ g$ is given by

$$F'(t) = -\frac{\partial f}{\partial x}(g(t)) \sin t + \frac{\partial f}{\partial y}(g(t)) \cos t = -4 \sin t \cos t + \sin^2 t - \cos^2 t.$$

Verify this by a direct substitution of g_1 and g_2 and by performing the differentiation. We come to the same result.

Example 34.9 Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function so that all partial derivatives exist and they are continuous functions. Take a vector $v \in \mathbb{R}^n$ and consider the function

$$g(t) = x + tv$$

where x is a fixed vector. Find the derivative of $F = f \circ g$.

Obviously $g'(t) = v$, and the Chain-rule tells us that

$$F'(t) = \langle f'(x + tv), v \rangle.$$

In particular, at $t = 0$ we have

$$F'(0) = \langle f'(x), v \rangle = \sum_{k=1}^n \frac{\partial f}{\partial x_k}(x) v_k$$

where the scalars v_k are the coordinates of the vector v .

34.4 Second order partial derivatives

Whenever the partial derivative function of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with respect to x_i is partially differentiable with respect to x_j at a given point x then we can consider the second order partial derivatives of f :

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(x) \text{ or } f''_{ij}(x) \quad \text{and in case } i = j: \quad \frac{\partial^2 f}{\partial x_i^2}(x) \text{ or } f''_{ii}(x)$$

The first is called mixed, the latter is called pure second order partial derivative.

Example 34.10 For instance, in the case of the function

$$f(x, y) = x^2 - 3x^2y^2 + 2y^3$$

the second order partial derivatives are

$$\frac{\partial^2 f}{\partial x \partial y}(x, y) = -12xy \quad \text{and} \quad \frac{\partial^2 f}{\partial y^2}(x) = -6x^2 + 12y$$

for each x and y .

Example 34.11 Reconsider the function F defined in Example 34.9 and find its second derivative $F''(x + tv)$.

Since we have

$$F'(t) = \langle f'(x + tv), v \rangle = \sum_{i=1}^n \frac{\partial f}{\partial x_i}(x + tv)v_i$$

then by carrying out the differentiation again, by the Chain-rule we get

$$F''(t) = \sum_{j=1}^n \sum_{i=1}^n \frac{\partial^2 f}{\partial x_j \partial x_i}(x + tv)v_i v_j$$

in particular, for $t = 0$ we have

$$F''(0) = \sum_{j=1}^n \sum_{i=1}^n \frac{\partial^2 f}{\partial x_j \partial x_i}(x)v_i v_j$$

This is precisely a quadratic form of the variable v , whose matrix is

$$A = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(x) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(x) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(x) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(x) & \frac{\partial^2 f}{\partial x_2^2}(x) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(x) & \frac{\partial^2 f}{\partial x_n \partial x_2}(x) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(x) \end{bmatrix}$$

By using this matrix, the above second derivative can be given in the form:

$$F''(0) = \langle v, Av \rangle$$

for every vector v .

Definition 34.12 The above $n \times n$ matrix A is called the *second derivative* of f at x . Its notation:

$$f''(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(x) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(x) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(x) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(x) & \frac{\partial^2 f}{\partial x_2^2}(x) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(x) & \frac{\partial^2 f}{\partial x_n \partial x_2}(x) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(x) \end{bmatrix}$$

Sometime the matrix $f''(x)$ is also called the *Hesse-matrix* of f at x .

34.5 Young's theorem

Example 34.13 As it is easy to see in the case of the function $f(x, y) = 2x^3 + 5x^2y^3 - \ln(xy^2)$ we have

$$\begin{aligned} \frac{\partial f}{\partial x}(x, y) &= 6x^2 + 10xy^3 - 1/x \\ \frac{\partial f}{\partial y}(x, y) &= 15x^2y^2 - 1/y^2 \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2}(x, y) &= 12x + 10y^3 - 1/x^2 \\ \frac{\partial^2 f}{\partial y^2}(x, y) &= 30x^2y - 2/y^3 \\ \frac{\partial^2 f}{\partial x \partial y}(x, y) &= \frac{\partial^2 f}{\partial y \partial x}(x, y) = 30xy^2 \end{aligned}$$

In the example above, we see that the mixed second order partial derivatives of f coincide. Our next theorem formulates that this is not a coincidence, it is always true under relatively general conditions.

Theorem 34.14 (Young) *If the second order partial derivatives of the function f with n variables exist, and they are continuous, then the Hesse-matrix $f''(x)$ is symmetric that is*

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(x) = \frac{\partial^2 f}{\partial x_j \partial x_i}(x)$$

for all indices $i, j = 1, 2, \dots, n$.

Proof. Clearly, it is enough to prove for two variables, and consider a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ that fulfills the conditions of the theorem. Take a fixed vector $v \in \mathbb{R}$ and introduce the functions

$$F(t) = f(t, y + v) - f(t, y), \quad G(t) = f(x + v, t) - f(x, t)$$

By our assumptions they are twice differentiable in a neighborhood of x and y respectively, and we have the identity:

$$F(x + v) - F(x) = G(y + v) - G(y). \quad (34.1)$$

By the Mean value theorem, there exists a number $0 < t < 1$ with

$$F(x + v) - F(x) = F'(x + tv)v,$$

and hence, in view of the definition of F we get

$$\begin{aligned} F(x + v) - F(x) &= (f'_1(x + tv, y + v) - f'_1(x + tv, y))v \\ &= (D_{12}f(x + tv, y)v + o(v))v. \end{aligned}$$

Exploiting the continuity of the second derivative, we obtain

$$\lim_{v \rightarrow 0} \frac{F(x + v) - F(x)}{v^2} = \frac{\partial^2 f}{\partial x \partial y}(x, y).$$

A completely similar argument gives us

$$\lim_{v \rightarrow 0} \frac{G(y + v) - G(y)}{v^2} = \frac{\partial^2 f}{\partial y \partial x}(x, y).$$

The identity (34.1) immediately implies

$$\frac{\partial^2 f}{\partial x \partial y}(x, y) = \frac{\partial^2 f}{\partial y \partial x}(x, y),$$

and we conclude that the second derivative is a symmetric matrix. \square

Example 34.15 Consider the function $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ defined by

$$f(x, y, z) = 2x^2y + xyz - y^2z^2$$

Check that the second derivative at an arbitrarily given point (x, y, z)

$$f''(x, y, z) = \begin{bmatrix} 4y & 4x + z & y \\ 4x + z & -2z^2 & x - 4yz \\ y & x - 4yz & -2y^2 \end{bmatrix}$$

is a symmetric matrix.

Example 34.16 Find the second derivative of the quadratic form

$$Q(x) = \langle x, Bx \rangle$$

where B is a given $n \times n$ symmetric matrix.

On the one hand we have $Q'(x) = 2Bx$, on the other hand this expression is linear, therefore the (symmetric) Hesse-matrix is

$$Q''(x) = 2B$$

at every point x .

Recitation and Exercises

1. Reading: Textbook-1, Sections 15.3, 15.4, 15.5, 15.6, 16.1 and 16.2
2. Homework: Textbook-1, Section 15.5, Exercises 1, 2, 3, 4, 5, 6, Section 16.1, Exercises 1, 2, 3, 4, 5 and 6.
3. Review: "Linear Algebra Exercises"

Chapter 35

Local extrema

35.1 Local extrema

Definition 35.1 The *unit ball* with center at the origin in the space \mathbb{R}^n is defined by

$$B = \{x \in \mathbb{R}^n : \|x\| \leq 1\}$$

Quite similarly, a ball with center at $a \in \mathbb{R}^n$ and with radius $r > 0$ is given by

$$a + rB = \{x \in \mathbb{R}^n : \|x - a\| \leq r\}$$

Consider a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. A point a in the domain of f is said to be a local minimum point of f , if there exists an $\varepsilon > 0$ so that

$$f(x) \geq f(a)$$

at every point x of the domain with $x \in a + \varepsilon B$ (that is $\|x - a\| \leq \varepsilon$).

The definition of a local maximum point can be formulated analogously. We speak about global minimum (or maximum) if the inequality holds on the entire domain.

35.2 First order necessary condition

In the following we assume that all partial derivatives of the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ exist and they are continuous in a neighborhood of the point a .

Theorem 35.2 If $a \in \mathbb{R}^n$ is a local minimum point of f , then

$$\frac{\partial f}{\partial x_1}(a) = \dots = \frac{\partial f}{\partial x_n}(a) = 0.$$

Proof. Consider the basis vectors $e_k \in \mathbb{R}^n$ and set

$$F(t) = f(a + te_k).$$

On the one hand, F has a local minimum point at $t = 0$, on the other hand, F is differentiable by the Chain-rule, namely

$$F'(t) = \langle f'(a + te_k), e_k \rangle.$$

Consequently, we get

$$0 = F'(0) = \langle f'(a), e_k \rangle = \frac{\partial f}{\partial x_k}(a)$$

for all indices $k = 1, \dots, n$. □

The above theorem tells us that we can find all extreme points of a function in the solution set of a system of equations for partial derivatives. However, this is only a necessary condition (just like in the one variable case). For instance, consider the function

$$f(x, y) = x^3y^2$$

then one solution to the system $f'_1(x, y) = f'_2(x, y) = 0$ is $(x, y) = (0, 0)$. Then

$$f(0, 0) = 0$$

which is neither a minimum nor a maximum. Indeed, in any neighborhood of the origin the function f takes both positive and negative values.

We need second order (necessary and/or sufficient) conditions as well.

35.3 Second order necessary condition

In this section we assume that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable in a neighborhood of the point a .

Theorem 35.3 *Let us suppose that a is a local minimum point of f . Then the Hesse-matrix of f at a is positive semidefinite.*

Proof. Take any vector $v \in \mathbb{R}^n$ and introduce the function

$$F(t) = f(a + tv)$$

By the assumptions F is twice continuously differentiable, and if a is a local minimum point of f , then $t = 0$ is a local minimum point of F , which implies $F''(0) \geq 0$. This means

$$0 \leq F''(0) = \langle f''(a)v, v \rangle$$

Since the vector $v \in \mathbb{R}^n$ was chosen arbitrarily, we conclude that the Hesse-matrix is positive semidefinite. \square

This result does not give a sufficient condition (only necessary) for the minimum. It is enough just to think of the function

$$f(x, y) = x^5 + y^4$$

At the point $(0, 0)$ both partial derivatives are zero, and the Hesse-matrix is the zero matrix (which is positive semidefinite), but the origin is not a local minimum point.

We can formulate an analogous statement for the local maximum, in that case the Hesse-matrix is negative semidefinite.

35.4 Sufficient condition for local extrema

In this section we assume again that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function, whose second order partial derivatives exist, and they are continuous in a neighborhood of the point a .

Theorem 35.4 *Let us suppose that at the point a all first order partial derivatives of f are zero, and the Hesse-matrix $f''(a)$ is positive definite. Then a is a local minimum points of f .*

Of course, the negative definite property of $f''(a)$ means local maximum.

The proof of this theorem goes technically somewhat beyond this course, and we skip it (it is not too hard though). We note however that it might be tempting to think that our assumptions imply that the function

$$F(t) = f(a + tv)$$

has a local minimum at $t = 0$ for every vector v . Indeed, the sufficient condition for this is that $F''(0) > 0$ for every $v \neq 0$, which is equivalent to $f''(a)$ being positive definite.

However, the trouble is that the fact that F has a local minimum at $t = 0$ for every vector v does not imply the existence of the local minimum of f at a . This unfortunate phenomenon is illustrated in the following example.

Example 35.5 Consider the following function on the plane:

$$f(x, y) = \begin{cases} x^2 & \text{if } y = x^2 \text{ and } x > 0 \\ -x^2 & \text{if } y = x^2 \text{ and } x < 0 \\ x^2 + y^2 & \text{elsewhere} \end{cases}$$

Clearly, the function f does not have a local minimum at the origin, because in any neighborhood of the origin it takes both positive and negative values.

However, if any nonzero planar vector v is given, then the function

$$F(t) = f((0,0) + tv)$$

has a strict local minimum at $t = 0$. Indeed, any straight line passing through the origin has a segment containing the origin that does not intersect the parabola with equation $y = x^2$.

It is highly recommended to create a picture!

35.5 Finding the extreme values

If we want to find the minimum and maximum points of a function with n variables, we have to perform the following steps:

1. Find all partial derivatives.
2. Make them equal to zero, and solve the system of equations.
3. At every such point determine the Hesse-matrix.
4. If at a point the Hesse-matrix is positive definite, then it is a local minimum point.
5. If at a point the Hesse-matrix is negative definite, then it is a local maximum point.
6. If at a point the Hesse-matrix is indefinite, then it is not a local extremum (so-called "saddle point").
7. If at a point the Hesse-matrix is semidefinite, then further examinations are needed. (Based on the derivatives this case is undecided.)

Example 35.6 Consider the following function on the plane:

$$f(x, y) = x^4 + y^2$$

The origin is the only critical point where both partial derivatives are zero. The Hesse-matrix at the origin is

$$H = \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix}$$

which is obviously positive semidefinite. Clearly, the origin is the (global) minimum point of the function.

If we slightly modify the function like:

$$f(x, y) = -x^4 + y^2$$

then the origin is still the only critical point, and the corresponding Hesse-matrix remains identical. However, the origin is no longer an extreme point, since in any neighborhood of the origin the function takes both positive and negative values. Such a point is called a saddle point of f .

Példa 35.7 Find all extreme values of the function:

$$f(x, y, z) = (x^2 - 4y)e^{-(x+y+z^2)}$$

By making the partial derivatives equal to zero, we get the following system of equations:

$$\begin{aligned} f'_1(x, y, z) &= (2x - x^2 + 4y)e^{-(x+y+z^2)} = 0 \\ f'_2(x, y, z) &= (-4 - x^2 + 4y)e^{-(x+y+z^2)} = 0 \\ f'_3(x, y, z) &= -2z(x^2 - 4y)e^{-(x+y+z^2)} = 0 \end{aligned}$$

whose only solution is $(x, y, z) = (-2, 2, 0)$

Apply the second order condition. The Hesse-matrix is:

$$f''(-2, 2, 0) = \begin{bmatrix} 6 & 4 & 0 \\ 4 & 4 & 0 \\ 0 & 0 & 8 \end{bmatrix}$$

and the corresponding quadratic form

$$6x_1^2 + 8x_1x_2 + 4x_2^2 + 8x_3^2 = 2x_1^2 + 4(x_1 + x_2)^2 + 8x_3^2.$$

is positive definite. Consequently, the function f has a local minimum at the critical point $(-2, 2, 0)$.

35.6 The special case of two variables

Our second order sufficient condition for an extremum can be reformulated for two variables in an "easy to use" way. The idea is that for two variables the definite property can easily be verified.

Let $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ be a function that fulfills the assumptions of Theorem 35.4, and take a point $a \in \mathbb{R}^2$ where $f'(a) = 0$ (i.e. a critical point). The Hesse-matrix of f at this point is

$$f''(a) = \begin{bmatrix} f''_{11}(a) & f''_{12}(a) \\ f''_{21}(a) & f''_{22}(a) \end{bmatrix}$$

The characteristic polynomial of the Hesse-matrix is given by

$$P(\lambda) = \lambda^2 - (f''_{11}(a) + f''_{22}(a))\lambda + f''_{11}(a)f''_{22}(a) - f''_{12}(a)^2$$

where we have taken into account the symmetry of the Hesse-matrix.

As we know, this quadratic polynomial has purely real roots (we refer to Theorem 32.11). By making use of the Viéte-formula, we can summarize our observations in the statement below.

Theorem 35.8

- If $f''_{11}(a)f''_{22}(a) - f''_{12}(a)^2 > 0$, then the critical point a is a local extreme point of the function f . This extreme point is a
 - local minimum point, if $f''_{11}(a) > 0$.
 - local maximum point, if $f''_{11}(a) < 0$.
- It is not an extreme point (saddle point), if $f''_{11}(a)f''_{22}(a) - f''_{12}(a)^2 < 0$.
- Undecided, if $f''_{11}(a)f''_{22}(a) - f''_{12}(a)^2 = 0$.

In the latter circumstances we need further investigations.

Recitation and Exercises

1. Reading: Textbook-1, Chapter 17.
2. Homework: Textbook-1, Section 17.4, Exercises 1, 2, 3, 4, 5, 6, 7, 8 and 9.
3. Review: "Linear Algebra Exercises"

Chapter 36

Least squares method, regression

In this section we exhibit an approximation method that provides the mathematical background for regression. A detailed discussion of regression is given in the statistics course.

36.1 Least squares method

Let us suppose that for the outcome of an experiment we carried out n observations, and at the different points x_1, \dots, x_n we obtained the values y_1, \dots, y_n . We have the idea that a linear model can be fitted on these experimental data. In other words, we are looking for a straight line with the equation $y = mx + b$ so that

$$mx_1 + b = y_1 \quad \dots \quad mx_n + b = y_n$$

In reality the data do not necessarily match our hypotheses, therefore, most of the time such a straight line does not exist. If we consider this as a "measurement error" and we are satisfied with a "good" approximation, then we may look for a line that fits the data the best possible way. By a good approximation we mean that the expression

$$f(m, b) = \sum_{i=1}^n (mx_i + b - y_i)^2$$

is minimal. This approximation procedure is called the *least squares method*.

ATTENTION! Why do not we use simply the sum of the gaps $mx_i - y_i$? In fact, we could use the sum of the absolute values $|mx_i - y_i|$. This would be theoretically absolutely correct, but it would make our computations a lot more difficult.

36.2 Analytic solution

Consider the function

$$f(m, b) = \sum_{i=1}^n (mx_i + b - y_i)^2 \quad (36.1)$$

with given values x_1, \dots, x_n and y_1, \dots, y_n respectively, and find the values of m and b so that f is a minimum.

By making the partial derivatives equal to zero, we get

$$\begin{aligned} \frac{\partial f}{\partial m}(m, b) &= \sum_{i=1}^n 2x_i(mx_i + b - y_i) = 0 \\ \frac{\partial f}{\partial b}(m, b) &= \sum_{i=1}^n 2(mx_i + b - y_i) = 0 \end{aligned}$$

This leads us to the following system of equations:

$$\begin{aligned} \sum_{i=1}^n x_i y_i &= m \sum_{i=1}^n x_i^2 + b \sum_{i=1}^n x_i \\ \sum_{i=1}^n y_i &= m \sum_{i=1}^n x_i + bn \end{aligned} \quad (36.2)$$

and from this linear system the solutions m and b can be determined uniquely. It is obvious that we really get a minimum point, since the function f is given as a sum of complete squares.

We just note that if we calculate the second order partial derivatives, then the Hesse-matrix of f is constant (independent of m and b), in particular

$$f''(m, b) = \begin{bmatrix} \sum_{i=1}^n 2x_i^2 & \sum_{i=1}^n 2x_i \\ \sum_{i=1}^n 2x_i & 2n \end{bmatrix}$$

We can easily see that the Hesse-matrix is positive definit, since the characteristic equation

$$\det(A - \lambda E) = \lambda^2 - \lambda \left(\sum_{i=1}^n 2x_i^2 + 2n \right) + 2n \sum_{i=1}^n 2x_i^2 - \left(\sum_{i=1}^n 2x_i \right)^2 = 0$$

has only positive roots. Indeed, by exploiting the inequality between the arithmetic and quadratic means (averages) we have

$$\frac{1}{n} \sum_{i=1}^n x_i < \sqrt{\sum_{i=1}^n \frac{1}{n} x_i^2}$$

since the numbers x_i are different. This implies that the constant term of the above quadratic equation is positive.

36.3 Algebraic solution

The minimum point that we obtained in the previous section, can be found by using purely algebraic machinery as well. If we introduce the notations

$$A = \begin{bmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad z = \begin{bmatrix} m \\ b \end{bmatrix}$$

then we can easily verify that

$$Az - y = \begin{bmatrix} mx_1 + b - y_1 \\ \vdots \\ mx_n + b - y_n \end{bmatrix}$$

Therefore, the function $f(m, b)$ at (36.1) can be written in the form:

$$f(z) = \|Az - y\|^2$$

We are looking for a vector z so that the distance between the vectors Az and y is the smallest possible. In other words: we are looking for a vector in the subspace $\text{im } A$ that is closest to the vector y . Obviously, this distance is the lowest possible if and only if the vector $Az - y$ is orthogonal to the subspace $\text{im } A$.

The orthogonality means that for both basis vectors e_i of the space \mathbb{R}^2

$$\langle Az - y, Ae_i \rangle = 0.$$

By a slight modification we get

$$\langle A^T Az, e_i \rangle = \langle A^T y, e_i \rangle$$

for $i = 1, 2$, which implies

$$A^T Az = A^T y.$$

Here the matrix $A^T A$ is clearly invertible, since it is a 2×2 matrix with a rank of 2. Consequently,

$$\begin{bmatrix} m \\ b \end{bmatrix} = z = (A^T A)^{-1} A^T y.$$

By carrying out the indicated matrix operations we can easily check that we retain the solution of the linear system (36.2). For practicing, perform this calculation.

36.4 Regression

Let X and Y be random variables, where the range of X is the set $\{x_1, \dots, x_n\}$. Assume that observed conditional expectations of the variable Y at the points x_1, \dots, x_n are given by

$$y_i = E(Y|X = x_i)$$

where $P(X = x_i) \neq 0$ for all indices $i = 1, \dots, n$.

By locating the points (x_i, y_i) in the coordinate system, we may have the hypothesis that they approximately lie on the graph of a given function. If for instance, this function is a straight line with equation

$$y = mx + b$$

then we want to choose the unknown parameters m and b so that this approximation should be the best possible. This means that the expected value

$$E((mX + b - Y)^2)$$

is a minimum (the smallest possible). For this purpose consider the function

$$\begin{aligned} g(m, b) &= E((mX + b - Y)^2) \\ &= E(m^2X^2 + 2bmX + b^2 - 2mXY - 2bY + Y^2) \\ &= m^2E(X^2) + 2bmE(X) + b^2 - 2mE(XY) - 2bE(Y) + E(Y^2) \end{aligned}$$

For the partial derivatives we have the following equations

$$\begin{aligned} \frac{\partial g}{\partial m}(m, b) &= 2mE(X^2) + 2bE(X) - 2E(XY) = 0 \\ \frac{\partial g}{\partial b}(m, b) &= 2mE(X) + 2b - 2E(Y) = 0 \end{aligned}$$

This linear system has a unique solution:

$$m = \frac{E(XY) - E(X)E(Y)}{E(X^2) - E(X)^2} = \frac{Cov(X, Y)}{Var(X)}$$

and

$$b = E(Y) - \frac{Cov(X, Y)}{Var(X)}E(X)$$

It is easy to see that we really get a minimum point, since g is the expansion of a complete square.

ATTENTION!

Verify that the Hesse-matrix of g is positive definite! (And incidentally, independent of the variables m and b .)

This function $y = mx + b$ is called the *linear regression function*. In statistics other types of regression functions (for instance quadratic, or more complicated) are also used.

Recitation and Exercises

1. Reading: Textbook-1, Chapter 17 and Textbook-2, Sections 11.1, 11.2 and 11.3.
2. Homework: Textbook-2, Exercises 11.1 through 11.14.
3. Review: "Linear Algebra Exercises"