ALGEBRA

Lecture notes for MA 630/631

Rudi Weikard

updated version of December 2010 $\,$

Contents

Chapter	1. The Language of Mathematics	1
1.1.	Propositional Calculus and Laws of Inference	1
1.2.	Predicate Logic and Basic Set Theory	4
1.3.	Further Topics in Set Theory	8
1.4.	Relations, Functions and Partial Orderings	10
1.5.	The Subjects of Algebra	13
1.6.	The Number Systems	14
Chapter	2. Groups	23
2.1.	Groups	23
Chapter	3. Rings	33
3.1.	General Ring Theory	33
3.2.	Ring Homomorphisms	36
3.3.	Unique Factorization	37
3.4.	Polynomials	42
3.5.	Algebraic Geometry	49
Chapter	: 4. Fields	53
4.1.	Field Extensions	53
4.2.	Some more concepts from group theory	57
4.3.	Galois Theory	58
	Radical Extensions	59
4.5.	The Theorem of Ruffini and Abel	61
Chapter	5. Vector Spaces	63
	Fundamentals	63
5.2.	Linear Transformations	65
5.3.	Finite-dimensional vector spaces	67
5.4.	Eigenvalues and Eigenvectors	71
5.5.	Spectral Theory in Finite-dimensional Complex Vector Spaces	72
5.6.	Multilinear Algebra	77
5.7.	Normed Spaces and Inner product Spaces	79
Index		85

CHAPTER 1

The Language of Mathematics

1.1. Propositional Calculus and Laws of Inference

1.1.1 Statements. A sentence (in the grammatical sense) for which it is (in principle) possible to determine whether it is true or false will be called a *statement*. Thus "Paris is the capital of France." and "Two plus three is six." are statements. "How are you?", "Let's see.", and " $2 + x^2 = 11$." are not. "True", denoted by T, and "false", denoted by F, are called truth values.

1.1.2 Connectives. Statements may be combined to obtain other statements: the combination is achieved using so called conjunctions like 'and', 'or', or 'but'. In logic and mathematics one uses the following conjunctions and calls them *connectives*:

- The conjunction 'and'; symbol \wedge .
- The disjunction 'or' (in the sense of 'p or q or both.'); symbol \vee .
- The conditional 'if-then'; symbol \rightarrow .
- The biconditional 'if and only if'; symbol \leftrightarrow .

The precise definitions of these connectives are given in the truth table below by stating what the truth value of such a combined statement is given the truth values of its parts. In the following p and q denote given statements.

p	q	p and q	p or q	p or q but not both	If p then q	p if and only if q
Т	Т	Т	Т	\mathbf{F}	Т	Т
Т	F	F	Т	Т	\mathbf{F}	\mathbf{F}
\mathbf{F}	Т	F	Т	Т	Т	\mathbf{F}
F	F	F	F	\mathbf{F}	Т	Т

Just as in arithmetic the use of parentheses might be necessary in order to indicate priorities in performing various operations when more complicated statements are formed symbolically.

Instead of the phrase 'If p then q.' one uses also 'q if p.' and 'p only if q.".

Together with the statement $p \to q$ one often considers also its *converse* $q \to p$ and its *contrapositive* $-q \to -p$.

1.1.3 Negation. Another way to manipulate a statement is to form its negation. The negation of p is called 'Not p.' and is denoted by -p. The statement -p is true if p is false and false if p is true.

1.1.4 Tautologies. A combined statement is called a *tautology* if it is always true no matter what the truth values of its parts are. The most primitive examples of tautologies are $p \vee (-p)$ which is called the law of the excluded middle and $-(p \wedge -p)$ which is called the law of non-contradiction.

A statement p is said to *imply* a statement q *logically*, if $p \to q$ is a tautology. This is denoted by $p \Rightarrow q$. For instance, $p \land (p \to q) \Rightarrow q$. This can be seen from a truth table.

Two statements p and q are called *logically equivalent* if $p \leftrightarrow q$ is a tautology. We write then $p \Leftrightarrow q$. Note that a logical equivalence is a logical implication both ways. The following examples are important:

- $p \to q \Leftrightarrow q \lor -p$,
- $p \to q \Leftrightarrow -q \to -p$,
- "Either p or q but not both." $\Leftrightarrow (p \land -q) \lor (-p \land q)$.

The second example, in particular, shows that a statement and its contrapositive are logically equivalent. Note, however, that a statement and its converse are not logically equivalent.

In mathematics (as well as in many other aspects of life) one wants to draw conclusions from certain pieces of information. The information is given in terms of statements which are taken for granted. They are called *premises*. The conclusion is another statement which we will always accept provided we accept the premises. Drawing a conclusion is (by definition) allowed if the conclusion is logically implied (in the above sense) by the premises. Thus the tautologies are our rules for drawing conclusions. For example, if we know (or take for granted) the truths of the statements "If it rains in Birmingham then Vulcan will get wet." and "On May 35 it rained in Birmingham all day." then we may conclude that "On May 35 Vulcan got wet.".

1.1.5 Laws of logic. In the following we list the most important logical implications and equivalences together with their names in logic. Note, however, that the list is not complete and that any tautology may be used to draw conclusions.

Law of detachment	$p \land (p \to q) \Rightarrow q$
Modus tollens	$-q \wedge (p \rightarrow q) \Rightarrow -p$
Law of disjunctive syllogism	$-p \land (p \lor q) \Rightarrow q$
Law of hypothetical syllogism	$(p \to q) \land (q \to r) \Rightarrow p \to r$
Law of simplification	$p \wedge q \Rightarrow p$
Law of addition	$p \Rightarrow p \lor q$
Law of absurdity	$p ightarrow (q \wedge -q) \Rightarrow -p$
Law of conditionalizing	$q \Rightarrow p \rightarrow q$
Law of double negation	$p \Leftrightarrow -(-p)$
Law of contraposition	$p ightarrow q \Leftrightarrow -q ightarrow -p$
Law of ex- and importation	$(p \land q) \to r \Leftrightarrow p \to (q \to r)$
Law of the conditional	$p ightarrow q \Leftrightarrow q \lor -p$
Law of the biconditional	$p \leftrightarrow q \Leftrightarrow (p \rightarrow q) \land (q \rightarrow p) \Leftrightarrow (p \land q) \lor (-p \land -q)$
De Morgan's laws	$-(p \wedge q) \Leftrightarrow -p \vee -q \text{ and } -(p \vee q) \Leftrightarrow -p \wedge -q$
Commutative laws	$p \wedge q \Leftrightarrow q \wedge p \text{ and } p \lor q \Leftrightarrow q \lor p$
Associative laws	$(p \land q) \land r \Leftrightarrow p \land (q \land r) \text{ and } (p \lor q) \lor r \Leftrightarrow p \lor (q \lor r)$
Distributive laws	$p \wedge (q \vee r) \Leftrightarrow (p \wedge q) \vee (p \wedge r)$ and
	$p \lor (q \land r) \Leftrightarrow (p \lor q) \land (p \lor r)$

1.1.6 A scheme for proofs. We will now introduce a very formal scheme for the proof of a statement from a number of premises. While later on we will not actually employ the scheme we will always employ its spirit. A proof according to this scheme consists of a list of lines. Each line has four entries: a number to enumerate it, a statement, a list of premises on which the statement depends, and an explanation. Lines are added to the list of lines by one of the following rules:

Rule P: Add a premise. (Put the current line number in the list of premises and write *Premise* as explanation.)

- **Rule L:** Add a statement which is logically implied by a preceding statement. (Copy the list of premises from that statement and state as explanation which logical implication was used.)
- **Rule C:** Add the conjunction of previous statements. (For the new list of premises put the union of the lists of premises of the appropriate previous statements. As explanation write which lines have been combined.)
- **Rule CP:** If statement S is in the list depending on the premises $P_1, ..., P_n$ and R, then statement $R \to S$ may be added as depending on the premises $P_1, ..., P_n$.

Consider the following example involving the statements

n= "Napoleon expels Snowball from the farm.",

w = "There will be a windmill.",

f = "There will be a higher food production."

Suppose the following scenario. If Napoleon expels Snowball from the farm, then there will be no higher food production. Resources are spent on a higher food production or a windmill. We also know for a fact that no windmill was built. We want to prove that Napoleon did not expel Snowball from the farm. The proof can be done in the following way:

(1)	$n \rightarrow -f$	$\{1\}$	premise
(2)	$f \lor w$	$\{2\}$	premise
(3)	$w \lor f$	$\{2\}$	from (2) by commutative law
(4)	-w	$\{4\}$	premise
(5)	$-w \wedge (w \vee f)$	$\{2,4\}$	from (3) and (4) by Rule C
(6)	f	$\{2,4\}$	from (5) by law of disjunctive syllogism
(7)	-(-f)	$\{2,4\}$	from (6) by law of double negation
(8)	$-(-f) \wedge (n \to -f)$	$\{1,2,4\}$	from (1) and (7) by Rule C
(9)	-n	$\{1,2,4\}$	from (8) by modus tollens

Using premises 1,2, and 4 we have logically deduced that Napoleon did not expel Snowball. As another example suppose we want to deduce $s \to -m$ from the premises $-s \lor \ell$ and $m \to -\ell$. We might proceed in the following way:

(1)	s	$\{1\}$	premise
(2)	-(-s)	{1}	from (1) by law of double negation
(3)	$-s \lor \ell$	$\{3\}$	premise
(4)	$-(-s) \land (-s \lor \ell)$	$\{1,3\}$	from (2) and (3) by Rule C
(5)	ℓ	$\{1,3\}$	from (4) by law of disjunctive syllogism
(6)	$-(-\ell)$	$\{1,3\}$	from (5) by law of double negation
(7)	$m \to -\ell$	$\{7\}$	premise
(8)	$-(-\ell) \wedge (m \to -\ell)$	$\{1,3,7\}$	from (6) and (7) by Rule C
(9)	-m	$\{1,3,7\}$	from (8) by modus tollens
(10)	$s \rightarrow -m$	$\{3,7\}$	from (9) by Rule CP

1.1.7 Obtaining logical implications by proofs. Suppose a line in a proof is of the form

(n) $p \to s$ {} ...

where p can, of course, be of the form $p_1 \wedge ... \wedge p_k$. Since the statement $p \to s$ does not depend on any premises (as indicated by {} in the third column) it is a tautology and therefore we have shown that s is logically implied by p or, symbolically, $p \Rightarrow s$. For instance the law of conditionalizing may be proved using a truth table but also in the following way:

(1) q {1} premise (2) $q \lor -p$ {1} from (1) by law of addition

1. THE LANGUAGE OF MATHEMATICS

(3)	$p \rightarrow q$	$\{1\}$	from (2) by logical equivalence of $q \vee -p$ and $p \to q$
(4)	$q \to (p \to q)$	{}	from (3) by Rule CP

Thus $q \Rightarrow p \rightarrow q$.

1.1.8 Inconsistent premises and proof by contradiction. Even if conclusions are drawn in a logically sound manner conclusions need not be true. Their truth still hinges on the truth of the premises. In particular, it may be possible to conclude both a statement and its negation from some set of premises. In this case the premises are called inconsistent. This is used in the method of proof called *proof by contradiction*. If one wants to conclude the statement *s* from the premise *p* one may be able to proceed in the following way:

(1)	p	$\{1\}$	premise
(2)	-s	$\{2\}$	premise
:			
(n)	r	$\{1,2\}$	
(n+1)	-r	$\{1,2\}$	
(n+2)	$r \wedge -r$	$\{1,2\}$	from (n) and $(n+1)$ by Rule C
(n+3)	$-s \rightarrow (r \wedge -r)$	$\{1\}$	from $(n+2)$ by Rule CP
(n+4)	-(-s)	$\{1\}$	from $(n+3)$ by law absurdity
(n+5)	8	$\{1\}$	from $(n+4)$ by law of double negation

The last line says that statement s has been logically deduced from the premise p.

1.2. Predicate Logic and Basic Set Theory

1.2.1 Predicates. So far the internal structure of statements did not enter our considerations. We now observe that in a statement something (the *predicate*) is said about something else (the subject) as in the statements "London is a city." and "Rome is the capital of France." Note here that the predicate "is the capital of France" has again an internal structure. We will distinguish one-place predicates like "is a city", two-place predicates like "is the capital of" and, generally, *n*-place predicates. If *P* is such an *n*-place predicate and $x_1, ..., x_n$ are individual objects taking those *n* places, then we will denote the corresponding the statement by $P(x_1, ..., x_n)$. For example, if *P* is the three-place predicate "________ is between _______ and ____" then P(5, 3, 9) stands for the statement "5 is between 3 and 9."

Note also that, if P is an n-place predicate and $x_2, ..., x_n$ are individual objects, then we may still think of $P(_, x_2, ..., x_n)$ as a one-place predicate and of $P(_, x_2, ..., x_{k-1}, ..., x_n)$, $x_{k+1}, ..., x_n$ as a two-place predicate etc.

1.2.2 Variables, constants, singular sentences, and formulas. A variable is a symbol meant to represent something unspecified. For instance, in the sentence " $x^2 + 16 = 25$." x represents a number but we do not specify which one. Consequently this sentence is not a statement. A constant, by contrast, represents a fixed object, e.g., "London" and "five" are constants. A sentence of the form $P(x_1, ..., x_n)$, where P is an n-place predicate and $x_1, ..., x_n$ are variables or constants is called a singular sentence. A sentence containing variables which becomes a statement after replacing the variables by constants is called a formula.

1.2.3 Sets and elements. Variables used in formulas are typically of one certain kind or another, e.g., they might represent numbers in one case or cities in another. It is often necessary to indicate precisely what is allowed to take the place of a variable. In order to do this set theory was invented. The present treatment of set theory is not rigorous. It is just supposed to enable us to deal with sets on a intuitive level. See **1.3.9** for a famous example of what might go wrong.

4

While the terms 'set' and 'element' will not be formally defined we think of a set (or collection) as a single entity which collects a variety of other entities, the elements (or members) of the set. A little more precisely, we assume that a set determines its elements and vice versa. We will use the following phrases: "x is an element (or a member) of (the set or collection) A", "x belongs to A", or A contains x" and this relationship is denoted by $x \in A$. Otherwise, if x does not belong to A, we write $x \notin A$.

1.2.4 Basic notation. There are essentially two ways to specify a set. Firstly, one can list all the elements of a set. One uses braces which include the list to indicate that the set rather than the elements is the object under consideration, e.g., $\{a, b, c\}$ is the set containing the letters a, b, and c. It is also customary to express a set by listing its first few elements and indicating the presence of more by a number of dots when it is specified in some other way which objects are contained in the set. For instance, the set of squares of natural numbers might be denoted (somewhat imprecisely) by $\{1, 4, 9, 16, ...\}$.

Secondly, a set might collect elements which all share a certain property or certain properties. Let P be a property and let P(x) stand for "x has the property P", e.g., if P is the property "is a city" then P(London) stands for "London is a city.". Then we denote by $\{x : P(x)\}$ the set of all objects which have the property P. For instance, $\{x : x \text{ is a city.}\}$ denotes the set of all cities.

1.2.5 Quantifiers. The sentence " $x^2 + 16 = 25$." can be made into a statement by replacing the variable x by a constant (number). However, there are two more possibilities: we can consider the statements "For all x we have that $x^2 + 16 = 25$." and "There is an x such that $x^2 + 16 = 25$." the first of which is false and the second of which is true. The phrase "For all x" is called a *universal quantifier* and is denoted by $\forall x :$. The phrase "There exists x such that " is called an *existential quantifier* and is denoted by $\exists x :$. Thus the previous two statements are denoted by $\forall x : x^2 + 16 = 25$.

When we say "For all x ..." it is to be understood that we agreed first on a set of which x has to be a member. This set is sometimes called the *domain of discourse*. For instance in " $\forall x : x^2 + 16 = 25$ " the domain of discourse could be the set of all real numbers, in "All atoms have nuclei." the domain of discourse is the set of atoms. It is often useful (or even necessary) to explicitly state what the domain of discourse is. If the domain of discourse is called A the quantifiers are given as $\forall x \in A :$ and $\exists x \in A :$.

Let P be a predicate and x a variable. Then x may be replaced by any other variable in both $\exists x : Px$ and $\forall x : Px$ without changing the content of the sentence provided the new variable is not used in P. We call this procedure renaming a variable.

1.2.6 Negation of quantified statements. Let F(x) be a formula involving the variable x. Then

$$\forall x : F(x) \Leftrightarrow -(\exists x : -F(x)), \\ \exists x : F(x) \Leftrightarrow -(\forall x : -F(x)).$$

1.2.7 Scope of a quantifier, bound and free variables. The formula to which a quantifier is applied is called the *scope of the quantifier*. An occurrence of a variable is called bound if it is in a quantifier or in the scope of a quantifier using this variable. Otherwise the occurrence of the variable is called free. Consider, for instance, the formulas

$$(\exists x : Px) \land Qy \text{ and } \exists x : Px \to Qy$$

where P and Q are predicates and x and y are variables. Here the two occurrences of x are bound while the occurrences of y are free. A variable in a formula is called *bound* (*free*) if there is a bound (free) occurrence of it. By renaming bound occurrences of variables

appropriately we can achieve that a variable is never both bound and free in the same formula and this will be assumed in the sequel.

Note that a formula which does not contain any free variables is a statement.

1.2.8 Rules of inference in predicate logic. In the proof scheme in 1.1.6 each line was either a premise or a statement derived from previous lines using a tautology. We now give four more rules by which a line may be added to the scheme. In particular we will allow formulas to be added to the scheme even as premises. In the following let F(x) be a formula involving freely the constant or variable x and perhaps other constants or variables. The formula or statement F(y) is obtained from F(x) by replacing every occurrence of x by y. The domain of discourse is denoted by A.

- (1) Universal Specification (US) From $\forall x \in A : F(x)$ we may infer F(y) where y is any variable or a constant from the domain of discourse. In particular, we may infer F(x). The list of premises remains unchanged. For example, if we have a line containing "All philosophers are Greek." we may add a line "Kant is Greek." since Kant is a member of the domain of discourse, i.e., the set of philosophers.
- (2) Existential Generalization (EG) If a is a constant we may infer $\exists x \in A : \tilde{F}(x)$ from F(a) where $\tilde{F}(x)$ is obtained by replacing some (may be all) of the occurrences of a in F(a) by x. The list of premises remains unchanged. For example, suppose the domain of discourse is the set of all cities and we have a line saying "London is the capital of England and London is a city in Europe." we may add any or all of the following lines to our proof:
 - "There is a city which is the capital of England and there is a city in Europe."
 - "London is the capital of England and there is a city in Europe."
 - "There is a city which is the capital of England and London is a city in Europe."
- (3) Existential Specification (ES) From $\exists x \in A : F(x)$ we may infer F(y) where y is a certain constant not previously used (a name we give to the thing we know exists). The list of premises remains unchanged. For example, let the domain of discourse be the set of senators in the U.S. Senate. Given the statement "There is a woman from Alaska in the Senate." we may infer "Ms. Y. is in the Senate." if we have not used before the name Y for anything else.
- (4) Universal Generalization (UG) From F(y) we may infer $\forall y \in A : F(y)$ provided that each of the following three conditions is satisfied.
 - y is a variable.
 - If y is free in a premise P(y) listed in the line containing F(y), then A must be chosen such that $\forall y \in A : P(y)$ becomes true.
 - y does not appear together with any constant introduced by an application of ES in the same singular sentence which is a part of a F(y).

The typical use of universal generalization is given by the following example: let F, G, and H be predicates.

(1) $\forall x : Fx \to Gx$	$\{1\}$	premise
(2) $\forall x : Gx \to Hx$	$\{2\}$	premise
(3) $Fy \to Gy$	$\{1\}$	from (1) by US
(4) $Gy \to Hy$	$\{2\}$	from (2) by US
(5) $(Fy \to Gy) \land (Gy \to H)$	$(y) \{1,2\}$	from (3) and (4) by C
(6) $Fy \to Hy$	$\{1,2\}$	from (5) by law of hypothetical syllogism
(7) $\forall y : Fy \to Hy$	$\{1,2\}$	from (6) by UG

However, universal generalization ought not to be too universal as shown in the following examples. Obviously, generalizing on a constant is wrong.

(1) $\pi < 4$ {1} premise (2) $\forall y : y < 4$ {1} from (1)

y < 4 {1} from (1) by a false UG

The next example shows that we have to watch for free variables in premises.

(1) n is even.	$\{1\}$	premise
(2) n^2 is even.	$\{1\}$	from (1) by the laws of arithmetic
(3) $\forall n : n^2$ is even.	$\{1\}$	from (2) by a false UG

The final example emphasizes that we have to be careful after using existential specification.

UG

$(1) \ \forall y : \exists x : y < x$	$\{1\}$	premise
$(2) \exists x : y < x$	$\{1\}$	from (1) by US
$(3) \ y < X$	$\{1\}$	from (2) by ES
$(4) \ \forall y : y < X$	$\{1\}$	from (3) by a false
$(5) \exists x : \forall y : y < x$	$\{1\}$	from (4) by EG
NT / / 1 / 1 /	1	• • 1 1•

Note, though, that generalizing on x in the line

(n) x is even \land (a < b) {...} ...

may be allowed even if, say, a was introduced by ES.

1.2.9 Axioms. In a given mathematical theory some statements are taken for granted. Such statements are called *axioms*. The axioms, in fact, characterize the theory. Changing an axiom means to consider a different theory. Euclidean geometry, for instance, relies on five axioms. One of them is the axiom of parallels. After trying for centuries to infer the parallel axiom from the other axioms mathematicians of the nineteenth century developed non-Euclidean geometries in which the parallel axiom is replaced by something else.

1.2.10 Definitions. A mathematical *definition* specifies the meaning of a word or phrase leaving no ambiguity. It may be considered an abbreviation. For instance, the sentence "A prime number is a natural number larger than one such that if it divides a product of two natural numbers it divides one of the factors." defines the word prime number. Note that we can not assign a truth value to this sentence. However, after having made the definition we treat the statement "p is a prime number if and only if p is a natural number larger than one which divides a product of two natural numbers only if it divides one of the factors." as always true, i.e., as a tautology. Therefore we may use it to obtain new lines in a proof just as we use the laws of logic. For instance, we may obtain from a line containing the statement "p is a prime number." another line containing the statement "p is larger than one." by the laws of detachment and simplification.

1.2.11 Theorems. A theorem is a true statement of a mathematical theory. It is usually of the form $p \to q$, i.e., the statement $p \to q$ is inferred from the axioms which act as premises. For example the theorem "If n is even, then n^2 is divisible by four." is of this form. Sometimes, when a statement hinges only on the axioms, the theorem could simply be an atomic (containing no connectives) statement, e.g., "2 is a prime number." is a true statement assuming the validity of the axioms of arithmetic.

1.2.12 References. For further study you might want to consider a textbook. The following is a very incomplete list.

- [1] Howard Kahane, *Logic and Philosophy: A Modern Introduction*, Wadsworth Publ. Co., Belmont, California
- [2] Patrick Suppes, Introduction to Logic, van Nostrand, Princeton, 1957.
- [3] Albert E. Blumberg, Logic. A First Course, Alfred E. Knopf, New York, 1976.

1. THE LANGUAGE OF MATHEMATICS

1.3. Further Topics in Set Theory

1.3.1 Equality. Two sets A and B are defined to be equal (denoted by A = B) if they contain the same elements. For example,

 $\{1, 4, 9, 16, \ldots\} = \{x : x \text{ is the square of a natural number}\}.$

In particular, the sets $A = \{a\}$ and $B = \{a, a\}$ are equal, i.e., A = B.

1.3.2 The empty set. There is one set called the *empty set* which contains no element at all. It is denoted by {}. It is introduced in order to be able to express some things more conveniently. For instance, $\{x \in \mathbb{R} : x^2 + px + q = 0\} = \{\}$ if $p^2 - 4q < 0$.

1.3.3 Subsets. If every element of a set A is also contained in the set B then we say that A is a subset of B or that B includes A. This is denoted by $A \subset B$, i.e.,

$$A \subset B \Leftrightarrow \forall x : (x \in A \to x \in B)$$

Note that $A \subset B$ is different from $A \in B$. In the latter case A is an element of B while in the former case the elements of A are also elements of B. In particular, $A \subset B$ allows A = B. If A is a subset of B but not equal to B, i.e., when B contains an element which is not contained in A then one says that A is a *proper subset* of B. The empty set is considered to be a subset of every set.

Two sets are equal if and only if both $A \subset B$ and $B \subset A$ hold (and equality of sets is often proven by proving mutual inclusion).

To every given set A there exists a set $\mathcal{P}(A)$, the *power set* of A, which contains precisely all the subsets of A (including the empty set and the set A itself). If A has finitely many, say n, elements then $\mathcal{P}(A)$ has 2^n elements.

1.3.4 Difference of sets. If A and B are sets then we introduce their *difference* by

$$A - B = \{ x \in A : x \notin B \}.$$

If B is a subset of A then A - B is also called the *complement* of B with respect to A.

1.3.5 Unions. Let C be a collection of sets. Then we define the *union* of the sets in C to be the set which contains all elements which belong to at least one of the members of C. This set is denoted by $\bigcup_{A \in C} A$, i.e.,

$$\bigcup_{A \in C} A = \{ x : x \in A \text{ for some } A \in C \} = \{ x : (\exists A \in C : x \in A) \}.$$

In particular, $\bigcup_{A \in \{\}} A = \{\}$ and $\bigcup_{A \in \{B\}} A = B$. The union of finitely many sets is also denoted by $A_1 \cup ... \cup A_n$. In particular, the union of the sets A and B is denoted by $A \cup B$.

1.3.6 Intersections. Now let C be a nonempty collection of sets. Then we define the *intersection* of the sets in C to be the set which contains all elements that belong to every member of C. This set is denoted by $\bigcap_{A \in C} A$, i.e.,

$$\bigcap_{A \in C} A = \{ x : x \in A \text{ for every } A \in C \} = \{ x : (\forall A \in C : x \in A) \}.$$

If C is empty then $\bigcap_{A \in C} A$ is not defined. Note that $\bigcap_{A \in \{B\}} A = B$. The intersection of finitely many sets is also denoted by $A_1 \cap ... \cap A_n$. In particular, the intersection of the sets A and B is denoted by $A \cap B$.

If $A \cap B = \{\}$ then A and B are called disjoint. A collection C of sets is called a set of *pairwise disjoint* sets when $A \cap B = \{\}$ for any two distinct elements A, B of C.

1.3.7 Properties of unions, intersections, and complements. The processes of taking unions and intersections obey the following laws:

(1) Commutative laws:

$$A \cup B = B \cup A$$
 and $A \cap B = B \cap A$.

(2) Associative laws:

$$(A\cup B)\cup C=A\cup (B\cup C) \quad \text{and} \quad (A\cap B)\cap C=A\cap (B\cap C).$$

(3) Distributive laws:

$$A \cup (\bigcap_{X \in S} X) = \bigcap_{X \in S} (A \cup X),$$
$$A \cap (\bigcup_{X \in S} X) = \bigcup_{X \in S} (A \cap X).$$

Here A, B, and C denote sets and S denotes a collection of sets which, in the first case, must be nonempty. Parentheses are used to denote priorities when performing the various operations of taking unions and intersections. When an associative law holds and parentheses become superfluous they are usually omitted.

Let E be a set and suppose all sets under consideration are subsets of E. One then can simplify notation by avoiding to refer to E when taking complements with respect to E. The difference E - A is abbreviated by A^c . Taking complements (with respect to E) obeys the following rules:

$$(A^c)^c = A,$$

$$\{\}^c = E, E^c = \{\},$$

$$A \cap A^c = \{\}, A \cup A^c = E,$$

$$A \subset B \text{ if and only if } B^c \subset A^c.$$

We will prove that $A \subset (A^c)^c$ in an exemplary way using the scheme for proofs from 1.1.6.

{1} (1) $x \in A$ premise (2) $-(-(x \in A))$ {1} from (1) by law of double negation (3) $-(x \not\in A)$ from (2) by definition of \notin $\{1\}$ (4) $-(x \in A^c)$ $\{1\}$ from (3) by definition of complement {1} (5) $x \not\in A^c$ from (4) by definition of $\not\in$ from (5) by definition of complement (6) $x \in (A^c)^c$ {1} $x\in A\to x\in (A^c)^c$ (7)from (6) by rule CP {} (8) $\forall x : x \in A \to x \in (A^c)^c \{\}$ from (7) by UG (9) $A \subset (A^c)^c$ {} from (8) by definition of subset

1.3.8 De Morgan's laws. Let E be a set and suppose that S is a nonempty collection of subsets of E. De Morgan's laws are

$$\left(\bigcup_{X\in S} X\right)^c = \bigcap_{X\in S} X^c,$$
$$\left(\bigcap_{X\in S} X\right)^c = \bigcup_{X\in S} X^c.$$

These formulas can be expressed concisely, if not precisely, as follows: The complement of a union is the intersection of the complements and the complement of an intersection is the union of the complements.

1.3.9 Axiomatic set theory. The discussion of set theory in this section can be labeled "Intuitive set theory". The presentation, however, has certain pitfalls. In particular, we have considered the set $A = \{x : P(x)\}$ when P was some property. Now let P be the property "is not an element of A" and consider whether $A \in A$. First assume that $A \in A$. Then P(A) is true, i.e., $A \notin A$ which is a contradiction. Hence assume that $A \notin A$. Then P(A) must be false, i.e., $A \in A$ again a contradiction. This paradox was first noticed by B. Russell and is known as Russell's paradox. In "Axiomatic set theory" this problem is avoided by restricting the use of the word set in a proper way (determined by the axioms). For a treatment of axiomatic set theory see, for instance:

- [1] Paul Halmos, Naive Set Theory, Springer, New York, 1973.
- [2] Patrick Suppes, Axiomatic Set Theory, van Nostrand, Princeton, 1960.
- [3] Martin Zuckerman, Sets and Transfinite Numbers, Macmillan Publishing Co., New York, 1974.

1.4. Relations, Functions and Partial Orderings

1.4.1 Ordered pairs and Cartesian products. Let A and B be two sets (not necessarily distinct) and assume $a \in A$ and $b \in B$. Then consider the *ordered pair* (a, b). The adjective 'ordered' emphasizes that (a, b) and (b, a) are different, in general. More precisely, (a, b) = (c, d) if and only if a = c and b = d. Thus (a, b) = (b, a) if and only if a = b.

The set of all ordered pairs which can be formed from A and B is called the *Cartesian* product of A and B and is denoted by $A \times B$, i.e.,

$$A \times B = \{(a, b) : a \in A \land b \in B\}.$$

One can next define the Cartesian product of three sets A, B, and C by $A \times B \times C = (A \times B) \times C$. Proceeding in this manner one can then define the Cartesian product of any finite number of sets.

1.4.2 Relations. Let X and Y be sets. A set of ordered pairs in $X \times Y$, i.e., a subset of $X \times Y$, is called a *relation*. If $(x, y) \in R$ we say that x is in relation R to y and denote this by xRy. Given a relation R we define the *domain* and the *image* of R by

$$dom(R) = \{x : \exists y : (x, y) \in R\},\$$
$$im(R) = \{y : \exists x : (x, y) \in R\}.$$

A set R of ordered pairs can always be thought of as a subset of the Cartesian product $dom(R) \times im(R)$. Conversely, every subset of the Cartesian product $X \times Y$ is a relation.

If R is a subset of $X \times Y$ we call R a relation from X to Y. If Y = X we call R a relation in X.

Let R be a relation from X to Y and A a subset of X. Then the set $R(A) = \{y \in Y : \exists x \in A : xRy\}$ is called the *image* of A under the relation R. In particular, im(R) = R(X). **1.4.3 Inverse and composite relations.** Given a relation R from X to Y the *inverse relation* R^{-1} is the relation from Y to X defined by

$$R^{-1} = \{(y, x) \in Y \times X : (x, y) \in R\}.$$

If B is a subset of Y then the image $R^{-1}(B)$ of B under the relation R^{-1} is also called the *preimage* of B under the relation R. In particular,

$$R^{-1}(B) = \{ x \in X : \exists y \in B : yR^{-1}x \} = \{ x \in X : \exists y \in B : xRy \}.$$

Given a relation R_1 from X to Y and a relation R_2 from Y to Z the *composite relation* $R_2 \circ R_1$ is the relation from X to Z defined by

$$R_2 \circ R_1 = \{ (x, z) \in X \times Z : \exists y \in Y : ((x, y) \in R_1 \land (y, z) \in R_2) \}.$$

1.4.4 Reflexivity, symmetry, and transitivity. A relation R in X, i.e., $R \subset X \times X$, is called *reflexive* if xRx for every $x \in X$. It is called *symmetric* if xRy implies yRx and it is called *transitive* if xRy and yRz imply xRz. Thus, a relation R in a set is symmetric if and only if $R = R^{-1}$ and it is transitive if and only if $R \circ R \subset R$.

Examples: Consider the set of all cities. The relation "is less than 300 miles away from" is reflexive and symmetric but not transitive. The relation "is smaller than or equal to" in the set of all real numbers is reflexive and transitive but not symmetric.

If the relation R in A is symmetric and transitive let $B = \{b \in A : (\exists a \in A : (b, a) \in R)\}$. Then R is symmetric, transitive, and reflexive in B.

1.4.5 Equivalence relations. A relation which is reflexive, symmetric and transitive is called an *equivalence relation*. If E is an equivalence relation in A and $x \in A$ then the equivalence class of x is defined to be the set

$$[x] = \{ y \in A : xEy \}.$$

x is then called a representative of [x]. The set of all such equivalence classes, i.e., the set $\{[x] : x \in A\}$, is called the set of equivalence classes of E and is denoted by A/E (pronounce A modulo E).

Example: Fix the natural number m. Consider the relation M in the set of integers defined by aMb if and only if a - b is divisible by m. Then M is an equivalence relation.

A partition of a nonempty set A is a set of pairwise disjoint nonempty subsets of A whose union is equal to A.

Theorem. If $E \subset A \times A$ is an equivalence relation in the nonempty set A then the set of all equivalence classes of E is a partition of A. Conversely, if P is a partition of A then there exists an equivalence relation E in A such that A/E = P.

1.4.6 Functions. Let f be a relation from X to Y. The relation f is called a *function* from X to Y if the following two conditions hold:

1. $\operatorname{dom}(f) = X$,

2. $\forall x \in X : \forall y_1, y_2 \in Y : (xfy_1 \land xfy_2) \rightarrow y_1 = y_2.$

In other words, for every $x \in X$ the exists a unique element of Y which is related to x by f. This element is called the image of x under f or the value which f assumes at x and is denoted by f(x). The point x is called an argument of f. The following notation for a function from X to Y is commonly used: $f : X \to Y : x \mapsto y = f(x)$. The words map, mapping, transformation, and operator are frequently used synonymously for function. By definition the set X is the domain of f. However, the set Y, which will be called the *range* of f, is to be distinguished from the image of X under f, i.e., the set f(X) = im(f) which is, in general, only a subset of Y.

If $f: X \to Y$ and $g: Y \to Z$ are functions then the composite relation $g \circ f$ is a function from X to Z. Composition of functions is an associative operation, i.e., $(f \circ g) \circ h = f \circ (g \circ h)$.

If $f : X \to Y$ is a function the inverse relation f^{-1} need not be a function. For example, if X = Y is the set of real numbers and $f = \{(x, x^2) : x \in X\}$ then f is a function with $f(x) = x^2$. On the other hand, both (4, 2) and (4, -2) are elements of f^{-1} and thus $\{y \in Y : 4f^{-1}y\}$ consists of two elements.

A function $f: X \to Y$ is called *onto* or *surjective* if $\operatorname{im}(f) = Y$. We then say that f is a function from X onto Y. A function $f: X \to Y$ is called *one-to-one* or *injective* if the relation f^{-1} is a function from $\operatorname{im}(f)$ (on)to X. Equivalently, f is injective if and only if f(x) = f(x') implies that x = x'. A function $f: X \to Y$ which is both injective and surjective is called *bijective* or a one-to-one correspondence between X and Y. If $f: X \to Y$ is a bijective function then $f^{-1}: Y \to X$ is also a bijective function.

Let f be a function from X to Y and A a subset of X. Then the function $g: A \to Y$ defined by g(x) = f(x) for all $x \in A$ is called the *restriction* of f to A and is denoted by $f|_A$. Conversely, the function f is called an *extension* of g.

A particular bijective function from a nonempty set X to itself is the identity function $id: x \mapsto x$.

1.4.7 Partial orderings. A relation R in X is called *antisymmetric* if xRy and yRx together imply that x = y. A reflexive, antisymmetric, transitive relation in X is called a *partial ordering* of X. We then say that X is partially ordered by R. Now let R be a partial ordering of X and A a subset of X. We define

- $u \in X$ is an *upper bound* of A if aRu for all $a \in A$,
- $l \in X$ is a *lower bound* of A if lRa for all $a \in A$,
- $\sup(A) \in X$ is a *least upper bound* or *supremum* of A if it is an upper bound of A and a lower bound of the set of all upper bounds of A,
- $\inf(A) \in X$ is a greatest lower bound or infimum of A if it is a lower bound of A and an upper bound of the set of all lower bounds of A,
- $\sup(A)$ is called a *maximum* if it is an element of A, similarly, if $\inf(A) \in A$ it is called a *minimum*,
- $m \in A$ is a maximal element if mRa implies m = a whenever $a \in A$, similarly, if $m \in A$ and if aRm implies m = a for all $a \in A$ then m is called a minimal element.

The supremum and infimum of a set are each unique if they exist. A maximum (or minimum) of A is always a maximal (or minimal) element of A.

Example: Consider the power set $\mathcal{P}(X)$ of a nonempty set X. Then the relation "is subset of" is a partial ordering of $\mathcal{P}(X)$. Now let $X = \{1, 2, 3\}, A = \{\{1\}, \{2\}, \{1, 2\}\},$ and $B = \{\{1, 2\}, \{2, 3\}\}$. Then $\max(A) = \{1, 2\}, \inf(A) = \{\}, \sup(B) = \{1, 2, 3\},$ and $\inf(B) = \{2\}$. However, B does not have a maximum and neither A nor B have a minimum. Both $\{1\}$ and $\{2\}$ are minimal elements of A.

1.4.8 Total orderings. A partial ordering R of X is called a *total (or linear) ordering* if xRy or yRx for any two $x, y \in X$, i.e., any two elements of X may be compared. In this case we call X totally ordered by R. If a is a maximal (minimal) element of $A \subset X$ then it is equal to $\max(A) (\min(A))$. In particular, maximal and minimal elements are unique if they exist.

Example: The set of real numbers is totally ordered by the relation "is smaller than or equal to".

Often a total ordering will be denoted by the symbol \leq . We then introduce also the symbols \geq , <, and > in the following way:

- $x \ge y$ if and only if $y \le x$,
- x < y if and only if $x \leq y$ but $x \neq y$, and
- x > y if and only if y < x.

1.4.9 Well orderings. A total ordering R of X is called a *well ordering* if every nonempty subset of X contains a minimum. In this case we call X well-ordered by R. In 1.6.4 we will show that the set of natural numbers is well-ordered by the relation "is smaller than or equal to".

1.4.10 The axiom of choice and some of its equivalents. If the two sets X and Y are not empty then X contains an element x and Y contains an element y. Therefore, the Cartesian product $X \times Y$ contains at least the element (x, y) and hence is not empty. In an axiomatic approach to set theory it is possible to generalize this result to collections of finitely many sets but not to arbitrary collections of sets without adding another axiom. This axiom is called the Axiom of Choice.

Axiom of Choice: Let C be a nonempty collection of nonempty sets. Then there exist a function $F: C \to \bigcup_{X \in C} X$, called a choice function, such that $F(X) \in X$ for every $X \in C$. In other words, it is possible to choose simultaneously an element from each member of the collection C and thus form an element of the Cartesian product of the elements of C.

The following two statements, Zorn's Lemma and the Well-ordering Theorem, can be shown to be equivalent to the Axiom of Choice. Even if the Axiom of Choice seems obvious and hardly worth mentioning, the Well-ordering Theorem is certainly not obvious.

Zorn's Lemma: Let X be a set which is partially ordered by R. If every subset which is totally ordered by R has an upper bound then X has a maximal element.

Well-ordering Theorem: For every set there exists a well-ordering.

1.5. The Subjects of Algebra

1.5.1 Binary Operations. A binary operation on a set A is a function from $A \times A$ to A. In other words, a binary operation assigns to each ordered pair of elements of A uniquely another element of A. One generally writes aFb = c instead of F((a,b)) = c. Familiar examples are $A \cup B$ and $A \cap B$ when A and B are sets and a + b and $a \cdot b$ when a and b are numbers.

Let F and G be two binary operations on a set A and a, b, and c arbitrary elements of A. Then

(1) F is called *associative* if always

$$(aFb)Fc = aF(bFc),$$

(2) F is called *commutative* if always

$$aFb = bFa$$
,

(3) F is called *left distributive over* G if always

$$aF(bGc) = (aFb)G(aFc),$$

(4) F is called *right distributive over* G if always

$$(aGb)Fc = (aFc)G(bFc),$$

and

(5) F is called *distributive over* G if it is both right and left distributive.

If F is associative parentheses as in (aFb)Fc may be omitted since no confusion can arise.

If F is commutative then the notions of left distributivity, right distributivity, and distributivity all coincide.

1.5.2 Identities and inverse elements. Let F be a binary operation on a set A. An element e is called a left (right) identity if eFa = a (aFe = a) for all $a \in A$. The element e is called an *identity* if it is both a left and a right identity. If there is a left identity e and a right identity e' then e = ee' = e'. In particular, an identity, if it exists, must be unique. However, it is instructive to consider the example where $A = \{a, b\}$ and aFa = bFa = a and bFb = aFb = b.

Let F be a binary operation on a set A. Let e be an identity or a left or right identity. An element b is called a left (right) inverse of an element a if bFa = e (aFb = e). The element b is called an *inverse element* of a (or just an inverse of a) if it is both a left and a right inverse of a.

If several binary operations are considered, reference to the operation in question must be made in statements about identities and inverses. **1.5.3 Groups, rings, fields, modules, vector spaces, and algebras.** Let G be a set and \cdot a binary operation on G. Then (G, \cdot) is called a *group* if the operation \cdot is associative, if G contains a left identity, and if every element of G possesses a left inverse. (G, \cdot) is called a commutative group if \cdot is commutative.

Let R be a set with two binary operations + and \cdot . Then $(R, +, \cdot)$ is called a *ring* if (R, +) is a commutative group and if \cdot is associative as well as left and right distributive over +. If \cdot is commutative then $(R, +, \cdot)$ is called a commutative ring.

A set F with two binary operations + and \cdot is called a *field* if (F, +) is a commutative group with identity element 0, if $(F - \{0\}, \cdot)$ is a commutative group, and if \cdot is distributive over +. The operations + and \cdot are called addition and multiplication, respectively.

Let (M, +) be a commutative group, R a commutative ring, and σ a function from $R \times M$ to M (called a scalar multiplication). Then $(M, R, +, \sigma)$ is called an R-module if the following conditions are satisfied:

- $\forall r, s \in R : \forall x \in M : (rs)x = r(sx),$
- $\forall r, s \in R : \forall x \in M : (r+s)x = rx + sx$,
- $\forall r \in R : \forall x, y \in M : r(x+y) = rx + ry.$

If R has a multiplicative identity 1, then it is also required that 1x = x for all $x \in M$. In this case the module is called a unitary module.

Let (V, +) be a commutative group, K a field, and σ a function from $K \times V$ to V (called a *scalar multiplication*). Then $(V, K, +, \sigma)$ is called a *vector space over* K if the following conditions are satisfied:

- $\forall r, s \in K : \forall x \in V : (rs)x = r(sx),$
- $\forall r, s \in K : \forall x \in V : (r+s)x = rx + sx$,
- $\forall r \in K : \forall x, y \in V : r(x+y) = rx + ry,$
- $\forall x \in V : 1x = x$ where 1 is the multiplicative identity in K.

In other words $(V, K, +, \sigma)$ is a vector space over K if it is a unitary K-module.

Finally $(A, K, +, \cdot, \sigma)$ is called an *associative algebra* if $(A, K, +, \sigma)$ is a vector space over K, if $(A, +, \cdot)$ is a ring, and if $(\alpha x)y = x(\alpha y) = \alpha(xy)$ for all $\alpha \in K$ and all $x, y \in A$.

1.6. The Number Systems

In this section we give a very quick overview over the construction of the various number systems starting from Peano's axioms. It should be remarked here that Peano's axioms may be derived from the usual axioms in (axiomatic) set theory. A certain familiarity with the real numbers is assumed. Anybody who is interested in more details may consult, for instance, the monographs listed at the end of the section.

1.6.1 The Peano axioms. There exists a set \mathbb{N} , called the set of natural numbers, and a function $s : \mathbb{N} \to \mathbb{N}$, called the *successor function*, with the following properties:

- (1) s is one-to-one, i.e., if s(n) = s(m) then n = m,
- (2) $\mathbb{N} \operatorname{im}(s)$ contains an element, called 1,
- (3) (Principle of Induction) Let X be a subset of N. If X contains 1 and if X contains s(n) whenever it contains n then $X = \mathbb{N}$.

If $n \in \mathbb{N}$, s(n) is called the *successor* of n. Note that 1 is not a successor of any natural number.

Theorem. $\mathbb{N} - im(s) = \{1\}$, i.e., every natural number except for 1 is the successor of some natural number.

Sketch of proof: This follows easily from the Principle of Induction.

The Principle of Induction is also the basis for the method of proof called Proof by Induction: One proves that a certain statement holds for 1 and that it holds for s(n)whenever it holds for n. Then one has proven that the statement holds for every natural number.

1.6.2 The Recursion Theorem. Let X be a nonempty set, f a function from X to X, and x_1 an element of X. Then there is one and only one function $u : \mathbb{N} \to X$ such that $u(1) = x_1$ and u(s(n)) = f(u(n)) for every $n \in \mathbb{N}$.

1. Existence: Let

Sketch of proof:

$$C = \{A \subset \mathbb{N} \times X : (1, x_1) \in A \land \forall (n, x) \in \mathbb{N} \times X : [(n, x) \in A \to (s(n), f(x)) \in A]\}.$$

Since $\mathbb{N} \times X \in C$ the collection C is not empty. Therefore

$$u = \left(\bigcap_{A \in C} A\right) \subset \mathbb{N} \times X$$

exists and is a relation from \mathbb{N} to X.

2. $u \in C$: This claim is proved in a straightforward manner.

3. dom $(u) = \mathbb{N}$: Let M = dom(u). Then $1 \in M$. Suppose $n \in M$, i.e., $\exists x : (n, x) \in u$. Then, since $u \in C$ we have also $(s(n), f(x)) \in u$, i.e., $s(n) \in M$. By the induction principle $M = \mathbb{N}$.

4. u is a function: Let

$$M = \{ n \in \mathbb{N} : \forall x, y \in X : ((n, x) \in u \land (n, y) \in u) \to x = y \}.$$

We first want to show that $1 \in M$. Assume $1 \notin M$, i.e., there exists $y_1 \in X$ such that $x_1 \neq y_1$ but $(1, y_1) \in u$. Then $u_0 = u - \{(1, y_1)\}$ is a proper subset of u and an element of C which is impossible. Hence 1 is indeed in M. Next suppose $n \in M$ and $s(n) \notin M$, i.e., there exists a unique x such that $(n, x) \in u$ but in addition to (s(n), f(x)) the set u contains also an element (s(n), y) where $y \neq f(x)$. Again $u_0 = u - \{(s(n), y)\}$ is a proper subset of u and an element of C which proves our assumption $s(n) \notin M$ wrong.

5. Uniqueness: Let $v : \mathbb{N} \to \mathbb{X}$ be a function such that $v(1) = x_1$ and v(s(n)) = f(v(n)). Also let $M = \{n \in \mathbb{N} : u(n) = v(n)\}$. Again the induction principle shows $M = \mathbb{N}$ and hence u = v.

The following is an important application of the Recursion Theorem. Let A be a nonempty set and F(A, A) the set of all functions from A to itself. Choosing X = F(A, A), $x_1 = g \in F(A, A)$, and $f = (h \mapsto g \circ h)$ the theorem shows the existence of a unique function $u : \mathbb{N} \to F(A, A)$ such that u(1) = g and $u(s(n)) = g \circ u(n)$. One may show by induction that $g \circ u(n) = u(n) \circ g$. It is customary to use the notation $u(n) = g^n$.

1.6.3 Ordering of the natural numbers. Define $u_1 = \text{id}$ and $u_{s(n)} = s^n$ as functions from N to itself. The functions u_n have the following properties:

1.
$$u_n \circ s = s \circ u_n = u_{s(n)}$$

2. $u_n(1) = n$. Use induction for $M = \{n \in \mathbb{N} : u_n(1) = n\}$.

3. $\operatorname{im}(u_n) \subset \operatorname{im}(u_m)$ if and only if $\operatorname{im}(u_{s(n)}) \subset \operatorname{im}(u_{s(m)})$. One of these conditionals follows since for every $j \in \mathbb{N}$ there exists an $\ell \in \mathbb{N}$ such that

$$u_{s(n)}(j) = s(u_n(j)) = s(u_m(\ell)) = u_{s(m)}(\ell).$$

The other one is proved similarly.

4. $\operatorname{im}(u_n) = \{n\} \cup \operatorname{im}(u_{s(n)})$. This follows from properties 1 and 2 and the definitions. 5. $n \notin \operatorname{im}(u_{s(n)})$. Use induction for $M = \{n \in \mathbb{N} : n \notin \operatorname{im}(u_{s(n)})\}$. 6. $\operatorname{im}(u_m) \subset \operatorname{im}(u_n)$ if and only if $m \in \operatorname{im}(u_n)$. Use induction for $M = \{k \in \mathbb{N} : u_m(k) \in \operatorname{im}(u_n)\}$.

Given two natural numbers n and m we introduce the relation 'less than or equal to', denoted by \leq : we say $n \leq m$ if and only if $im(u_m) \subset im(u_n)$.

Theorem. The relation \leq is a total ordering of \mathbb{N} .

Sketch of proof: Obviously the relation is reflexive and transitive. Assume that $n \leq m$ and $m \leq n$. This implies that $im(u_n) = im(u_m)$. Since

$$n \in im(u_n) = im(u_m) = \{m\} \cup im(u_{s(m)}) = \{m\} \cup im(u_{s(n)})$$

property 5 forces n = m, i.e., \leq is antisymmetric and hence a partial ordering of N.

Assume now that $n \neq m$ and let $M = \{j \in \mathbb{N} : n \in \operatorname{im}(u_j) \land m \in \operatorname{im}(u_j)\}$. Then $1 \in M$ but $M \neq \mathbb{N}$. Therefore there must be a $k \in \mathbb{N}$ such that n and m are both in $\operatorname{im}(u_k)$ but at least one of them, say n, is not in $\operatorname{im}(u_{s(k)})$. Since $\operatorname{im}(u_k) = \{k\} \cup \operatorname{im}(u_{s(k)})$ it follows that k = n and $m \in \operatorname{im}(u_{s(n)})$. Thus $n \leq s(n) \leq m$.

Given a total ordering on \mathbb{N} we may introduce the notation

$$\{k, ..., m\} = \{n \in \mathbb{N} : k \le n \le m\}$$

1.6.4 Induction and well ordering of \mathbb{N} **.** \mathbb{N} is actually well ordered by the relation 'less than or equal to'. In fact, one has the following theorem.

Theorem. Each of the following two statements is equivalent to the induction principle (IP):

Well ordering principle (WOP): \mathbb{N} is well-ordered with respect to the relation \leq . In particular, every nonempty subset of \mathbb{N} has a minimum (or first element).

Second induction principle (SIP): If S is a subset of \mathbb{N} , $1 \in S$, and $\{1, ..., n\} \subset S \Rightarrow s(n) \in S$ then $S = \mathbb{N}$.

Sketch of proof: (IP to WOP): Suppose T is a nonempty subset of \mathbb{N} . Let $M = \{n \in \mathbb{N} : (\forall t \in T : n \leq t)\}$, i.e., the set of all lower bounds of T. Then $1 \in M$. Since T is not empty it contains a number k. Then $s(k) \notin M$ and hence $M \neq \mathbb{N}$. By the induction principle there must be a number $m \in M$ such that $s(m) \notin M$, i.e., $m \leq t < s(m)$ for some $t \in T$. Then

$$\operatorname{im}(u_{s(m)}) \subset \operatorname{im}(u_t) \subset \operatorname{im}(u_m) = \{m\} \cup \operatorname{im}(u_{s(m)})$$

implies that the second inclusion is not proper, i.e., $t = m \in T$. This shows that $m = \max M = \min T$.

(WOP to SIP): Suppose S is a subset of N which contains 1 and that $s(n) \in S$ whenever $\{1, ..., n\} \subset S$. Let $T = \mathbb{N} - S$ and assume that $T \neq \{\}$. Then T has a minimum larger than 1, say s(k). Then $\{1, ..., k\} \cap T$ is empty, i.e., $\{1, ..., k\} \subset S$. But then $s(k) \in S \cap T = \{\}$. This is impossible and hence $T = \{\}$ and $S = \mathbb{N}$.

(SIP to IP): Assume $S \subset \mathbb{N}$, $1 \in S$, and $s(n) \in S$ whenever $n \in S$. Then, obviously, $\{1, ..., n\} \subset S$ implies $n \in S$ and hence $s(n) \in S$. Thus by the second induction principle $S = \mathbb{N}$.

1.6.5 Finite, countable, and uncountable sets. Let $\{1, ..., n\}$ denote the set of natural numbers which are less than or equal to n. One may then prove by induction

 $M = \{n \in \mathbb{N} : \forall k \in \{1, ..., n\} : (\exists \text{bijection } f : \{1, ..., n\} \rightarrow \{1, ..., k\}) \rightarrow k = n\} = \mathbb{N}.$

Note also that there is never a bijection from a nonempty set to the empty set.

Therefore it makes sense to give the following definition: a set A has n elements or has cardinality n if there exists a natural number n and a bijective function $f : A \to \{1, ..., n\}$. Also, if A is empty or if there is a natural number n such that A has n elements, then A is

16

called a *finite set*. Otherwise, A is called an *infinite* set. In particular, one shows that \mathbb{N} is infinite using a proof by contradiction.

Now let A be any infinite set. If there exists a bijective function $f : A \to \mathbb{N}$ then A is called *countably infinite* and otherwise *uncountable*. A set is called *countable* if it is finite or countably infinite.

Examples: The set of even natural numbers and the set of rational numbers are all countable (in fact, countably infinite). The set of real numbers, however, is uncountable.

1.6.6 Addition and multiplication. Using the Recursion Theorem we introduced the functions $s^n : \mathbb{N} \to \mathbb{N}$. We now define the binary operation of addition (+) on \mathbb{N} by letting

$$n+m=s^n(m).$$

One may then show that the operation + is associative and commutative. Note, in particular, that s(m) = 1 + m = m + 1. The law of cancellation holds, i.e., n + m = n + k if and only if m = k. Also $n \le n + m$.

We can now also define an addition on the set $F(\mathbb{N}, \mathbb{N})$ of all functions from \mathbb{N} to \mathbb{N} by letting $(f+g) = (k \mapsto f(k) + g(k))$. This is also an associative and commutative operation.

To define multiplication (·) we let $X = F(\mathbb{N}, \mathbb{N})$, $x_1 = id$, and $f = (X \to X : g \mapsto g + id)$. Then, by the Recursion Theorem there exists a unique function $t : \mathbb{N} \to F(\mathbb{N}, \mathbb{N})$ such that t(1) = id and t(s(n)) = t(n) + id. Using it we define

$$n \cdot m = (t(n))(m).$$

The operation \cdot is commutative and associative. Also, \cdot is distributive over +. The number 1 is an identity with respect to multiplication. The law of cancellation holds, i.e., $n \cdot m = n \cdot k$ if and only if m = k. If m > 1 then $n \cdot m > n$.

1.6.7 The whole numbers or integers. Call $(a, b), (c, d) \in \mathbb{N} \times \mathbb{N}$ equivalent if a+d = b+c. This definition introduces an equivalence relation Z in $\mathbb{N} \times \mathbb{N}$. Define

$$\mathbb{Z} = (\mathbb{N} \times \mathbb{N})/Z,$$

the set of integers or whole numbers. The equivalence class of (a, b) is denoted by a - b (pronounce, for now, a dash b).

The set \mathbb{Z} is totally ordered by the relation \leq defined by: $(a - b) \leq (c - d)$ if and only if $a + d \leq b + c$ where the last occurrence of the relation \leq refers to the ordering of natural numbers. This definition does not depend on the representatives chosen for the classes a - band c - d. One says that the relation \leq is *well-defined*.

We define addition and multiplication of integers in the following way. Let x = a - band y = c - d be integers. Then

$$x + y = (a - b) + (c - d) = (a + c) - (b + d),$$

$$x \cdot y = (a - b) \cdot (c - d) = (a \cdot c + b \cdot d) - (a \cdot d + b \cdot c).$$

Note that each of the symbols + and \cdot are used here to represent two different operations, namely operations on integers as well as operations on natural numbers.

The operations of addition and multiplication of integers are well-defined (independent from the representatives chosen), associative and commutative. Also multiplication is distributive over addition.

The integer a - b is called positive if b < a and negative if a < b. The integer a - a is called zero, i.e., a - a = 0. The integers a - a and (a + 1) - a are identities with respect to addition and multiplication, respectively. Also every integer a - b has an inverse with respect to addition (called an additive inverse) which is denoted by -(a - b) (pronounced

minus a dash b). In fact the inverse of a - b is b - a since (a - b) + (b - a) = 0. Note also that -(-(a - b)) = a - b. Hence we have the following theorem.

Theorem. $(\mathbb{Z}, +, \cdot)$ is a commutative ring.

The map $j : \mathbb{N} \to \mathbb{Z} : n \mapsto (b+n) - b$ is a bijection from the natural numbers to the positive integers. Suppose that j(n) = a - b and j(m) = c - d. Then j(n) + j(m) = (a+c) - (b+d) = (b+n+d+m) - (b+d) = j(n+m). Similarly, $j(n) \cdot j(m) = j(n \cdot m)$ and $n \leq m$ if and only if $j(n) \leq j(m)$. Since the map $n \mapsto b - (b+n)$ is also a bijection from the natural numbers to the negative integers we have that any integer may be represented by either 0, j(n), or -j(n) when n is a suitable natural number. These considerations show that \mathbb{N} may be treated as a subset of \mathbb{Z} . Hence we will drop the usage of j and just identify the numbers n and (b+n) - b. Also we will use -n instead of b - (b+n). Finally, if we introduce (as a luxury) the binary operation of subtraction on \mathbb{Z} as addition of an additive inverse we obtain that the symbol - (dash) represents subtraction. The identification of the natural numbers with the positive integers also justifies the use of the same symbol for the addition of natural numbers and integers. We will use the symbol \mathbb{N}_0 to denote the set $\mathbb{N} \cup \{0\} \subset \mathbb{Z}$.

We have equipped the set \mathbb{N} with various structures, namely a total ordering, an addition, and a multiplication. These structures are preserved when \mathbb{N} is treated as a subset of \mathbb{Z} . We say then that $(\mathbb{N}, \leq, +, \cdot)$ is embedded in $(\mathbb{Z}, \leq, +, \cdot)$.

The absolute value |a| of an integer a is defined to be equal to a or -a depending on whether a itself is nonnegative or negative.

1.6.8 Division theorem for integers. The following theorem is called the division theorem:

Theorem. Let $a, b \in \mathbb{Z}$, $a \neq 0$. Then there exists a unique pair of integers q and r (called quotient and remainder, respectively) such that $0 \leq r < |a|$ and b = aq + r.

Sketch of proof: Let

$$S = \{ n \in \mathbb{N} : (\exists q, r \in \mathbb{N}_0 : n - 1 = |a|q + r \land r < |a|) \}.$$

Choosing q = 0 and r = 0 shows that $1 \in S$. Now assume that $\{1, ..., k\} \subset S$ and consider k + 1. If k < |a| then choosing q = 0 and r = k shows $k + 1 \in S$. If $k \ge |a|$ then $k + 1 - |a| \in \{1, ..., k\}$. Hence there exist $q, r \in \mathbb{N}_0$ such that (k + 1 - |a|) - 1 = |a|q + r and r < |a|. This implies (k + 1) - 1 = |a|(q + 1) + r, i.e., $k + 1 \in S$. By the second induction principle $S = \mathbb{N}$. Therefore quotient and remainder with respect to |a| exist for every nonnegative number. If b is negative one shows existence of quotient and remainder by considering -b. To show uniqueness assume b = aq + r = aq' + r'. Then |r - r'| = |a||q' - q| and |r - r'| < |a|. Hence q = q' and r = r'.

1.6.9 Greatest common divisors. An integer $a \neq 0$ is called a divisor or a factor of an integer b if there exists an integer q such that b = aq. One also says that a divides b or that b is a multiple of a. If a is a divisor of each of the integers $b_1, ..., b_n$ it is called a common divisors of these. The numbers 1 and -1 are divisors of any integer and every nonzero integer is a divisor of zero. If a is a divisor of b then $-|b| \leq a \leq |b|$ unless b = 0. Hence the set of common divisors of $b_1, ..., b_n$ is a nonempty finite set unless $b_1 = ... = b_n = 0$. The largest member of this set (which is positive) is called the greatest common divisor of $b_1, ..., b_n$. Note that gcd(a, 0) = |a| when $a \neq 0$.

1.6.10 Euclid's algorithm. Let $a, b \in \mathbb{Z}$ and $a \neq 0$. Then there exist unique $q, r \in \mathbb{Z}$ such that $0 \leq r < |a|$ and b = aq + r. If c is a common divisor of a and b then it divides r. Also, if c is a common divisor of a and r then it divides b. Hence gcd(a, b) = gcd(a, r).

Euclid's algorithm may be used to compute the greatest common divisor of two integers a_1, b_1 . It works by recursion: Let $X = \mathbb{Z}^2$ and $x_1 = (a_1, b_1)$ where we may assume, without loss of generality, that $a_1 > 0$. Define $f : \mathbb{Z}^2 \to \mathbb{Z}^2$ by

$$f((a,b)) = \begin{cases} (r(b,a), a) & \text{if } a \neq 0\\ (0,0) & \text{if } a = 0 \end{cases}$$

where $r(b, a) \in \{0, ..., |a| - 1\}$ is the remainder of b after division by a as defined according to the division theorem. By the Recursion Theorem there is a unique function $u : \mathbb{N} \to \mathbb{Z}^2$ such that $u(1) = (a_1, b_1)$ and $u(n+1) = f(u(n)) = (a_{n+1}, b_{n+1})$.

Now note that $a_n \ge 0$ and, if $a_n > 0$, that $a_{n+1} = r(b_n, a_n) < a_n$. Let $M = \{k \in \mathbb{N} : a_k > 0\}$. Then M is a finite set and we denote $\max(M)$ by m. Therefore $u_{m+1} = (0, a_m)$. Since, by the above remark, $\gcd(a_{k+1}, b_{k+1}) = \gcd(a_k, b_k)$ if $a_k \ne 0$ we obtain

 $gcd(a_1, b_1) = gcd(a_2, b_2) = \dots = gcd(0, a_m) = a_m.$

1.6.11 The GCD identity. Given two integers a, b where $a \neq 0$ consider the set

$$S = \{ax + by : x, y \in \mathbb{Z}\} \cap \mathbb{N}.$$

S is not empty and hence contains a smallest element which will be denoted by d. Thus there exist $x_0, y_0 \in \mathbb{Z}$ such that $d = ax_0 + by_0$ and $q, r \in \mathbb{Z}$ such that a = qd + r and $0 \leq r < d$. Therefore $r = a - qd = a(1 - qx_0) + b(-qy_0)$ is in $S \cup \{0\}$. Since r cannot be in S we get that r = 0 and hence that d divides a. Similarly one shows that d divides b, i.e., d is a common divisor of a and b. Now suppose that c > d is also a common divisor of a and b. Then there exist integers n, m such that a = nc and b = mc. Hence $d/c = nx_0 + my_0$ is a positive integer strictly smaller than one which is impossible. Hence $d = \gcd(a, b)$. We have proven

Theorem. If $a, b \in \mathbb{Z}$ and $a \neq 0$ then there exist $x, y \in \mathbb{Z}$ such that gcd(a, b) = ax + by.

The numbers x and y can be computed by running Euclid's algorithm backwards.

1.6.12 Prime and irreducible numbers. An integer p for which |p| > 1 is called *prime* or a *prime number* if, whenever p divides ab, then p divides a or b.

An integer p for which |p| > 1 is called *irreducible* if p = ab implies that |a| = 1 or |b| = 1. In other words p is irreducible if its only divisors are 1, -1, p, and -p.

Theorem. An integer is prime if and only if it is irreducible.

Sketch of proof: Suppose p is prime and p = ab. Then (without loss of generality) p divides a, i.e., there exists $n \in \mathbb{Z}$ such that a = np. Hence p = npb, i.e., nb = 1 which shows that p is irreducible. Next suppose p is irreducible and that p divides ab. If p divides a nothing is to be proven and we assume therefore that p does not divide a. Since p is irreducible its only positive divisors are 1 and |p|. Since p does not divide a we get that gcd(p, a) = 1. By the GCD identity there exist integers x, y such that 1 = ax + py and hence b = abx + pby. This shows that p divides b.

1.6.13 Unique factorization theorem for integers. The following well-known theorem is also called the Fundamental Theorem of Arithmetic.

Theorem. Every integer x other than 0 and ± 1 is either an irreducible or a product of irreducibles. Moreover, this product is essentially unique in the sense that, when $x = a_1...a_n = b_1...b_m$ where $a_1,...,b_m$ are irreducibles, then n = m and the b_j may be rearranged so that $a_i = \pm b_i$ for i = 1, ..., n.

Sketch of proof: For convenience we will consider an irreducible as a product with one factor. We first proof the existence of a factorization by induction. Let $S = \{n : n + 1 \text{ is a product of irreducibles}\}$. Then $1 \in S$. Assume that $\{1, ..., n\} \subset S$. If n + 2is not irreducible then n + 2 = ab where $a, b \in \{2, ..., n + 1\}$. Hence, both a and b are products of irreducibles and hence n + 2 is a product of irreducibles, too. Uniqueness is also proven by induction: Let S be the set of all natural numbers n which satisfy the property that every product of n irreducibles is essentially unique. Suppose $a_1 = b_1...b_m$ where $a_1, b_1..., b_m$ are irreducibles. Since a_1 is prime it divides one of $b_1, ..., b_m$, say b_1 . Since b_1 is irreducible we get $a_1 = \pm b_1$ and hence $b_2...b_m = \pm 1$. This is impossible showing that m = 1, $a_1 = \pm b_1$, and thus $1 \in S$. Next assume that $n \in S$ and that $a_1...a_{n+1} = b_1...b_m$ where $a_1, ..., b_m$ are irreducibles. Now $a_2...a_{n+1} = (\pm b_2)b_3...b_m$. Since $a_2...a_{n+1}$ has n factors the induction hypothesis shows that m - 1 = n and, after a suitable rearrangement, $a_j = \pm b_j$ for j = 2, ..., n + 1.

1.6.14 The rational numbers. Now call $(a, b), (c, d) \in \mathbb{Z} \times \mathbb{N}$ equivalent if $a \cdot d = b \cdot c$. Again this definition introduces an equivalence relation Q in $\mathbb{Z} \times \mathbb{N}$. Define

$$\mathbb{Q} = (\mathbb{Z} \times \mathbb{N})/Q$$

the set of rational numbers. The equivalence class of the pair (a, b) is abbreviated by a/b (pronounce a slash b).

A total ordering (again, of course, denoted by \leq) is introduced on \mathbb{Q} in the following way: $a/b \leq c/d$ if and only if $a \cdot d \leq b \cdot c$ using the total ordering of the integers. a/b is called positive (negative) if a is a positive (negative) integer.

We define addition and multiplication of rational numbers in the following way. Let x = a/b and y = c/d be rational numbers. Then

$$\begin{aligned} x + y &= (a/b) + (c/d) = ((a \cdot d) + (b \cdot c))/(b \cdot d), \\ x \cdot y &= (a/b) \cdot (c/d) = (a \cdot c)/(b \cdot d). \end{aligned}$$

These operations are well-defined, associative and commutative. Multiplication is distributive over addition.

The number 0/1 is an additive identity while 1/1 is a multiplicative identity. Every rational number a/b has an additive inverse (-a)/b, simply denoted by -a/b, and every rational number but zero has a multiplicative inverse, denoted by $(a/b)^{-1}$. In fact, $(a/b)^{-1} = b/a$ if a is positive and $(a/b)^{-1} = (-b)/(-a)$ if a is negative. Altogether the following theorem holds.

Theorem. $(\mathbb{Q}, +, \cdot)$ is a field.

Note that $(\mathbb{Z}, \leq, +, \cdot)$ and hence $(\mathbb{N}, \leq, +, \cdot)$ are embedded in $(\mathbb{Q}, \leq, +, \cdot)$ by identifying the integer *n* with the rational number [(n, 1)] = n/1.

Another important property of rational numbers is that they are dense, i.e., between any two rational numbers there are infinitely many other rational numbers. Even though the set of rational numbers is countable.

1.6.15 The real numbers. We define real numbers by the method of Dedekind cuts. Let L be a subset of \mathbb{Q} with the following properties:

- (1) neither L nor L^c is empty, i.e., $\{L, L^c\}$ is a partition of \mathbb{Q} ,
- (2) if $x \in L$ and $y \in L^c$ then x < y,
- (3) for every $x_1 \in L$ there exists $x_2 \in L$ such that $x_1 < x_2$.

The set of all such L is denoted by \mathbb{R} and is called the set of real numbers. The name Dedekind cut reflects the fact that the partition (L, L^c) cuts the standard real number line in two.

A total ordering \leq on \mathbb{R} is defined through: $L_1 \leq L_2$ if and only if $L_1 \subset L_2$. A real number L is called positive if L contains a positive rational number and negative if L^c contains a negative rational number. The real number $\{q \in \mathbb{Q} : q < 0\}$ is called zero and is also denoted by 0. Every nonzero real number is either positive or negative but not both. Zero is neither positive nor negative.

Addition of real numbers is defined as follows:

$$L_1 + L_2 = \{ x_1 + x_2 : x_1 \in L_1, x_2 \in L_2 \}.$$

It is associative and commutative. The real number zero is an additive identity. Every real number L has an additive inverse

$$-L = \{ x \in \mathbb{Q} : (\exists y \in L^c : x + y < 0) \}.$$

If L is negative then -L is positive and vice versa.

To facilitate notation we introduce the set \mathbb{Q}_0^- of nonpositive rational numbers. Also, L^+ denotes the set of positive elements of L.

Multiplication of two positive real numbers is defined by

$$L_1 \cdot L_2 = \mathbb{Q}_0^- \cup \{x_1 x_2 : x_1 \in L_1^+, x_2 \in L_2^+\}.$$

Next one defines products of arbitrary real numbers by

$$L_1 \cdot L_2 = -((-L_1) \cdot L_2) \text{ if } L_1 < 0 \text{ and } L_2 > 0,$$

$$L_1 \cdot L_2 = -(L_1 \cdot (-L_2)) \text{ if } L_1 > 0 \text{ and } L_2 < 0,$$

$$L_1 \cdot L_2 = (-L_1) \cdot (-L_2) \text{ if } L_1 < 0 \text{ and } L_2 < 0, \text{ and }$$

$$L_1 \cdot L_2 = 0 \text{ if } L_1 = 0 \text{ or } L_2 = 0.$$

Multiplication is associative and commutative. Also it is distributive over addition. The number $1 = \{q \in \mathbb{Q} : q < 1\}$ is the multiplicative identity. Every nonzero real number L has a multiplicative inverse L^{-1} . If L > 0 then $L^{-1} = \{x \in \mathbb{Q} : (\exists y \in L^c : xy < 1)\}$. If L < 0 then $L^{-1} = -(-L)^{-1}$. Combining these facts we arrive at the following theorem.

Theorem. $(\mathbb{R}, +, \cdot)$ is a field.

Note that L does not contain a maximal element. Thus if $\sup(L)$ exists then it is in L^c , in fact it is the minimum of L^c , and hence a rational number. We now embed $(\mathbb{Q}, \leq, +, \cdot)$ into $(\mathbb{R}, \leq, +, \cdot)$ by identifying the rational number $\sup(L)$, if this exists, with the real number L. If $\sup(L)$ does not exist we call L an irrational number.

Examples: Let $L = \{q \in \mathbb{Q} : q < 3/4\}$. Then the real number L is identified with the rational number 3/4. Let $L = \{q \in \mathbb{Q} : (q < 0 \lor q^2 < 2)\}$. Then L is a real number but not a rational number since $\sup(L)$ does not exist. Of course, L is usually denoted by $\sqrt{2}$.

1.6.16 The least upper bound property of \mathbb{R} . Every subset of \mathbb{R} which has an upper bound has in fact a least upper bound. This fact accounts for the importance of the real numbers and sets them apart from the rational numbers, which do not have a least upper bound property.

Sketch of proof: Let $\{\} \neq \Sigma \subset \mathbb{R}$ and suppose Σ has an upper bound. Define $S = \bigcup_{L \in \Sigma} L$. Then $S \in \mathbb{R}$ and, for all $L \in \Sigma$ we have $L \leq S$, i.e., S is an upper bound of Σ . Now assume T < S is also an upper bound of Σ . Then there is a rational number $q \in S - T$ but also an $L \in \Sigma$ such that $q \in L$. This is impossible.

1.6.17 Roots. Denote the set of positive real numbers by \mathbb{R}^+ .

Theorem. Let n be a natural number. Then every positive real number has a unique n-th positive root, i.e., if $y \in \mathbb{R}^+$ then there is a unique $x \in \mathbb{R}^+$ such that $x^n = y$.

Sketch of proof: Consider the set $S = \{s \in \mathbb{R} : s > 0 \land s^n < y\}$. By the least upper bound property of \mathbb{R} the number $x = \sup(S)$ exists. One may show that $x^n = y$. Uniqueness follows from the fact that $x_1^n - x_2^n = (x_1 - x_2)(x_1^{n-1} + x_1^{n-2}x_2 + ... + x_2^{n-2})$.

1.6.18 The complex numbers. The set $\mathbb{R} \times \mathbb{R}$ is called the set of complex numbers. A complex number z = (x, y) is usually denoted by z = x + iy where x and y are real numbers. $x = \operatorname{Re}(z)$ is then called the real part and $y = \operatorname{Im}(z)$ the imaginary part of the complex number z = x + iy.

We define addition and multiplication:

$$(a+ib) + (c+id) = (a+c) + i(b+d),$$

 $(a+ib) \cdot (c+id) = (a \cdot c - b \cdot d) + i(b \cdot c + a \cdot d).$

The set of complex numbers is not equipped with an ordering. However, by identifying the real number a with the complex number a + i0 we embed $(\mathbb{R}, +, \cdot)$ into $(\mathbb{C}, +, \cdot)$, i.e., addition and multiplication in \mathbb{R} transfer in the natural way to \mathbb{C} .

Just as a + i0 is abbreviated by a the expression 0 + ib is abbreviated by ib and 0 + i1 by i. Since $i \cdot b = (0 + i1) \cdot (b + i0) = 0 + ib = ib$ we find that ib may be considered to be a product. (It is customary, in fact, to leave off the dot in all the products we have discussed, i.e., $a \cdot b = ab$.) Moreover, since also $i^2 = (0 + i1)(0 + i1) = -1 + i0 = -1$ the definitions of addition and multiplication follow formally from the rules of addition and multiplication of real numbers.

The numbers 0 = 0 + i0 and 1 = 1 + i0 are the additive and multiplicative identity, respectively. Every complex number a+ib has an additive inverse -(a+ib) = -a+i(-b) and every nonzero complex number has a multiplicative inverse $(a+ib)^{-1} = (a-ib)(a^2+b^2)^{-1}$.

To every complex number z = x + iy, $x, y \in \mathbb{R}$ one assigns a real number, called the absolute value of z and denoted by |z|, through $|z| = \sqrt{x^2 + y^2}$ and another complex number, called the complex conjugate of z and denoted by \overline{z} , through $\overline{z} = x - iy$. In particular, then, $z^{-1} = \overline{z}|z|^{-1/2}$.

Any two complex numbers u and v (and hence also two real numbers, two rational numbers etc.) satisfy the following inequalities, called triangle inequalities:

$$|u+v| \le |u| + |v|,$$

 $|u+v| \ge |u| - |v|.$

For $\theta \in \mathbb{R}$ let $\exp(i\theta) = \cos(\theta) + i\sin(\theta)$. From this one proves that $\exp(0) = 1$ and $\exp(i\alpha) \exp(i\beta) = \exp(i(\alpha + \beta))$. (There exists a deep relationship between the exponential function and the trigonometric functions which is studied in a course on complex analysis.)

Since \mathbb{C} is identified with $\mathbb{R} \times \mathbb{R}$ we may represent complex numbers by points in a two-dimensional plane. Every complex number has then a so-called polar representation: $z = x + iy = r \exp(i\theta)$ where r denotes the distance of the point (x, y) from the origin and θ is the oriented angle between the real axis and the line through (x, y) and the origin. In particular, $x = r \cos(\theta)$, $y = r \sin(\theta)$, and $r = \sqrt{x^2 + y^2}$.

1.6.19 References. For further study the following books are suggested:

- [1] Paul Halmos, Naive Set Theory, Springer, New York, 1973.
- [2] Steven G. Krantz, The Elements of Advanced Mathematics, CRC Press, Boca Raton, 1995.
- [3] Martin Zuckerman, Sets and Transfinite Numbers, Macmillan Publishing Co., New York, 1974.

22

CHAPTER 2

Groups

2.1. Groups

2.1.1 Groups. Let G be a set and \cdot a binary operation on G. Then (G, \cdot) is called a *group* if the operation \cdot is associative, if G contains a left identity and every element of G possesses a left inverse. More explicitly, (G, \cdot) is called a group if

- (1) $a \cdot (b \cdot c) = (a \cdot b) \cdot c$ for all $a, b, c \in G$,
- (2) there exists an element $1 \in G$ such that $1 \cdot a = a$ for all $a \in G$, and
- (3) for every $a \in G$ there exists $b \in G$ such that $b \cdot a = 1$.

The dot is usually omitted, i.e., $a \cdot b$ is simply written as ab. We will also say that G is a group under the operation \cdot or, if no confusion can arise just that G is a group.

A group is called commutative or abelian if its binary operation is commutative, i.e., if ab = ba for all $a, b \in G$.

2.1.2 Basic properties. Let G be group with left identity 1.

1. Any left inverse of a is also a right inverse of a.

Proof: Suppose ba = 1 and cb = 1. Then b = (ba)b and hence ab = (cb)(ab) = c(ba)b = cb = 1.

2. 1 is the unique identity in G.

Proof: Because of 1.5.2 we only have to show that 1 is a right identity. Now let $a \in G$ and suppose ba = ab = 1. Then a = (ab)a = a(ba) = a1.

3. Every element $a \in G$ has a unique inverse, denoted by a^{-1} . In particular, $(a^{-1})^{-1} = a$.

Proof: Suppose ba = 1 and b'a = 1. Then b' = b'(ab) = (b'a)b = b.

4. The inverse of the product ab is $(ab)^{-1} = b^{-1}a^{-1}$.

5. Cancellation: If ca = cb or ac = bc then a = b.

6. For every $a, b \in G$ the equations ax = b and ya = b have unique (possible different) solutions, i.e., there exists one and only one $x \in G$ such that ax = b and one and only one $y \in G$ such that ya = b.

2.1.3 Examples of abelian groups. A lot of familiar examples of groups are abelian groups:

1. $(\mathbb{Z}, +), (\mathbb{Q}, +), (\mathbb{R}, +), (\mathbb{C}, +),$

2. $(\mathbb{Q}_+, \cdot), (\mathbb{R}_+, \cdot), \text{ where } \mathbb{Q}_+ \text{ and } \mathbb{R}_+ \text{ denoted the positive rational and real, respectively,}$

3. $(\mathbb{Q} - \{0\}, \cdot), (\mathbb{R} - \{0\}, \cdot), (\mathbb{C} - \{0\}, \cdot),$

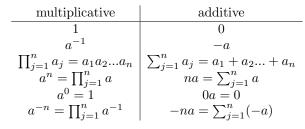
4.
$$(\{1\}, \cdot), (\{0\}, +), (\{1, -1\}, \cdot),$$

5. $(\{\exp(i\theta): \theta \in \mathbb{R}\}, \cdot).$

2.1.4 Notation. If the group G has only finitely many elements then we denote the number of its elements by |G|. This number is also called the *order* of the group. If G has infinitely many elements we say that G is of infinite order.

Nonabelian groups are usually written multiplicatively, i.e., using \cdot or juxtaposition to denote the binary operation. Abelian groups, however, are usually written additively, i.e.,

using + to denote the binary operation. This affects then also various other notational conventions.



If $n, m \in \mathbb{Z}$ and $a \in G$ then na + ma = (n + m)a and m(na) = (nm)a for a additively written group $a^n a^m = a^{n+m}$ and $(a^n)^m = a^{nm}$ for a multiplicatively written group.

2.1.5 Groups of bijections. Let X be a set and consider the set of all bijections $f : X \to X$. Composition is an associative binary operation in this set. The mapping $x \mapsto x$, called the identity map, is an identity (with respect to composition). Given any bijection $f : X \to X$ then the f^{-1} is also a bijection from X to X and it is the inverse element of f (with respect to composition).

Bijection groups need not be abelian. For example, let $X = \mathbb{R}$ and consider the bijections $x \mapsto f(x) = 2x$ and $x \mapsto g(x) = x + 1$. Then $(f \circ g)(x) = 2x + 2 \neq 2x + 1 = (g \circ f)(x)$.

2.1.6 Symmetry groups. A *rigid motion* is a bijection of the two-dimensional plane (or three-dimensional space) to itself such that distances between points are preserved. Rigid motions form a group under composition. A rigid motion which maps some subset A of the plane onto itself is called a symmetry of A. The set of symmetries of a subset A of the plane is again a group under composition. Symmetry groups need not be abelian. This is seen in the following example.

Consider the equilateral triangle two of whose vertices lie on the x-axis while the third lies on the positive y-axis. Denote the clockwise and counterclockwise rotation through 120° about the center of the triangle by ρ_{-} and ρ_{+} , respectively. Also, denote the reflection across the median from the upper vertex by ϕ_u , the reflection across the median from the lower left vertex by ϕ_l , and the reflection across the median from the lower right vertex by ϕ_r . All these operations are symmetries of the triangle. Finally, denote the identity by ι . All possible compositions of these symmetries are given by the following multiplication table where the first factor of a product is taken from the left and the second from the top:

	ι	ρ_{-}	ρ_+	ϕ_l	ϕ_r	ϕ_u
ι	ι	ρ_{-}	ρ_+	ϕ_l	ϕ_r	ϕ_u
ρ_{-}	ρ_{-}	ρ_+	$ \begin{array}{c} \rho_{+} \\ \iota \\ \rho_{-} \\ \phi_{r} \\ \phi_{u} \\ \phi_{l} \end{array} $	ϕ_r	ϕ_u	ϕ_l
ρ_+	ρ_+	ι	ρ_{-}	ϕ_u	ϕ_l	ϕ_r
ϕ_l	ϕ_l	ϕ_u	ϕ_r	ι	ρ_+	ρ_{-}
ϕ_r	ϕ_r	ϕ_l	ϕ_u	ρ_{-}	ι	ρ_+
ϕ_u	ϕ_u	ϕ_r	ϕ_l	ρ_+	ρ_{-}	ι

This table shows that compositions of the six symmetries of an equilateral triangle introduced above do not yield any new symmetries. Are there any other symmetries of the equilateral triangle? Let α be any symmetry. Then α maps vertices to vertices. Let β be the restriction of α to the set of vertices, i.e., $\beta = \alpha|_{\{1,2,3\}}$. Each such function β is a bijection from $\{1,2,3\}$ to itself. Since there are six such bijections there are at most six symmetries. Hence

$$D_3 = \{\iota, \rho_-, \rho_+, \phi_l, \phi_r, \phi_u\}$$

is the set of all symmetries of an equilateral triangle. More generally, D_n is the symmetry group of a regular *n*-gon.

2.1.7 Permutations and symmetric groups. Consider again the bijections β of 2.1.6. Any β just rearranges or permutes the "letters" 1,2, and 3. Therefore β is called a *permutation*. More generally, let $M = \{1, 2, ..., n\}$ and $S_n = \{s : M \to M : s \text{ is a bijection}\}$. Then S_n is a group under composition. (S_n, \circ) (or just S_n) is called the group of permutations of n letters (symbols) or the symmetric group of n letters. Any element of S_n is called a permutation. One usually uses juxtaposition to denote the group operation, i.e., $\alpha \circ \beta$ is abbreviated by $\alpha\beta$.

Note that S_n has n! elements since the number of bijections between two sets of n elements is n!.

A convenient notation for a permutation $\beta \in S_n$ is $\beta = \begin{pmatrix} 1 & 2 & \dots & n \\ \beta(1) & \beta(2) & \dots & \beta(n) \end{pmatrix}$. If $\beta \in S_n$ we call the set $\{k : \beta(k) \neq k\}$ the support of β denoted by $\operatorname{supp}(\beta)$. Two

If $\beta \in S_n$ we call the set $\{k : \beta(k) \neq k\}$ the support of β denoted by $\operatorname{supp}(\beta)$. Two permutations are called disjoint if their supports are disjoint. Note that $\alpha(\operatorname{supp}(\alpha)) = \operatorname{supp}(\alpha)$.

Suppose that α and β are disjoint permutations. If k is not in the support of β then $\alpha(k)$ is also not in the support of β and hence $(\alpha\beta)(k) = \alpha(k) = (\beta\alpha)(k)$. A similar argument applies when k is not in the support of α . Hence we have shown that two disjoint permutations are commutative.

2.1.8 Cycles. A permutation $\beta \in S_n$ is called a *cycle of length* k if there are distinct integers $a_1, ..., a_k \in \{1, ...n\}$ such that $\beta(a_1) = a_2, \beta(a_2) = a_3, ..., \beta(a_k) = a_1$, and such that β leaves all other n - k elements of $\{1, ..., n\}$ fixed. Such a cycle will be denoted by $(a_1, a_2, ..., a_k)$. For example the permutation

$$\left(\begin{array}{rrrrr} 1 & 2 & 3 & 4 & 5 \\ 3 & 5 & 2 & 4 & 1 \end{array}\right)$$

is a cycle of length 4 and can also be denoted by (1, 3, 2, 5).

Theorem. Every permutation, except the identity, is either a cycle or can be written as a composition (product) of disjoint cycles. This factorization is unique up to the order of the cycles.

Sketch of proof: Fix n and consider S_n . Let

 $C = \{ \alpha \in S_n : \alpha \text{ is the identity, a cycle, or a composition of disjoint cycles} \}.$

and

$$X = \{k \in \mathbb{N} : \forall \alpha \in S_n : \# \operatorname{supp}(\alpha) = k \to \alpha \in C\}.$$

Then $1 \in X$ since the support of any permutation different from the identity has at least two elements. Suppose $\{1, ..., k\} \subset X$. If $k \ge n$ then $k + 1 \in X$ since the support of any permutation has at most n elements. If k < n consider a permutation α whose support has k + 1 elements and hence is not empty. Pick $a_1 \in \text{supp}(\alpha)$ and define recursively $a_{j+1} = \alpha(a_j) = \alpha^j(a_1)$. Let $J = \{j : \exists m < j : a_m = a_j\}$. Since J is not empty it contains a smallest element $\ell + 1$. Since $a_{\ell+1} = \alpha(a_\ell) = a_1$ one obtains that $(a_1, ..., a_\ell)$ is a cycle which we denote by β . If $\ell = k + 1$ then $\alpha = \beta$ is a cycle. If $\ell \le k$ let $\gamma = \beta^{-1}\alpha$. The permutations β and γ are disjoint and $\text{supp}(\gamma)$ has no more than k elements. Hence, by induction hypothesis, γ is in C and so is α . Altogether we have $k + 1 \in X$ and the second induction principle gives now that every element of S_n except the identity is a cycle or a product of disjoint cycles. 2. GROUPS

To prove uniqueness assume that $B = \{\beta_1, ..., \beta_m\}$ and $C = \{\gamma_1, ..., \gamma_\ell\}$ each form a set of pairwise disjoint cycles and that $\alpha = \prod_{j=1}^m \beta_j = \prod_{j=1}^\ell \gamma_j$. Suppose $\beta \in B$ and choose $k \in \text{supp}(\beta)$. Then $k \in \text{supp}(\alpha)$ and hence there must be a precisely one element γ in Csuch that $k \in \text{supp}(\gamma)$. Induction shows then $\beta = \gamma$ (let $X = \{j : \beta^{j-1}(k) = \gamma^{j-1}(k)\}$) and hence $B \subset C$. Similarly one proves $C \subset B$ and this implies uniqueness.

2.1.9 Transpositions and the parity of a permutation. A cycle of length 2 is called a *transposition*. Every permutation is a product of transpositions. This follows immediately from

$$(a_1, ..., a_k) = (a_1, a_k)(a_1, a_{k-1})...(a_1, a_2).$$

Factorizations into transpositions need not be unique, for example $(1, 2) = (1, 2)^3$. However, if a permutation has a factorization into an even (odd) number of transpositions then all its factorizations into transpositions have an even (odd) number of factors, as will be shown below.

Definition. The *parity* of the identity is even, the parity of a cycle of length k is the parity of the number k - 1, and the parity of a product of disjoint cycles of lengths k_1, \ldots, k_m is the parity of the number $\sum_{j=1}^{m} (k_j - 1)$.

Theorem. If a permutation is factored into transpositions then the parity of the number of factors equals the parity of the permutation.

Sketch of proof: First note that the parity of a transposition is odd. Let $a, b, c_1, ..., c_r, d_1, ..., d_s$, where $r, s \ge 0$, be pairwise distinct and define $f = (a, c_1, ..., c_r, b, d_1, ..., d_s)$, $g = (b, c_1, ..., c_r)(a, d_1, ..., d_s)$, and $h = (a, c_1, ..., c_r)$ (here we agree that both (a) and (b) represent the identity). Since

$$\begin{aligned} (a,c_1,...,c_r,b,d_1,...,d_s)(a,b) &= (b,c_1,...,c_r)(a,d_1,...,d_s),\\ (b,c_1,...,c_r)(a,d_1,...,d_s)(a,b) &= (a,c_1,...,c_r,b,d_1,...,d_s),\\ (a,c_1,...,c_r)(a,b) &= (a,b,c_1,...,c_r) \end{aligned}$$

we have that the parity of f(a, b), g(a, b), and h(a, b) is opposite of that of f, g, and h, respectively. Hence, if $\alpha = \prod_{j=1}^{m} \beta_j$, where the β_j are pairwise disjoint cycles, then the parities of α and $\alpha(a, b)$ are opposite. Induction shows then that the parity of a product of transpositions equals the parity of the number of factors.

Note that it follows now at once that the parity of a product of permutations equals the sum of the parities of the factors. In particular, the parity of a product of two permutations is even if and only if the parities of the factors are the same (both even or both odd). Otherwise the parity of the product is odd.

2.1.10 The alternating groups. The identity of S_n has even parity. The parity of each element of S_n equals the parity of its inverse. Since the product of two even permutations is even we see that the set of all even permutations in S_n forms a group, the alternating group A_n . If $n \ge 2$, the group A_n has n!/2 elements since any transposition is a bijection from A_n to $S_n - A_n$.

2.1.11 Subgroups. Let (G, \cdot) be a group and A a subset of G. If (A, \cdot) is a group then it is called a subgroup of G.

Theorem. If (G, \cdot) is a group and A a nonempty subset of G then (A, \cdot) is a subgroup of G if and only if $a \cdot b^{-1} \in A$ for every $a, b \in A$.

Sketch of proof:

Necessity: Let (A, \cdot) be a subgroup and e the identity of A. Then $e^2 = e$. Since $e \in G$ it

26

2.1. GROUPS

has an inverse $e^{-1} \in G$. Hence $e = e^{-1}e^2 = e^{-1}e = 1$ when 1 is the identity of G. Let $a, b \in A$. Then b^{-1} , the inverse of b in G is also the inverse in A and hence an element of A. Therefore $a \cdot b^{-1}$ is in A, too.

Sufficiency: Since A is not empty there is an element $a \in A$. Then $1 = a \cdot a^{-1} \in A$. Also, for every $b \in A$ we have that $b^{-1} = 1 \cdot b^{-1} \in A$. Finally, $a \cdot b = a \cdot (b^{-1})^{-1} \in A$ whenever $a, b \in A$. Hence \cdot is an associative binary operation on A.

Examples: For any group G with identity 1 the sets $\{1\}$ and G are both subgroups of G. The former is called the trivial subgroup of G. If $k \leq n$ then S_k is a subgroup of S_n . The group of integers, $(\mathbb{Z}, +)$, is a subgroup of $(\mathbb{Q}, +)$. Rigid motions of the plane form a subgroup of the group of all bijections from \mathbb{R}^2 to \mathbb{R}^2 . The subsets $\{\iota, \rho_-, \rho_+\}$ and $\{\iota, \phi_k\}$, k = 1, 2, 3, of D_3 in Section 2.1.6 are subgroups of D_3 . The alternating group A_n is a subgroup of S_n .

If U_1 and U_2 are subgroups of a group G then $U_1 \cap U_2$ is also a subgroup of G, or, more generally, if C is any nonempty collection of subgroups of G then $\bigcap_{U \in C} U$ is also a subgroup of G.

2.1.12 Subgroups generated by subsets of a group. Let M be a nonempty subset of a group G and C the collection of all subgroups of G which include M. Then $\langle M \rangle = \bigcap_{U \in C} U$ is called the *group generated by* M.

The group generated by M is the minimum of C when the set of all subgroups of G is partially ordered by inclusion. It is therefore also called the smallest subgroup of G including M.

Theorem. Let G be a group and M a nonempty subset of G. Then

$$\langle M \rangle = \{ \prod_{j=1}^{N} a_j^{n_j} : N \in \mathbb{N}, \, a_j \in M, \, n_j \in \mathbb{Z} \text{ for } j = 1, \dots, N \}.$$

Sketch of proof: Let A be the set on the right hand side. Then $\{\} \neq M \subset A$. Suppose $a, b \in A$. Then $ab^{-1} \in A$, too. Hence A is a subgroup of G which includes M, i.e., $\langle M \rangle \subset A$. Now let $a \in A$, i.e., $a = \prod_{j=1}^{N} a_j^{n_j}$ where $a_j \in M$ and $n_j \in \mathbb{Z}$. Then $a_j \in \langle M \rangle$ and therefore $a \in \langle M \rangle$, too.

2.1.13 Centralizer and center. Let (G, \cdot) be a group and M a subset of G. Then define

 $C(M) = \{g \in G : gm = mg \text{ for every } m \in M\},\$

the centralizer of M. The centralizer C(G) of the group itself is also called the center of G. $(C(M), \cdot)$ is a subgroup of (G, \cdot) and $(C(G), \cdot)$ is abelian.

2.1.14 Order of a group element. The *order* ord *a* of an element $a \in G$ is defined to be the smallest natural number *k* such that $a^k = 1$ if this exists. Otherwise *a* is said to have infinite order (ord $a = \infty$).

2.1.15 Cyclic groups. Let $a \in G$. Then $\langle a \rangle = \{a^n : n \in \mathbb{Z}\}$ is called the *cyclic subgroup* of G generated by a. A group G such that $\langle a \rangle = G$ for some $a \in G$ is called a *cyclic group*. A cyclic group is always abelian. If a has infinite order then the elements a^n of $\langle a \rangle$ are all distinct. If a has order k then $\langle a \rangle = \{1, a, a^2, \ldots, a^{k-1}\}$.

Examples: $(\mathbb{Z}, +)$ is cyclic, in fact $\mathbb{Z} = \langle 1 \rangle$. The subgroup $\{\iota, \rho_-, \rho_+\}$ of D_3 equals $\langle \rho_+ \rangle = \langle \rho_- \rangle$ and hence is cyclic. $\{\exp(in) : n \in \mathbb{Z}\}$ and $\{\exp(2\pi i n/k) : n \in \{1, ..., k\}\}$ are cyclic groups under multiplication of order infinity and k, respectively.

Theorem. The following statements hold:

1. All subgroups of a cyclic group are cyclic.

2. For an infinite cyclic group all subgroups but the trivial one have infinitely many elements.

2. GROUPS

Sketch of proof: 1. Let S be a nontrivial subgroup of $G = \langle a \rangle$ and $A = \{n \in \mathbb{N} : a^n \in S\}$. Let $m = \min(A)$ and $b = a^m$. For any $c \in S$ we have $c = a^\ell$ where $\ell = qm + r$ with $q \in \mathbb{Z}$ and $r \in \{0, ..., m-1\}$. Then $c = b^q a^r$ and $a^r = b^{-q} c \in S$. Hence r = 0 and $S = \langle b \rangle$.

2. $\langle a \rangle$ has infinite order if and only if *a* does. Hence $b^q = a^{mq}$ are all different. **2.1.16 Subset multiplication in groups.** Let *A* and *B* be subsets of a group *G* (written multiplicatively). Then the product of *A* and *B* is defined to be the set

$$AB = \{ab : a \in A, b \in B\} \subset G.$$

Hence subset multiplication is a binary operation on the power set of G. The associative law of G induces an associative law for subset multiplication, i.e., (AB)C = A(BC). When $A = \{a\}$ we also write aB instead of AB (see the definition of cosets).

If G is written additively one uses also additive notation for subset multiplication, i.e., $A + B = \{a + b : a \in A, b \in B\}.$

2.1.17 Cosets. Let U be a subgroup of G. Define a relation R_U on $G \times G$ by defining aR_Ub if and only if $a^{-1}b \in U$. This relation is an equivalence relation. The equivalence classes have the form $[a] = aU = \{g \in G : g = au$ for some $u \in U\}$ and are called *left cosets* of U. In particular [1] = U. The set G/R_U of equivalence classes of G with respect to the relation R_U is usually denoted by G/U.

By Theorem 1.4.5 the left cosets form a partition of G. The number of left cosets is called the *index* of U in G.

Let aU and bU be two cosets then $g \mapsto f(g) = ba^{-1}g$ defines a bijective function from aU to bU. Hence, any two cosets contain either the same number of elements, are both countably infinite, or are both uncountable.

2.1.18 Lagrange's theorem. Let G be a group of finite order and U a subgroup of G. The cosets aU, $a \in G$, form a partition of G such that each one has the same number, namely |U|, of elements. In particular, the order of G is the product of the order of U and the number distinct cosets, i.e., the number of elements of G/U. Thus we have the

Theorem. The order of a finite group is divisible by the order of any of its subgroups.

Corollary. If G is a finite group, then |G| is divisible by the order of any of its elements. In particular $a^{|G|} = 1$.

The results of this section can also be found by working with the equivalence relation $aRb \Leftrightarrow ab^{-1} \in U$ having right cosets Ua as equivalence classes. In general, it is not true that aU = Ua for all $a \in G$.

2.1.19 Normal subgroups. A subgroup N of a group G is called a *normal subgroup* if the left cosets of N in G are equal to the right cosets of N in G, i.e., if gN = Ng for all $g \in G$. We write $N \triangleleft G$.

Theorem. Equivalent conditions for a subgroup N of G to be normal are given by:

1.
$$\forall g \in G : Ng \subset gN$$
,
2. $\forall g \in G : gN \subset Ng$,

3. $\forall g \in G : gNg^{-1} \subset N$.

Sketch of proof: The definition implies the first statement trivially. To prove the second from the first choose $g \in G$ and $b \in gN$. Then there exits $n \in N$ such that b = gn. Note that $Ng^{-1} \subset g^{-1}N$, in particular, there is an $m \in N$ such that $ng^{-1} = g^{-1}m$, i.e., gn = mg. Hence $b = mg \in Ng$. Since b was arbitrary $gN \subset Ng$. Since g was arbitrary the second statement follows from the first. The remaining proofs are similar.

Examples: $\{1\}$ and G are normal subgroups of G. Every subgroup of an abelian group is normal. The alternating group A_n is a normal subgroup of S_n .

28

2.1.20 Quotient groups. The following theorem is true.

Theorem. If N is a normal subgroup of G then G/N is a group under the binary operation given by the subset product.

Sketch of proof: Let $aN, bN \in G/N$. Then $(aN)(bN) = \{anbm : n, m \in N\}$. Since N is normal we find that for all $n \in N$ there exists $n' \in N$ such that nb = bn'. Hence $(aN)(bN) = \{abn'm : n', m \in N\} = (ab)N$ is a left coset of N. Choosing different representatives from aN and bN will yield the same product set so that multiplication of cosets is well defined. Hence subset multiplication is an associative binary operation on G/N. Also, N is the identity element in G/N and $a^{-1}N$ the inverse element of aN.

The group G/N of left cosets of a normal subgroup N of G is called a *quotient group* or *factor group*.

If G is finite, then we have |G/N| = |G|/|N|. If G is infinite and N is finite, then G/N is infinite.

2.1.21 Residue classes. Let m be an integer. Then the set $m\mathbb{Z} = \{mk : k \in \mathbb{Z}\}$ is a subgroup of $(\mathbb{Z}, +)$. The cosets of $m\mathbb{Z}$ are called *residue classes mod* m. For $a + (m\mathbb{Z}) = (m\mathbb{Z}) + a$ we will write a_m . The residue classes mod m are explicitly given by

$$a_m = \{a + km : k \in \mathbb{Z}\}$$

Since $(m\mathbb{Z}, +)$ is a normal subgroup of $(\mathbb{Z}, +)$ we get that $\mathbb{Z}_m = \mathbb{Z}/(m\mathbb{Z})$ is a group (under subset addition).

2.1.22 Direct products. Let G and H be groups. Then introduce a binary operation on $G \times H$ by $(g_1, h_1)(g_2, h_2) = (g_1g_2, h_1h_2)$. Under this binary operation $G \times H$ is a group called the direct product of G and H.

2.1.23 Homomorphisms. Let G and H be groups. A mapping $\eta : G \to H$ is called a *(group) homomorphism* if $\eta(ab) = \eta(a)\eta(b)$ for all $a, b \in G$.

Examples: The identity map from G to G is a homomorphism. The map $G \to \{e\} : g \mapsto e$, where e is the identity of any group, is a homomorphism. For $m \in \mathbb{N}$ the map $\mathbb{Z} \to m\mathbb{Z}$, $k \mapsto mk$, is a homomorphism.

For a homomorphism $\eta : G \to H$ we define the kernel ker η of η to be the set ker $\eta = \{g \in G : \eta(g) = 1_H\}$. Recall that the image of G under η is the set $\eta(G) = \{\eta(g) : g \in G\}$.

Let $\eta: G \to H$ and $\nu: H \to K$ be group homomorphisms. Then the following basic facts hold:

1. $\eta(1_G) = 1_H$.

2. $\eta(a^{-1}) = \eta(a)^{-1}$.

3. If $g \in G$ has finite order then the order of $\eta(g) \in H$ divides the order of g.

4. $\nu \circ \eta : G \to K$ is a homomorphism.

Theorem. For any homomorphism $\eta : G \to H$ the set $\eta(G)$ is a subgroup of H which is abelian if G is abelian. The set ker η is a normal subgroup of G. A homomorphism $\eta : G \to H$ is injective if and only if ker $\eta = \{1_G\}$.

Sketch of proof: $1_H = \eta(1_G) \in \eta(G)$. Suppose $\eta(x), \eta(y) \in \eta(G)$. Then $\eta(x)\eta(y)^{-1} = \eta(xy^{-1}) \in \eta(G)$. If G is abelian then $\eta(x)\eta(y) = \eta(xy) = \eta(yx) = \eta(y)\eta(x)$.

 $1_G \in \ker \eta$. If $x, y \in \ker \eta$ then $\eta(xy^{-1}) = \eta(x)\eta(y)^{-1} = 1_H$ and hence $xy^{-1} \in \ker \eta$. Let *a* be any element in *G* and *x* any element in $\ker \eta$. Then $axa^{-1} \in \ker \eta$, too. Hence $\ker \eta \triangleleft G$.

Suppose ker $\eta = \{1_G\}$ and $\eta(g) = \eta(g')$. Then $\eta(g^{-1}g') = 1_H$ and hence g = g'. **2.1.24 Canonical homomorphisms.** Let $N \triangleleft G$. Then the mapping $G \rightarrow G/N : a \mapsto aN$ is a homomorphism which is called the *canonical homomorphism* from G to G/N. **2.1.25 Embeddings and projections.** Let $G \times H$ be a direct product of groups. Then the functions from G or H to $G \times H$ defined by $g \mapsto (g, 1)$ and $h \mapsto (1, h)$, respectively, are injective homomorphisms. They are called embeddings.

The functions from $G \times H$ to G or H defined by $(g, h) \mapsto g$ and $(g, h) \mapsto h$, respectively, are surjective homomorphisms. They are called projections.

2.1.26 Isomorphisms. If a homomorphism $\eta : G \to H$ is bijective then η is called an *isomorphism*. In this case $\eta^{-1} : H \to G$ is also an isomorphism. Two groups G and H are called *isomorphic* if there exists an isomorphism from G to H (or from H to G). If G and H are isomorphic we write $G \cong H$.

Theorem. The relation \cong in the set of all groups is an equivalence relation.

Examples: If G is cyclic and |G| = n, then $G \cong \mathbb{Z}_n$. If G is cyclic and of infinite order then $G \cong \mathbb{Z}$. If |G| = p is a prime number, then $G \cong \mathbb{Z}_p$. If |G| = 4, then either $G \cong \mathbb{Z}_4$ or $G \cong \mathbb{Z}_2 \times \mathbb{Z}_2$. The map $x \mapsto \exp(x)$ from $(\mathbb{R}, +)$ to (\mathbb{R}_+, \cdot) is an isomorphism.

2.1.27 The fundamental isomorphism theorem for groups. Any group homomorphism on G is related to quotient groups of G.

Theorem. Let $\eta : G \to H$ be a surjective homomorphism and let ν be the canonical homomorphism from G to $G/\ker \eta$. Then there is an isomorphism $\mu : G/\ker \eta \to H$ such that $\eta = \mu \circ \nu$. In particular, $H \cong G/\ker \eta$.

Sketch of proof: Define $\mu : G/\ker\eta \to H$ by $g \ker\eta \mapsto \eta(g)$. However, this defines a function only if $\eta(g) = \eta(g')$ for all $g' \in g \ker\eta$. Hence assume that $g^{-1}g' \in \ker\eta$. Then $1_H = \eta(g^{-1}g') = \eta(g)^{-1}\eta(g')$, the required equality.

 μ is surjective since η is. Assume that $\mu(g \ker \eta) = \mu(g' \ker \eta)$, i.e., $\eta(g) = \eta(g')$, then $g^{-1}g' \in \ker \eta$ which implies injectivity. Finally, $\mu((g \ker \eta)(g' \ker \eta)) = \mu(gg' \ker \eta) = \eta(gg') = \mu(g \ker \eta)\mu(g' \ker \eta)$ which shows that μ is a homomorphism. \Box

Corollary. Every normal subgroup of G is the kernel of a homomorphism on G.

2.1.28 Automorphisms. Let G be a group. A homomorphism $\eta : G \to G$ is called an *endomorphism* in G. If, in addition, η is an isomorphism it is called an *automorphism* of G. By Aut(G) we denote the set of all automorphisms of G.

Theorem. Aut(G) is a group under composition.

Sketch of proof: $\operatorname{Aut}(G)$ is a subset of the group of all bijections of G onto itself. The identity is in $\operatorname{Aut}(G)$. Suppose $\eta, \nu \in \operatorname{Aut}(G)$. It is easily checked that $\eta \circ \nu^{-1}$ is a homomorphism. Since it is bijective it is in fact an automorphism. \Box

Example: Let G be a group, $a \in G$. Then the map from G to G defined by $x \mapsto axa^{-1}$ is an automorphism of G. An element $\eta \in \operatorname{Aut}(\mathbb{Z}_m)$ is uniquely defined by $\eta(1_m) = k_m$ where k = 1 or $k \in \{2, ..., m-1\}$ but has no factor in common with m.

2.1.29 Inner and outer automorphisms. An automorphism of G of the form $x \mapsto axa^{-1}$ for some $a \in G$ is called an *inner automorphism* of G. Every other automorphism is called an *outer automorphism* of G.

The set of all inner automorphisms of a group G is a subgroup of Aut(G).

Examples: The identity is the only inner automorphism of an abelian group. S_3 has six inner but no outer automorphisms. If p is prime then $\operatorname{Aut}(\mathbb{Z}_p)$ has p-1 elements. p-2 of these are outer automorphisms.

2.1.30 Conjugate subgroups and elements. Automorphisms of G map subgroups to subgroups.

2.1. GROUPS

Two subgroups U_1 and U_2 of G are called *conjugate subgroups* if there exists an inner automorphism σ_a such that $\sigma_a(U_1) = U_2$. Similarly, two elements x and y of G are called *conjugate* if there is an $a \in G$ with $y = axa^{-1}$.

The relation of being conjugated is an equivalence relation in G. The equivalence classes of this relation are called the *conjugacy classes* of G.

2.1.31 Free groups. Let S be a finite or infinite set of symbols: $S = \{a, b, c, ...\}$. Consider the set $F = F^* \cup \{1\}$ where $1 \notin S$ and

$$F^* = \{x_1^{n_1} \dots x_N^{n_N} : N \in \mathbb{N}, x_j \in S, n_j \in \mathbb{Z} - \{0\}, x_j \neq x_{j+1}\}.$$

Introduce on F a binary operation by 1f = f1 = f for all $f \in F$ and (recursively) by $(x^n)(x^{-n}) = 1$,

$$(x_1^{n_1}...x_N^{n_N})(y_1^{m_1}...y_M^{m_M}) = \begin{cases} x_1^{n_1}...x_N^{n_N}y_1^{m_1}...y_M^{m_M} & \text{if } x_N \neq y_1, \\ x_1^{n_1}...x_N^{n_N+m_1}y_2^{m_2}...y_M^{m_M} & \text{if } x_N = y_1 \text{ and } n_N + m_1 \neq 0, \\ (x_1^{n_1}...x_{N-1}^{n_{N-1}})(y_2^{m_2}...y_M^{m_M}) & \text{if } x_N = y_1 \text{ and } n_N + m_1 = 0. \end{cases}$$

Then F is a group called the *free group* on S.

Theorem. Let F be the free group on a set S, G an arbitrary group, and φ a function from S to G. Then φ extends in a unique way to a group homomorphism from F to G.

Sketch of proof: This follows immediately once you define

$$\varphi(x_1^{n_1}\dots x_N^{n_N}) = \varphi(x_1)^{n_1}\dots\varphi(x_N)^{n_N}$$

Let S be a subset of a group G which does not contain 1. Then S defines a homomorphism φ from the free group on S to G by letting $\varphi(x) = x$ for all $x \in S$. If φ is surjective then S is said to generate G. The elements of S are then called the generators of G. This definition is consistent with the definition in 2.1.12.

2.1.32 Relations among generators. Let *S* be a subset of a group *G* which generates *G* and such that $1 \notin S$. Let *F* the free group on *S*, and φ the unique homomorphism from *F* to *G* which is the identity on *S*. Then Theorem 2.1.27 implies that $F/\ker \varphi$ is isomorphic to *G*. Any element *x* of ker φ satisfies $\varphi(x) = 1_G$ and these equations (and sometimes just the elements $x \in \ker \varphi$) are called *relations among the generators*. If ker $\varphi = \{1\}$ then *F* and *G* are isomorphic. *G* is then also called a free group.

Example: Consider the group $(\mathbb{Z}_m, +)$, $m \in \mathbb{N}$ and $S = \{[1]_m\}$. Then $F \cong (\mathbb{Z}, +)$. The associated homomorphism $\varphi : \mathbb{Z} \to \mathbb{Z}_m$ is given by $n = qm + r \mapsto \varphi(n) = r$ where $r \in \{0, ..., m-1\}$ and $q \in \mathbb{Z}$. The relations are of the form $\varphi(m) = 0$, $\varphi(2m) = 0$ etc.

Let F be a free group on S, R a subset of F and N the smallest normal subgroup including R. Let φ be the canonical homomorphism from F to F/N. Then F/N (and any group isomorphic to it) is called the group generated by S with defining relations R. For example, $(\mathbb{Z}_m, +)$ is the group generated by $\{1\}$ with defining relation m1 = 0.

Theorem. Any group G is the image of a free group F on a set S under a homomorphism φ . Any subset of F which generates the kernel of φ is a set of defining relations among the elements of S.

2.1.33 Free abelian groups. Let F be the free group on S and consider the relations ab = ba for all $a, b \in S$, i.e., let N be the smallest normal subgroup of F generated by $\{a^{-1}b^{-1}ab: a, b \in S\}$. Then F/N is abelian. It is called the *free abelian group* on S.

CHAPTER 3

Rings

3.1. General Ring Theory

3.1.1 Rings. Let R be a set. Suppose there are two binary operations + and \cdot on R such that the following properties are satisfied:

(a) (R, +) is an abelian (commutative) group and

(b) \cdot is associative and (left and right) distributive over +.

Then $(R, +, \cdot)$ is called a *ring*.

The ring is called commutative if multiplication is commutative.

The identity element of (R, +) is denoted by 0.

Examples: \mathbb{Z} , \mathbb{Q} , \mathbb{R} , and \mathbb{C} are commutative rings. Let R be a ring and X a nonempty set, then the set of all R-valued functions on forms a ring, which is commutative if R is. The set of all $n \times n$ matrices is a non-commutative ring. The set $\{0\}$ can be considered a ring, the zero ring.

3.1.2 Elementary properties. Let $(R, +, \cdot)$ be a ring. The additive identity is denoted by 0 and the additive inverse of $a \in R$ is denoted by -a. The following properties hold then in any ring:

1. 0a = a0 = 0 for all $a \in R$,

2. (-a)b = -(ab) = a(-b).

3.1.3 Identity. A (multiplicative) identity or unity in a ring is a nonzero element 1 such that 1a = a1 = a for all elements a in the ring. If there is such an element the ring is called a ring with identity (or ring with unity).

An identity, if it exists, is unique, for 1e = 1 = e if both 1 and e are identities.

3.1.4 Direct products. Let R and S be rings. We introduce two binary operations on $R \times S$ by $(r_1, s_1) + (r_2, s_2) = (r_1 + r_2, s_1 + s_2)$ and $(r_1, s_1) \cdot (r_2, s_2) = (r_1 \cdot r_2, s_1 \cdot s_2)$. Under these binary operation $R \times S$ is a ring called the direct product of R and S. It is commutative if and only if R and S are commutative.

3.1.5 Integral Domains. If a = 0 or b = 0 then ab = 0. If $a, b \neq 0$ but ab = 0 then a is called a *left zero divisor* and b a *right zero divisor*.

Examples: Every nilpotent matrix is a zero divisor in the ring of matrices. In the ring of real-valued functions on X any function which vanishes on some nonempty subset of X but not on all of X is a zero divisor.

A commutative ring with identity but without zero divisors is called an *(integral) do-main*.

Hence if ab = 0 in an integral domain then a = 0 or b = 0. In an integral domain the so called cancellation law holds. This law asserts that ab = ac and $a \neq 0$ imply b = c.

Example: \mathbb{Z} is an integral domain: given $a, b \in \mathbb{Z} - \{0\}$ and ab = 0 we may assume (by multiplying, perhaps, by -1) that $a, b \in \mathbb{N}$ which is impossible since multiplication is a binary operation in \mathbb{N} .

3.1.6 Units. An element *a* of a ring with identity is called a *unit* if there exists a *b* in the ring such that ab = ba = 1. Also, if *a* is a unit and ab = ba = 1 then *b* is a unit, too. The set of all units in a ring *R* is denoted by R^* .

The zero element can never be a unit since $0b = b0 = 0 \neq 1$ for all b in the ring.

Example: In the ring \mathbb{Z} the number 1 is an identity while both +1 and -1 are units.

Theorem. If R is a ring with identity then (R^*, \cdot) is a group.

Sketch of proof: Let $a, b \in R^*$. Then there exist $\hat{a}, \hat{b} \in R$ such that $\hat{a}a = a\hat{a} = 1$ and $\hat{b}b = b\hat{b} = 1$. Hence $\hat{b}\hat{a}ab = ab\hat{b}\hat{a} = 1$, i.e., $ab \in R^*$ and \cdot is a binary operation on R^* . It inherits associativity from R. $1 \in R^*$ is the identity. Inverses exist by the very definition of R^* .

Since inverses in groups are uniquely determined we have that for every unit a in a ring R there exists exactly one element $b \in R$, denoted by a^{-1} , such that ab = ba = 1.

3.1.7 Fields. A commutative ring with identity is called a *field* if every nonzero element is a unit, i.e., if $R^* = R - \{0\}$.

This definition is compatible with our earlier definition. In particular, a set R with two binary operations + and \cdot is a field if and only if (R, +) and $(R - \{0\}, \cdot)$ are commutative groups and if \cdot is distributive over +.

Any field is an integral domain since ab = 0 and $a \neq 0$ imply that $b = a^{-1}ab = 0$.

3.1.8 Fraction fields. Let R be an integral domain and E the relation on $R \times (R - \{0\})$ defined by $(a, b)E(\alpha, \beta)$ if and only if $a\beta = b\alpha$. The relation E is an equivalence relation. Therefore we may define the set $F = R \times (R - \{0\})/E$. The equivalence class of (a, b) will be denoted by a/b. On F addition and multiplication is now defined by

$$\frac{a/b + c/d}{a/b \cdot c/d} = \frac{(ad + bc)}{(bd)}$$

These definitions are possible since they are independent of the representatives chosen.

Theorem. If R is an integral domain the set F defined above is a field called the fraction field of R.

Sketch of proof: Since $bd \neq 0$ addition and multiplication are binary operations. One checks easily that they are associative and commutative and that multiplication is distributive over addition. The zero element is given by 0/1 and the additive inverse of a/b is given by (-a)/b. An identity of F is given by 1/1. If $a/b \neq 0$ then it has a multiplicative inverse given by b/a.

Examples: The fraction field of the integers are the rational numbers and the fraction field of the ring of polynomial functions on \mathbb{R} (or \mathbb{C}) is the field of rational functions.

3.1.9 Subrings. A subset of a ring R is called a *subring* if it is a ring itself. Similarly, a subset of a field K is called a *subfield* if it is a field itself.

Theorem. A nonempty subset S of a ring R is a subring if and only if a - b and ab are in S whenever $a, b \in S$. If R is abelian then so is S.

Sketch of proof: Assume S is a subring. Then ab and a - b = a + (-b) are in S. Now assume that S is not empty and that $a - b, ab \in S$ for all $a, b \in S$. Then (S, +) is a group by the subgroup criterion. It inherits commutativity from (R, +). Also S is closed under multiplication, i.e., \cdot is a binary operation on S. Associativity and distributivity of multiplication over addition are inherited from R. So is, if necessary, commutativity of multiplication.

Examples: $m\mathbb{Z}$ is a subring of \mathbb{Z} . $C^1(\mathbb{R})$ is a subring of $C^0(\mathbb{R})$.

34

Suppose R is a ring and S is a subring of R. It may now be the case that neither R nor S have an identity, that R has an identity but S does not, that S has an identity but R does not, and that both R and S have an identity. In the last case it may be that the identity of R is also the identity S but R and S can also have different identities.

If R has an identity 1 and if $1 \in S$ then S is called a *unitary subring* of R. **3.1.10 Ideals.** A nonempty subset J of a commutative ring R is called an *ideal* if $a - b \in J$ for all $a, b \in J$ and $ar = ra \in J$ for all $a \in J$ and $r \in R$. This latter property is called the *multiplicative absorption property* of ideals. One may also introduce left and right ideals on noncommutative rings. An ideal is called a *proper* ideal if it is a proper subset of R.

Theorem. Any ideal is a subring.

Examples: If *m* is an integer the set $m\mathbb{Z}$ is an ideal in \mathbb{Z} . $J = \{0\}$ is an ideal, the zero ideal. Also J = R is an ideal, the improper ideal. However, $C^1(\mathbb{R})$ is not an ideal in $C^0(\mathbb{R})$. **3.1.11 Ideals generated by subsets.** Let *S* be a nonempty set of ideals in a commutative ring *R*. Then $\bigcap_{J \in S} J$, the intersection of all elements of *S*, is again an ideal in *R*.

Let R be a commutative ring and M a nonempty subset of R. Then $\langle M \rangle$ denotes the intersection of all ideals in R which contain M. Thus $\langle M \rangle$ is the smallest ideal which contains M.

Theorem. Let R be a commutative ring and M a nonempty subset of R. Then

(1)
$$\langle M \rangle = \{ \sum_{j=1}^{N} (r_j a_j + n_j a_j) : N \in \mathbb{N}, a_j \in M, r_j \in R, n_j \in \mathbb{Z} \text{ for } j = 1, ..., N \}.$$

Sketch of proof: Let J be the set on the right hand side of equation (1). Then J is an ideal containing M. Hence $\langle M \rangle \subset J$. Consider the element $x = \sum_{j=1}^{N} (r_j a_j + n_j a_j) \in J$. By the multiplicative absorption property $\langle M \rangle$ contains $r_j a_j$. Since $\langle M \rangle$ is a group with respect to addition it must contain also the element $n_j a_j$. Again by the group property $\langle M \rangle$ thus contains x, i.e., $J \subset \langle M \rangle$.

If R is a ring with identity then na = n1a and thus ra + na = (r + n1)a. Therefore, in this case,

$$\langle M \rangle = \{ \sum_{j=1}^{N} r_j a_j : N \in \mathbb{N}, a_j \in M, r_j \in R \text{ for } j = 1, ..., N \}.$$

An ideal J is called *finitely generated* if there exists a finite set M such that $J = \langle M \rangle$. If $M = \{a_1, ..., a_n\}$ we will also use the notation $\langle M \rangle = \langle a_1, ..., a_n \rangle$. Note that

$$\langle a_1, ..., a_n \rangle = \{ \sum_{j=1}^n r_j a_j : r_j \in R \text{ for } j = 1, ..., n \}.$$

3.1.12 Principal ideals and principal ideal domains. An ideal is called principal if it is generated by a single element. The zero ideal is always principal and the improper ideal is principal if R has an identity. In fact, in this case $R = \langle u \rangle$ for any unit u in R.

An integral domain in which every ideal is principal is called a *principal ideal domain* (PID).

Example: \mathbb{Z} is a principal ideal domain. To see this let $J \neq \{0\}$ be an ideal in \mathbb{Z} and a the smallest positive integer in J. Then $\langle a \rangle \subset J$. Let $b \in J$. Then b = qa + r for suitable $q, r \in \mathbb{Z}$ with $0 \leq r < a$. Since $qa \in \langle a \rangle \subset J$ we have that $r = b - qa \in J$. Thus r = 0 since there is no positive integer in J smaller than a. This proves $b = qa \in \langle a \rangle$ and hence $J = \langle a \rangle$.

3. RINGS

Theorem. A commutative ring R with identity is a field if and only if its only ideals are the zero and the improper ideal. In particular, every field is a principal ideal domain.

Sketch of proof: Suppose R is a field. Let J be any ideal other than the zero ideal. Then J contains a nonzero element a and, since a is a unit, $R = \langle a \rangle \subset J \subset R$, i.e., J = R. Next suppose the only ideals of R are the zero and the improper ideal. Consider a nonzero element r in R. Then $\langle r \rangle = R$ and hence 1 = rs for some $s \in R$, i.e., r is a unit. Hence all nonzero elements of R are units, i.e., R is a field.

3.1.13 Residue class rings. Let S be a subring of the ring R. Then (S, +) is a subgroup of (R, +). Recall that the cosets $a + S = \{a + s : s \in S\} = S + a$ of S form a partition of R. The set of all cosets of S is denoted by R/S. Since (R, +) is abelian (S, +) is a normal subgroup of (R, +). Hence R/S is an abelian group with respect to subset addition (i.e., $M + N = \{m + n : m \in M, n \in N\}$ for any $M, N \subset R$).

Theorem. Let S be an ideal in the commutative ring R. Define addition and multiplication of cosets by (a + S) + (b + S) = (a + b) + S and (a + S)(b + S) = ab + S. Then $(R/S, +, \cdot)$ is a commutative ring called the quotient ring or residue class ring of R modulo S.

Sketch of proof: Since addition of cosets coincides with subset addition we get (as explained above) that (R/S, +) is an abelian group. Next we prove that multiplication is well defined. If a + S = c + S and b + S = d + S then $a - c, b - d \in S$. Hence $ab - cd = a(b - d) + (a - c)d \in S$ since S is an ideal. Therefore ab + S = cd + S, i.e., multiplication of cosets is well defined. The validity of the associative, commutative and distributive laws follows now in a straightforward manner.

3.1.14 Modular arithmetic. For any natural number m the set $m\mathbb{Z}$ is an ideal in \mathbb{Z} . Hence $\mathbb{Z}_m = \mathbb{Z}/m\mathbb{Z}$ is the residue class ring of \mathbb{Z} modulo $m\mathbb{Z}$. We will denote the element of $a + m\mathbb{Z} \in \mathbb{Z}_m$ by a_m and call it the *residue class* mod m of a. If $a_m = b_m$ we say that aand b are *congruent* modulo m. This happens if and only if a - b is an integer multiple of m. Note that 0_m and 1_m are the additive and multiplicative identity, respectively, and that $(-a)_m$ is the additive inverse of a_m .

Theorem. \mathbb{Z}_m is a field if and only if *m* is a prime number.

Sketch of proof: If m is not prime let $a, b \in \{2, ..., m-1\}$ be such that ab = m. Then $a_m b_m = m_m = 0_m$, i.e., \mathbb{Z}_m has zero divisors and hence is not a field.

If m is prime consider $x \in \{1, ..., m-1\}$, i.e., $x_m \neq 0$ in \mathbb{Z}_m . Then gcd(x, m) = 1 and thus there exist integers k, j such that 1 = kx + jm. Now $k_m x_m = (kx)_m = (1 - jm)_m = 1_m$, i.e., $x_m^{-1} = k_m$. Hence every nonzero element in \mathbb{Z}_m is a unit, i.e., \mathbb{Z}_m is a field. \square **3.1.15 Prime ideals.** An ideal J in a commutative ring R with identity is called a prime ideal if $ab \in J$ implies $a \in J$ or $b \in J$.

Theorem. J is a prime ideal in R if and only if R/J is an integral domain.

The improper ideal R is always prime, and the zero ideal is prime if and only R itself is an integral domain. Proper maximal ideals are prime. In this case R/J is in fact a field.

3.2. Ring Homomorphisms

3.2.1 Ring homomorphisms. Let R_1 and R_2 be rings. A mapping $f : R_1 \to R_2$ is called a *(ring) homomorphism* if f(a + b) = f(a) + f(b) and f(ab) = f(a)f(b) for all $a, b \in R_1$. It is called a *(ring) isomorphism* if it is a bijective homomorphism and a *(ring) automorphism* if it is a ring isomorphism from R_1 to itself.

Examples: $f : \mathbb{Z} \to \mathbb{Z}_m, a \mapsto [a]_m$ is a ring homomorphism. The function $f : R_1 \to R_2, a \mapsto 0$ for all $a \in R_1$ is a homomorphism, called the zero homomorphism.

3.2.2 Basic properties of homomorphisms. The composition of homomorphisms is a homomorphism. Moreover, if $f : R_1 \to R_2$ is a ring homomorphism, the following properties hold:

- $f(0) = 0, \ f(-a) = -f(a).$
- $f(R_1)$ is a subring of R_2 .
- ker $f = \{x \in R_1 : f(x) = 0\}$ is an ideal of R_1 .
- f is injective if and only if ker $f = \{0\}$.
- If R_1 has an identity, if f is onto, and if $R_2 \neq \{0\}$ then R_2 has an identity, the image of every unit in R_1 is a unit in R_2 and f(1) = 1.
- If R_1 has an identity, if $f \neq 0$, and if R_2 is an integral domain then the image of every unit in R_1 is a unit in R_2 and f(1) = 1.

3.2.3 Domains as subsets of their fraction fields. Let R be an integral domain and F its fraction field. The mapping $\phi : R \to F$ defined by $a \mapsto a/1$ is an injective homomorphism. One often identifies R with its image $\phi(R) \subset F$, i.e., one usually considers R as a subset of F.

3.2.4 Isomorphy of rings. Two rings R_1 and R_2 are called isomorphic if there exists a ring isomorphism from R_1 to R_2 . The relation "is isomorphic to" is an equivalence relation on the set of all rings.

3.2.5 Canonical homomorphism. Given an ideal S of a commutative ring R the map $\phi: R \to R/S, a \mapsto a + S$ is called the canonical (ring) homomorphism from R to R/S. The kernel of ϕ is precisely the set S.

3.2.6 The fundamental isomorphism theorem for commutative rings. The following theorem relates homomorphisms with quotient rings.

Theorem. If $\phi : R_1 \to R_2$ is a surjective homomorphism between commutative rings and if $\nu : R_1 \to R_1/(\ker \phi)$ is a canonical homomorphism then there exists an isomorphism $\mu : R_1/(\ker \phi) \to R_2$ such that $\mu \circ \nu = \phi$.

Example: Suppose $a \in \mathbb{R}$ and let $\phi : C^0(\mathbb{R}) \to \mathbb{C}$, $f \mapsto f(a)$. Then ϕ is a surjective homomorphism whose kernel is the ideal ker $\phi = \{f \in C^0(\mathbb{R}) : f(a) = 0\}$. The theorem asserts now that $C^0(\mathbb{R})/(\ker \phi)$ is isomorphic to \mathbb{C} .

3.3. Unique Factorization

3.3.1 Divisors and multiples. Let R be a commutative ring. A nonzero element $a \in R$ is called a *divisor* or a *factor* of $b \in R$ if there exists $q \in R$ such that b = aq. We will also say that a divides b or that b is a multiple of a.

If R is an integral domain the element a is called a proper divisor of b if b = aq and neither a nor q are units. Otherwise it is called an improper divisor. If b is a unit it has only improper divisors namely all units of R. Every nonzero element of the ring is a divisor of zero.

3.3.2 Associates. Let R be an integral domain and $a, b \in R$. Then we shall say a is an associate of b or that a is associated with b if there exists a unit u such that a = ub. The relation "is an associate of" is an equivalence relation. Therefore we shall also say that a and b are associates. The equivalence class of 0 contains only 0. The equivalence class of 1 is the set of all units of R.

Example: The class of associates of $n \in \mathbb{Z}$ is $\{n, -n\}$.

3.3.3 The greatest common divisors. Let R be an integral domain. An element $d \in R$ is called a greatest common divisor of the elements $a, b \in R$ if d is a divisor of both a and b and if any common divisor of a and b is also a divisor of d. If d is a greatest common divisor

of a and b then so is any associate of d and any two greatest common divisor of a and b are associates. The notation gcd(a, b) is used to denote any greatest common divisor of a and b.

3.3.4 Prime and irreducible Elements. A nonzero element p of an integral domain is called *prime* if it is not a unit and if p divides a or b whenever p divides ab.

A nonzero element p of an integral domain is called *irreducible* if it is not a unit and has only improper divisors. In other words an irreducible element p allows only "trivial" factorizations p = up' where u is a unit.

Theorem. A prime element is irreducible.

Sketch of proof: Assume that p = ab is prime. Then p divides a or b, say a, i.e., a = pc. Hence p = ab = pcb which implies cb = 1 since cancellation is allowed in an integral domain. Thus b is a unit and hence p is irreducible.

3.3.5 Divisibility and ideals. Divisibility is strongly related to containment of ideals:

Theorem. Let R be an integral domain and $a, b \in R, b \neq 0$. Then

1. $\langle a \rangle = R$ if and only if a is a unit.

2. $\langle a \rangle \subset \langle b \rangle$ if and only if b is a divisor of a. Also $\langle a \rangle = \langle b \rangle$ if and only if b is an improper divisor of a, i.e., if and only if a and b are associates.

3. *a* is irreducible if and only if $\langle a \rangle \neq \{0\}$ is maximal among all proper principal ideals, i.e., there is no principal ideal but *R* and $\langle a \rangle$ itself which contains $\langle a \rangle$.

Let J_1 and J_2 be ideals in the integral domain R. One calls an ideal J_1 a multiple of J_2 or J_2 a factor or divisor of J_1 if $J_1 \subset J_2$. The greatest common divisor of J_1 and J_2 , denoted by $gcd(J_1, J_2)$, is the ideal generated by $J_1 \cup J_2$. The least common multiple of J_1 and J_2 is the ideal $J_1 \cap J_2$.

3.3.6 Unique factorization domains. Let R be an integral domain and a a nonzero element of R which is not a unit. Two factorizations $a = a_1...a_n = b_1...b_m$ into irreducibles $a_1, ..., b_m$ are called equivalent if n = m and if there is a permutation $\sigma \in S_n$ such that a_j is associated with $b_{\sigma j}$ for j = 1, ..., n. If all factorizations of a into irreducibles are equivalent one says that the factorization of a is essentially unique.

An integral domain R is called a *unique factorization domain* (UFD) if it satisfies the following conditions:

1. every nonzero element which is not a unit is a finite product of irreducibles,

2. for every nonzero element which is not a unit all factorizations into irreducibles are essentially unique.

Any two elements in a UFD which are not both equal to zero have a greatest common divisor.

3.3.7 Primality and irreducibility. In Theorem **3.3.4** we have shown that a prime element is always irreducible. If the converse is also true then factorizations are essentially unique. A more precise statement is the following

Theorem. Let R be an integral domain in which every nonzero non-unit has a factorization into finitely many irreducibles. Then these factorizations are essentially unique (i.e., R is a unique factorization domain) if and only if every irreducible element is prime.

Sketch of proof: Assume that R is a unique factorization domain. Let $p \in R$ be irreducible and assume p divides ab, i.e., ab = pc. First consider the case that one of a, b, say b is a unit. Then $a = pcb^{-1}$ and hence p divides a. Hence assume that neither a nor b is a unit. Then let $a_1...a_m$ and $b_1...b_n$ be a factorizations of a and b, respectively, into irreducibles. Since pc = ab the product $a_1...a_m b_1...b_n$ is (an essentially unique) factorization

38

of pc into irreducibles. The unique factorization property implies now that p is associated with one of the factors $a_1, ..., b_n$ and hence divides a or b. Thus p is prime.

The proof showing that factorizations are essentially unique when every irreducible element of R is prime follows very closely the analogous part of the proof in the Fundamental Theorem of Arithmetic.

Corollary. The terms prime and irreducible coincide in unique factorization domains.

3.3.8 Quadratic extensions of the integers. An integer n different from zero or one is called *square-free* if no square number other than 1 is a factor of n. In other words an integer is square-free if it is equal to -1, if it is irreducible, or if it is a product of irreducibles in which no factor occurs more often than once.

Let *n* be a square-free integer. Then $\mathbb{Z}[\sqrt{n}] = \{a + b\sqrt{n} : a, b \in \mathbb{Z}\}$ is a subring of \mathbb{C} . Note that $\mathbb{Z}[\sqrt{n}]$ is commutative and contains the multiplicative identity $1 + 0\sqrt{n}$. Note that $a + b\sqrt{n} = c + d\sqrt{n}$ if and only if a = c and b = d. As a subring of a field $\mathbb{Z}[\sqrt{n}]$ can not have zero divisors. Hence $\mathbb{Z}[\sqrt{n}]$ is an integral domain.

3.3.9 Factorization for quadratic extensions of the integers. Let n be a square-free integer and introduce the function $N : \mathbb{Z}[\sqrt{n}] \to \mathbb{N}_0$, $a + b\sqrt{n} \mapsto |a^2 - nb^2|$. This function has the following properties:

1. N(r) = 0 if and only if r = 0.

2. N(rs) = N(r)N(s) for all $r, s \in R$.

3. N(u) = 1 if and only if u is a unit in R.

4. If N(r) is prime in \mathbb{Z} then r is an irreducible in $\mathbb{Z}[\sqrt{n}]$.

Theorem. Let *n* be a square-free integer and *r* an element of $\mathbb{Z}[\sqrt{n}]$ which is neither zero nor a unit. Then *r* is irreducible or a finite product of irreducibles.

Sketch of proof: Let P be the set of all elements in $\mathbb{Z}[\sqrt{n}]$ which are irreducible or can be expressed as a finite product of irreducibles. Then define

$$S = \{k \in \mathbb{N} : \forall r : N(r) = k + 1 \to r \in P\}.$$

Then $1 \in S$. Assume that $\{1, ..., n\} \subset S$. Let r be such that N(r) = n + 2. Then r is irreducible or r = st where neither s nor t is a unit. Therefore $N(s), N(t) \in \{2, ..., n + 1\}$ and $s, t \in P$. Thus $r \in P$ and $n + 1 \in S$. This proves that $S = \mathbb{N}$ and that every nonzero element which is not a unit is in P.

However, $\mathbb{Z}[\sqrt{n}]$ need not be a UFD. For example, if n = -5 then $6 = (1 + \sqrt{-5})(1 - \sqrt{-5}) = 2 \cdot 3$. Using the function N one may show that the elements 2, 3, and $1 \pm \sqrt{-5}$ are irreducible and that $1 + \sqrt{-5}$ is not associated with either 2 or 3. Hence we have found inequivalent factorizations of $6 \in \mathbb{Z}[\sqrt{-5}]$. This implies also that the terms prime and irreducible do not coincide in $\mathbb{Z}[\sqrt{-5}]$. Indeed 2 divides $6 = (1 + \sqrt{-5})(1 - \sqrt{-5})$ but it divides neither of the factors.

3.3.10 Noetherian Domains. An integral domain is called a *Noetherian domain* if every ideal in the domain is finitely generated. In particular, principal ideal domains are Noetherian domains.

Theorem. The following conditions on an integral domain R are equivalent:

- (1) R is Noetherian.
- (2) Every ascending chain $J_1 \subset J_2 \subset ...$ of ideals in R terminates, i.e., $J_N = J_{N+1} = ...$ for some suitable natural number N.
- (3) Every nonempty set Σ of ideals in R has a maximal element, i.e., there is an element $M \in \Sigma$ such that whenever $J \in \Sigma$ includes M then J = M.

Sketch of proof: Given the chain $J_1 \subset J_2 \subset ...$ define $J = \bigcup_{k=1}^{\infty} J_k$. Then J is an ideal. If R is Noetherian then $J = \langle a_1, ..., a_n \rangle$ for suitable elements $a_1, ..., a_n$. However, there exists N such that $a_j \in J_N$ for all j = 1, ..., n and hence $J = J_N$. Thus 1. implies 2.

Assume Σ contains no maximal element. Then, for every $J \in \Sigma$ the set $X(J) = \{K \in \Sigma : J \subset K, J \neq K\}$ is not empty. Therefore the set $C = \{X(J) : J \in \Sigma\}$ is a nonempty collection of nonempty sets. By the axiom of choice there exists a function $F : X(J) \mapsto K \in X(J)$ and hence there exists a function $f : \Sigma \to \Sigma : J \mapsto F(X(J))$. Since Σ is not empty there exists $J_1 \in \Sigma$. By the recursion theorem there is a unique function $u : \mathbb{N} \to \Sigma$ such that $u(1) = J_1$ and u(n+1) = f(u(n)). The sequence $u(1), u(2), \dots$ is an ascending chain of ideals in R. Assuming statement 2. shows that there exists $N \in \mathbb{N}$ such that u(N+1) = u(N) which is impossible. Hence Σ must contain a maximal element and thus 2. implies 3.

Finally, let J be any nontrivial ideal in R and Σ the set of all finitely generated ideals which are contained in J. Then $\Sigma \neq \{\}$ and hence Σ has a maximal element $J_0 = \langle a_1, ..., a_n \rangle$. Of course $J_0 \subset J$. Let $a \in J - J_0$. Then $\langle a_1, ..., a_n, a \rangle$ is in Σ but strictly bigger than J_0 which is impossible. Hence a does not exist and $J - J_0$ is empty, i.e., $J = J_0$ which shows that 3. implies 1.

3.3.11.

Theorem. If n is a negative square-free integer then $\mathbb{Z}[\sqrt{n}]$ is Noetherian.

Sketch of proof: Let $J \neq \{0\}$ be an ideal in $\mathbb{Z}[\sqrt{n}]$. If $x + y\sqrt{n}$ is a nonzero element of J then $x^2 - ny^2 = (x + y\sqrt{n})(x - y\sqrt{n})$ is also an element of J which lies on the positive real axis. Let a be the smallest element of $J \cap \mathbb{N}$. Next note that J contains elements which lie in the upper half plane H. The set $S = \{\operatorname{Im}(z)/\sqrt{-n} : z \in J \cap H\}$ is a nonempty subset of \mathbb{N} and hence has a minimum. Let b be the element of J whose imaginary part equals $\sqrt{-n}\min(S)$ and whose real part has the smallest possible nonnegative value. Next let $L = \{ma + kb : m, k \in \mathbb{Z}\}$ (such sets are called lattices). Since $a, b \in J$ we get that $L \subset J$. We now prove that $J \subset L$. Let $z \in J$. By solving two linear equations in two unknowns we may determine two real numbers γ and δ such that $z = \gamma a + \delta b$. Hence z = a(m+r) + b(k+s) where $m, k \in \mathbb{Z}$ and $0 \leq r, s < 1$. Therefore $z - (ma + kb) = ra + sb \in J \cap H$ which shows that r = s = 0 and $z \in L$.

Let $P = \{ta + sb \in \mathbb{C} : 0 \le s, t < 1\}$ (called a fundamental parallelogram associated with L). The area of P is given by $a \operatorname{Im}(b)$ and is therefore an integer multiple of $\sqrt{-n}$.

Now let $J_1 \,\subset J_2 \,\subset \ldots$ be an ascending chain of ideals in $\mathbb{Z}[\sqrt{n}]$ and P_1, P_2, \ldots the associated fundamental parallelograms respectively spanned by $(a_1, b_1), (a_2, b_2), \ldots$ If J_{k+1} is strictly larger than J_k then P_k contains an element of J_{k+1} other than zero. Consequently, $a_{k+1} < a_k$ or $\operatorname{Im}(b_{k+1}) < \operatorname{Im}(b_k)$ and the area of P_{k+1} is strictly smaller than the area of P_k . When A(P) denotes the area of the parallelogram P then $A(P_1)/\sqrt{-n}, A(P_2)/\sqrt{-n}, \ldots$ is a nonincreasing sequence of natural numbers which must converge. This implies that the sequence is eventually constant and hence that eventually there can be no more strict inclusions among the ideals J_k , i.e., the chain $J_1 \subset J_2 \subset \ldots$ terminates and $\mathbb{Z}[\sqrt{n}]$ is Noetherian.

This proof fails when n > 0 since the set $\mathbb{Z}[\sqrt{n}]$ is then a subset of the real line.

3.3.12 Factorization in Noetherian domains. The importance of Noetherian domains is due to the following fact.

Theorem. In a Noetherian domain every nonzero element which is not a unit is either irreducible or a product of finitely many irreducibles.

Sketch of proof: Assume the contrary were true. Then there exists a nonzero non-unit a_1 which is not irreducible nor can it be factored into a product of finitely many irreducibles. Since a_1 is not irreducible $a_1 = a_2b_2$ where neither a_2 nor b_2 is a unit. At least one of these, say a_2 , is not irreducible nor can it be factored into a product of finitely many irreducibles since otherwise a_1 could be factored in this manner. Thus

(2)
$$a_1 = a_2 b_2, \quad a_2 = a_3 b_3, \quad a_3 = a_4 b_4, \quad \dots$$

where all a_j, b_j are nonzero non-units and none of the a_j is irreducible or can be factored into a product of finitely many irreducibles. Equations (2) show that $\langle a_1 \rangle \subset \langle a_2 \rangle \subset \dots$. By Theorem 3.3.10 this chain terminates, i.e., $\langle a_N \rangle = \langle a_{N+1} \rangle$ for some N. Therefore $a_{N+1} = ca_N = cb_{N+1}a_{N+1}$. Since an integral domain allows cancellation we have that $cb_{N+1} = 1$, i.e., that b_{N+1} is a unit, the desired contradiction.

3.3.13 Every PID is a UFD. Since a PID is Noetherian factorization of nonzero nonunits is always possible. We show now that it is essentially unique. Suppose p is an irreducible element in the PID R and that p divides the product ab. If p does not divide a then $a \notin \langle p \rangle$. Since $\langle p \rangle$ is maximal among all proper principal ideals and hence among all proper ideals we obtain that $\langle a, p \rangle = R$. Hence $1 \in \langle a, p \rangle$, i.e., there exist $x, y \in R$ such that 1 = ax + py and b = abx + pby which shows that p divides b. This implies that p is prime, i.e., in a PID every irreducible element is prime and hence we have the

Theorem. Every PID is a UFD.

3.3.14 Euclidean domains. An integral domain R is called a *Euclidean domain* if there exists a function $N : R \to \mathbb{N}_0$ with the following properties:

1. N(r) = 0 if and only if r = 0.

2. N(rs) = N(r)N(s) for all $r, s \in R$.

3. If $a, b \in R$ and $a \neq 0$ then there exist q and r in R such that b = aq + r and N(r) < N(a). Note that $N(1) = N(1)^2 = 1$. Let u be a unit in the Euclidean domain R. Then $1 = N(uu^{-1}) = N(u)N(u^{-1})$ and hence N(u) = 1. Conversely, if N(u) = 1 then there exists a $q \in R$ such that 1 = qu, i.e., q is a unit.

Example: \mathbb{Z} is a Euclidean domain where N(n) = |n|. Every field is a Euclidean domain when one defines N(r) = 1 for all $r \neq 0$.

Imitating the proof for \mathbb{Z} yields the following

Theorem. A Euclidean domain is a PID and hence a UFD.

We know that $\mathbb{Z}[\sqrt{-5}]$ is not a UFD. Therefore there must ideals which are not principal. Indeed, let $J = \langle 2, 1 + \sqrt{-5} \rangle$. Then $x + y\sqrt{-5}$ is in J if and only if x - y is an even integer. Therefore J is a proper ideal. Now assume J is principal and generated by an element r. Then both 2 and $1 + \sqrt{-5}$ are multiples of r. Therefore N(r) divides both N(2) = 4 and $N(1 + \sqrt{-5}) = 6$. Since r is not a unit N(r) = 2. However, N(z) is different from 2 for every $z \in \mathbb{Z}[\sqrt{-5}]$. This gives a contradiction to the assumption that J is principal.

3.3.15 Gaussian integers. The integral domain $\mathbb{Z}[i] = \{a + bi : a, b \in \mathbb{Z}\}$ is called the set of *Gaussian integers*. Since $i = \sqrt{-1}$ the Gaussian integers are a quadratic extension of the integers. Its units are ± 1 and $\pm i$.

Theorem. The Gaussian integers form a Euclidean domain.

Sketch of proof: Let $N : \mathbb{Z}[i] \to \mathbb{N}_0$, $a+ib \mapsto a^2+b^2$. We only have to prove that division with a remainder is possible. Let $\alpha, \beta \in \mathbb{Z}[i]$ and suppose $\alpha = a+ib \neq 0$. Then the complex number β/α can be written as m+r+i(n+s) where $m, n \in \mathbb{Z}$ and $-1/2 \leq r, s \leq 1/2$. Therefore $z_0 = \beta - (m+in)\alpha = (r+is)(a+ib) = (ar-bs)+i(rb+as)$ is a Gaussian integer. But $N(z_0) = (a^2+b^2)(r^2+s^2) \leq N(\alpha)/2 < N(\alpha)$.

3. RINGS

3.4. Polynomials

3.4.1 Polynomials. Let R be a ring and

$$P = \{ f : \mathbb{N}_0 \to R : (\exists N : \forall n > N : f(n) = 0) \}.$$

Hence $f \in P$ can be represented by a sequence $(a_0, a_1, a_2, ...)$ of elements of R where only finitely many of the elements a_j are different from zero. We define addition and multiplication in P in the following way: let $f, g \in P$ then

$$(f+g)(n) = f(n) + g(n),$$

 $(fg)(n) = \sum_{j+k=n} f(j)g(k).$

Thus addition and multiplication are binary operations. It is easy to check that they are associative. Also addition is commutative and multiplication is distributive over addition. The element (0, 0, ...) is an additive identity and for any $f = (a_0, a_1, ...) \in P$ the element $-f = (-a_0, -a_1, ...)$ is an additive inverse of f. Therefore P is a ring called the *ring of polynomials over* R. An element of P is called a *polynomial (with coefficients in* R).

3.4.2 Basic properties of polynomial rings. The ring P of polynomials over R is commutative if and only if R is commutative. If R has an identity 1, then so does P namely the element (1, 0, 0, ...). Conversely, if $(b_0, b_1, ...)$ is an identity of P then, for any $a \in R$ we have $(a, 0, 0, ...)(b_0, b_1, ...) = (ab_0, ab_1, ...)$ which shows that b_0 is an identity in R (and that $b_1 = b_2 = ... = 0$).

The mapping $\iota : R \to P$, $a \mapsto (a, 0, ...)$ is an isomorphism between R and $\iota(R) \subset P$. Therefore we will subsequently abbreviate the element (a, 0, ...) by a when no confusion can arise.

If R has an identity 1 let x = (0, 1, 0, 0, ...). Then $x^2 = (0, 0, 1, 0, ...)$ and, if $n \in \mathbb{N}_0$, $x^n = (0, ..., 0, 1, 0, ...)$ where 1 is in the $(n + 1)^{\text{st}}$ slot. Therefore

$$(a_0, a_1, \dots, a_n, 0, \dots) = a_0 + a_1 x + \dots + a_n x^n$$

which yields the familiar expression for a polynomial. It must be stressed, however, that a polynomial is not to be considered a function from $R \to R$, i.e., x is not to be thought of as a variable element of R. The element $x = (0, 1, 0, ...) \in P$ is called an *indeterminate* and, for k = 0, ..., n, the element a_k is called a *coefficient (of* x^k). The set P will be denoted by R[x] and will be called the ring of polynomials in one indeterminate over R.

3.4.3 Polynomial functions. Let P be a ring with identity, let R be a unitary subring of P, and $f = a_0 + ... + a_n x^n \in R[x]$ a polynomial over R. Then f gives rise to a function $\hat{f} : P \to P$ through the definition $\hat{f} : p \mapsto a_0 + ... + a_n p^n$. The function \hat{f} is called a polynomial function over R (or P).

The polynomial function $\hat{f} : R \to R$ must not be confused with the polynomial $f : \mathbb{N}_0 \to R$ as the following example shows: let $R = \mathbb{Z}_2$, $f = x^2 + x + 1$, and $g = x^4 + x^3 + x^2 + x + 1$. Then $f \neq g$ but $\hat{f} = \hat{g}$.

Suppose P is commutative and fix $p \in P$. Then $f \mapsto \hat{f}(p)$ is a ring homomorphism from R[x] to P. Therefore $f \mapsto \hat{f}$ is a surjective ring homomorphism from R[x] to the set of polynomial functions over R on P which is therefore a ring, too, called the ring of polynomial functions over R on P.

3.4.4 The degree of a polynomial. Let $f \neq 0$ be a polynomial. Then $\deg(f) = \max\{n : f(n) \neq 0\}$ is called the *degree* of f. No degree is assigned to the zero polynomial. If $\deg(f) = n$ then f(n) is called the *leading coefficient* of f. The polynomial f is called *monic*

3.4. POLYNOMIALS

if it has leading coefficient 1. The polynomial f is called *constant* if it has degree 0 or if f = 0.

If $\deg(f) < \deg(g)$ then $\deg(f+g) = \deg(g)$. If $\deg(f) \le \deg(g)$ then $\deg(f+g) \le \deg(g)$. If one of the leading coefficients of f and g is not a zero divisor then $\deg(fg) = \deg(f) + \deg(g)$ otherwise fg = 0 or $\deg(fg) < \deg(f) + \deg(g)$. If R is an integral domain then we always have $\deg(fg) = \deg(f) + \deg(g)$.

3.4.5 Polynomials in several indeterminates. When R is a ring then

$$R[x_1, ..., x_n] = \{f : \mathbb{N}_0^n \to R : f(k) = 0 \text{ for all but finitely many } k \in \mathbb{N}_0^n\}$$

becomes also a ring upon introduction of the binary operations

$$(f+g)(k_1,...,k_n) = f(k_1,...,k_n) + g(k_1,...,k_n)$$

$$(fg)(k_1,...,k_n) = \sum_{j_1+l_1=k_1} \dots \sum_{j_n+l_n=k_n} f(j_1,...,j_n)g(l_1,...,l_n).$$

 $R[x_1, ..., x_n]$ is called the polynomial ring in n indeterminates over R. For j = 1, ..., n define $x_j \in R[x_1, ..., x_n]$ by

$$x_j(k_1, ..., k_n) = \begin{cases} 1 & \text{if } k_j = 1 \text{ and } k_l = 0 \text{ for } l \neq j \\ 0 & \text{otherwise.} \end{cases}$$

It is then easy to check that $x_j x_k = x_k x_j$. Therefore $f \in R[x_1, ..., x_n]$ may be expressed as

$$f = \sum_{k_n} \dots \sum_{k_1} f(k_1, \dots, k_n) x_1^{k_1} \dots x_n^{k_n}.$$

From this latter equation it is clear that $R[x_1, ..., x_n]$ is isomorphic to $(R[x_1, ..., x_{n-1}])[x_n]$, the ring of polynomials in one indeterminate over $R[x_1, ..., x_{n-1}]$.

3.4.6 The elementary symmetric polynomials. Let R be an integral domain and consider the polynomial

$$f = (x - a_1)...(x - a_n) = x^n + u_{n-1}x^{n-1} + ... + u_0.$$

It is the easy to check that the coefficients u_j are polynomial functions of the roots a_k . In particular

$$u_0 = (-1)^n a_1 \dots a_n$$
 and $u_{n-1} = -a_1 - \dots - a_n$.

Each of these polynomials is independent of the labeling of the roots, i.e., if σ is a permutation in S_n then $u_j(a_1, ..., a_n) = u_j(a_{\sigma(1)}, ..., a_{\sigma(n)})$. These polynomials are therefore called the elementary symmetric polynomials.

3.4.7 Derivatives. Let R be a ring and $f \in R[x]$. With f we associate the polynomial f' defined by f'(k) = (k+1)f(k+1), called the derivative of f. The derivative of a constant is zero. If $\deg(f) \ge 1$ then $\deg(f') \le \deg(f) - 1$; equality holds if R is an integral domain of characteristic zero.

Similarly, if $f \in R[x_1, ..., x_n]$ and $j \in \{1, ..., n\}$ we associate with f the polynomials $D_j f$ defined by $D_j f(k_1, ..., k_n) = (k_j + 1)f(k_1, ..., k_{j-1}, k_j + 1, k_{j+1}, ..., k_n)$, which is called a partial derivative of f.

These definitions are made without taking any resort to the limit concept. However, if $R = \mathbb{R}$ or $R = \mathbb{C}$ then the polynomial function associated with a derivative of f agrees with the corresponding analytic derivative of the polynomial function \hat{f} .

Theorem. The following statements hold in R[x] respectively $R[x_1, ..., x_n]$:

(1) (f+g)' = f' + g' and $D_i(f+g) = D_i f + D_i g$.

l

(2) If $a \in R$ then (af)' = af' and $D_j(af) = aD_jf$.

3. RINGS

(3)
$$(fg)' = fg' + f'g$$
 and $D_i(fg) = (D_i f)g + f(D_i g)$.

3.4.8 The division algorithm. We now turn to polynomial division.

Theorem. Let R be a commutative ring with identity and f, g polynomials over R. Suppose that $\deg(g) = n$ and that a is the leading coefficient of g. Then there exist polynomials q and r and a nonnegative integer k such that $a^k f = qg + r$ and r = 0 or $\deg(r) < n$.

Sketch of proof: If f = 0 choose q = r = 0. For any nonzero polynomial $h \in R[x]$ let $\gamma(h) = \max\{\deg(h) - n + 1, 0\}$. Then define

$$P = \{ f \in R[x] : (\exists q, r \in R[x] : a^{\gamma(f)}f = qg + r \land (r = 0 \lor \deg(r) < n)) \}$$

and

$$S = \{ j \in \mathbb{N} : \deg(f) = j - 1 \Rightarrow f \in P \}.$$

First we show that $1 \in S$. Let f have degree zero. If $\deg(g) = 0$ choose r = 0, q = f. If $\deg(g) > 0$ choose q = 0 and r = f. Hence $1 \in S$. Now assume that $\{1, ..., j - 1\} \subset S$. If $j \leq n$ one may choose q = 0 and r = f to show that $j \in S$. If j > n denote the leading coefficient f(j-1) of f by b. Then $bx^{j-1-n}g$ is a polynomial of degree j - 1 with leading coefficient ab and hence $h = af - bx^{j-1-n}g$ has degree at most j - 2. By induction hypothesis $h \in P$ and hence there exist $q, r \in R[x]$ such that $a^{\gamma(h)}h = qg + r$ and r = 0 or $\deg(r) < n$. Since $\gamma(f) = j - n > 0$ and $\gamma(h) \leq \max\{j - n - 1, 0\} = j - n - 1$ we find

$$\begin{aligned} a^{\gamma(f)}f &= a^{j-n-1}af = a^{j-n-1}(h+bx^{j-1-n}g) = a^{j-n-1-\gamma(h)}(qg+r) + b(ax)^{j-n-1}g \\ &= (a^{j-n-1-\gamma(h)}q + b(ax)^{j-n-1})g + a^{j-n-1-\gamma(h)}r. \end{aligned}$$

Hence $f \in P$ and $j \in S$. This proves the theorem.

Corollary. If R is a field then R[x] is a Euclidean domain and hence a UFD.

Sketch of proof: Define $N : R[x] \to \mathbb{N}_0$ by N(0) = 0 and $N(f) = 2^{\deg(f)}$ if $f \neq 0$. **3.4.9 Zeros of polynomial functions.** Let R be a ring with identity, f a polynomial over R, and \hat{f} the associated polynomial function. An element $\alpha \in R$ is called a zero or root of \hat{f} if $\hat{f}(\alpha) = 0$.

Theorem. Let R be a commutative ring with identity, $f \in R[x]$, and \hat{f} the associated polynomial function. Then α is a zero of \hat{f} if and only if f is divisible by $x - \alpha$. Moreover, if R is an integral domain then the number of distinct roots of \hat{f} is at most equal to the degree of f.

Sketch of proof: By the division algorithm there exist polynomials q and r such that $f = q(x-\alpha)+r$ and r is a constant. But $0 = \hat{f}(\alpha) = \hat{r}(\alpha)$ since $g \mapsto g(\alpha)$ is a homomorphism. This shows that r = 0. Conversely, if $f = q(x - \alpha)$ then $\hat{f}(\alpha) = 0$.

Now suppose that R is an integral domain, that $\deg(f) = n$, and that $\alpha_1, ..., \alpha_{n+1}$ are distinct zeros of \hat{f} . Then, by induction, there exists a polynomial q such that $f = q(x - \alpha_1)...(x - \alpha_{n+1})$. This implies that $\deg(f) \ge n + 1$, a contradiction.

A zero α of \hat{f} is called a zero of *multiplicity* k if $(x - \alpha)^k$ divides f but $(x - \alpha)^{k+1}$ does not.

3.4.10 Polynomials over integral domains. Let R be an integral domain. Then R[x] and $R[x_1, ..., x_n]$ are integral domains also since, if $0 \neq f, g \in R[x]$, then $\deg(fg) = \deg(f) + \deg(g)$ exists, i.e., $fg \neq 0$. The polynomial $f \in R[x]$ is a unit if and only if $f = ux^0$ and u is a unit of R. Suppose $f = rx^0$ has degree zero. Then f is irreducible in R[x] if and only if r is irreducible in R. Every polynomial of degree one is irreducible.

44

3.4. POLYNOMIALS

3.4.11 Primitive polynomials and Gauss's lemma. Let R be a UFD and $f \in R[x]$. Then f is called *primitive* if its coefficients have no common divisors other than units.

Let f be any nonzero polynomial in R[x]. Then there exists $a \in R$ and a primitive polynomial g such that $f = ax^0g$. The element a is determined up to unit multiples. The class of associates of a (or simply a) is called the content of f and is denoted by c(f). Note that a polynomial is primitive if and only if its content is a unit. In particular, every monic polynomial is primitive.

Theorem Lemma of Gauss. If R is a UFD and $f, g \in R[x]$ then c(fg) = c(f)c(g). In particular, the product of primitive polynomials is primitive.

Sketch of proof: It is only necessary to prove the last assertion of the lemma. Assume that f, g are primitive but that fg is not. Let p be a prime factor of c(fg). If $f = \sum_{k=0}^{n} a_k x^k$ and $g = \sum_{j=0}^{m} b_j x^j$ let s be such that a_k is a multiple of p for k < s but that a_s is not a multiple of p. Similarly, let t be the smallest index such that b_t is not a multiple of p. Then the coefficient of x^{s+t} in fg is $\sum_{j+k=s+t} a_k b_j$. This coefficient and each summand except for $a_s b_t$ is a multiple of p. This is impossible.

Theorem. If R is a unique factorization domain then so are R[x] and $R[x_1, ..., x_n]$.

Sketch of proof: We consider only R[x]. We first show that every nonzero polynomial is a unit or a product of irreducibles. Let

 $P = \{f \in R[x] : f \in R[x]^* \text{ or } f \text{ is irreducible or a product of finitely many irreducibles} \}$ and

$$S = \{ n \in \mathbb{N}_0 : \deg(f) = n \Rightarrow f \in P \}.$$

Let $f = ax^0$ be a constant polynomial which is not a unit and suppose $a = p_1...p_k$ is a factorization into irreducibles. Then $f = (p_1x^0)...(p_kx^0)$ is also a product of irreducibles. Therefore $f \in P$ and $0 \in S$. Next assume that n > 0, that $\{0, ..., n-1\} \subset S$, and that f is a polynomial of degree n. Then $f = cf_1$ for some constant polynomial c and some primitive polynomial f_1 . If f_1 is irreducible then $f \in P$. Hence assume that f_1 is not irreducible and that $f_1 = gh$. Then the degrees of both g and h are positive but smaller than n. Hence g, h and therefore f are in P and $n \in S$. This shows that $S = \mathbb{N}_0$ and that every nonzero polynomial is a unit, an irreducible or a product of finitely many irreducibles.

Next let J be a nontrivial ideal in R[x] and h an element of smallest degree in J. If α is the leading coefficient of h and $f \in J$ then there exists a nonnegative integer k such that h divides $\alpha^k f$. Let $h = \gamma \tilde{h}$ where $\gamma \in R$ and \tilde{h} is primitive. Gauss's lemma then implies that \tilde{h} divides f.

We now show that every irreducible is a prime. Hence let p be irreducible and suppose that p divides fg but not f. If $p = \hat{p}x^0$ then \hat{p} is prime in R and p must divide g. Next assume that deg(p) > 0. Then p is primitive. Let $J = \langle p, f \rangle$ and h = af + bp, a polynomial of smallest degree in J. Suppose $h = \gamma \tilde{h}$ where $\gamma \in R$ and \tilde{h} is primitive. Then, according to what was just proved, \tilde{h} divides p. Suppose $p = s\tilde{h}$. Since p is irreducible \tilde{h} or s is a unit. If s were a unit then p would divide f which it does not. Hence \tilde{h} is a unit and thus h is constant. Now hg = afg + bgp is divided by p, i.e., hg = tp. Using Gauss's lemma again, hc(g) = c(t), i.e., $t = h\tilde{t}$ where \tilde{t} is primitive and this shows $g = \tilde{t}p$, the desired conclusion.

3.4.13 Polynomials over Noetherian domains. Let R be a PID and J an ideal in R[x]. Define

$$L_k = \{ r \in R : (\exists f \in J : \deg(f) = k, f(k) = r) \} \cup \{0\}.$$

3. RINGS

Then $L_0, L_1, ...$ is an ascending chain of ideals in R. Therefore there exists an index q such that $L_k = L_q$ for all $k \ge q$. Since R is a PID we have $L_k = \langle a_k \rangle$ for k = 0, ..., q. For each such k there exists a polynomial $f_k \in J$ of degree k with leading coefficient a_k . Let

$$P = \langle f_0, ..., f_q \rangle \subset J$$

and

$$S = \{ n \in \mathbb{N}_0 : (g \in J, \deg(g) = n) \Rightarrow g \in P \}.$$

Suppose $g = \alpha x^0 \in J$. Then $\alpha \in L_0$ and hence there exists $r \in R$ such that $\alpha = ra_0$. Therefore $g = (rx^0)f_0 \in P$ and $0 \in S$. Next suppose k > 0 and $\{0, ..., k-1\} \subset S$. Let g be a polynomial in J and suppose $\deg(g) = k$ and the leading coefficient of g is α . Then $\alpha \in L_k$ and hence there exists an $r \in R$ such that $\alpha = ra_m$ where $m = \min\{k, q\}$. The polynomial $h = g - rx^{k-m}f_m$ is an element of J with degree less than k. Hence $h \in P$. This shows that $g = h + rx^{k-m}f_m$ is in P, too, and hence that $k \in S$ and $S = \mathbb{N}_0$. Thus every nonzero element of J is in P, i.e., $J = \langle f_0, ..., f_q \rangle$. Hence, if R is a PID then R[x] is Noetherian.

This proof generalizes to the case where R is Noetherian with a little more notational effort. It is performed, for instance, in Zariski and Samuel, Commutative Algebra, Vol. I, Springer 1979, p. 201f. we have thus

Theorem Hilbert's basis theorem. If R is a Noetherian domain then so are R[x] and $R[x_1, ..., x_n]$.

3.4.14 Polynomials with rational coefficients. $\mathbb{Z}[x]$ is a Noetherian UFD. It is, however, not a PID since, for example, the ideal $\langle 2, x \rangle$ is not principal. $\mathbb{Q}[x]$ is a Euclidean domain and hence a PID.

Let $f \in \mathbb{Z}[x] \subset \mathbb{Q}[x]$ and $g, h \in \mathbb{Q}[x]$ such that f = gh. Then there exist integers $\alpha, \beta, \gamma, \delta$ such that $g = (\alpha/\beta)g_1$ and $h = (\gamma/\delta)h_1$ where g_1 and h_1 are primitive polynomials in $\mathbb{Z}[x]$. Hence $\beta\delta f = \alpha\gamma g_1h_1$ holds in $\mathbb{Z}[x]$. Gauss's lemma shows now that $\beta\delta c(f) = \alpha\gamma$ and hence $f = c(f)g_1h_1$, i.e., we have the following

Theorem. If a polynomial with integer coefficients factors in $\mathbb{Q}[x]$ then it factors also in $\mathbb{Z}[x]$. In other words, if a polynomial is irreducible in $\mathbb{Z}[x]$ then it is also irreducible in $\mathbb{Q}[x]$. However, if a polynomial $f \in \mathbb{Z}[x]$ is irreducible in $\mathbb{Q}[x]$ and if f = gh with $g, h \in \mathbb{Z}[x]$ then one of g and h is a constant.

Consider, for example, the polynomial 2x. It is irreducible as an element of $\mathbb{Q}[x]$ since 2 is a unit in $\mathbb{Q}[x]$. In $\mathbb{Z}[x]$, however, neither 2 nor x is a unit.

3.4.15 Polynomials with complex coefficients. $\mathbb{C}[x]$ is a Euclidean domain and hence a PID.

Theorem The fundamental theorem of algebra. If f is a polynomial in $\mathbb{C}[x]$ of degree n then \hat{f} has precisely n roots (counting multiplicities).

Despite its name this theorem is not a purely algebraic inasmuch no purely algebraic proof is known. It can be proven with some knowledge of complex analysis.

Corollary. A polynomial in $\mathbb{C}[x]$ is irreducible if and only if it has degree 1.

A field K with the property that the irreducible elements of K[x] are precisely the polynomials of degree 1 is called algebraically complete.

3.4.16 Polynomials with real coefficients. $\mathbb{R}[x]$ is a Euclidean domain and hence a PID. Every constant polynomial is a unit and every polynomial of degree 1 is irreducible. If f = gh has degree two and neither g nor h are units then both g and h have degree one. Since a polynomial function of degree 1 has a real root \hat{f} has real roots also. Hence

 $f = ax^2 + bx + c$ is irreducible if and only if the discriminant of f, i.e., $b^2 - 4ac$, is less than zero. Next let f be a polynomial of degree larger than 2. If \hat{f} has a real root then f is not irreducible by Theorem 3.4.9. Hence assume that \hat{f} has no real root. Let $\alpha = s + it$ be a complex root and define $g = (x - \alpha)(x - \overline{\alpha}) = x^2 - 2sx + s^2 + t^2 \in \mathbb{R}[x]$. By the division theorem there exist polynomials q and r such that f = qg + r and r = 0 or deg $(r) \leq 1$, i.e., r = cx + d with $c, d \in \mathbb{R}$. However, $0 = \hat{f}(\alpha) = \hat{r}(\alpha)$. Hence, if $c \neq 0$ then \hat{r} has two roots α and -d/c which is impossible. Therefore c = d = 0 and f = qg where neither q nor g is a unit. In summary we have the

Theorem. The following statements hold in $\mathbb{R}[x]$. A nonzero polynomial is a unit if and only if it is constant. A polynomial is irreducible if and only if it has either degree 1 or else degree 2 and a negative discriminant.

3.4.17 Common factors of polynomials. We now consider the problem of determining whether two polynomials have a common factor.

Proposition. Suppose $f, g \in R[x]$ where R is a UFD. Then f and g have a non-constant common factor if and only if there exist nonzero polynomials ϕ and ψ such that $\psi f = \phi g$, $\deg \psi < \deg g$, and $\deg \phi < \deg f$.

Sketch of proof: If f and g both have a factor h then $f = h\phi$ and $g = h\psi$ for suitable polynomials ϕ and ψ . Conversely assume that $\psi f = \phi g$. Not all the prime factors of g can also be factors of ψ if deg $\psi < \deg g$. Hence at least one of them must be a factor of f. \Box **3.4.18 Elimination.** Let R be a UFD and Q the associated quotient field. Suppose that

 $f = a_0 + \dots + a_n x^n$ and $g = b_0 + \dots + b_m x^m$

where $n, m \ge 1$ and a_n and b_m are different from zero. Then, letting

$$\psi = \alpha_1 + \dots + \alpha_m x^{m-1}$$
 and $\phi = \alpha_{m+1} + \dots + \alpha_{n+m} x^{n-1}$,

we obtain

$$\psi f - \phi g = \sum_{j=1}^{n+m} \sum_{k=1}^{n+m} M_{j,k} \alpha_k x^{j-1}$$

where, agreeing that $a_j = 0$ if j < 0 or j > n and $b_j = 0$ if j < 0 or j > m,

$$M_{j,k} = \begin{cases} a_{j-k} & \text{if } k \le m \\ b_{m+j-k} & \text{if } k > m. \end{cases}$$

The requirement that $\psi f = \phi g$ is therefore equivalent with $M\alpha = 0$ when M is the matrix formed by the $M_{j,k}$ and α is the vector formed by the α_j . Now, if f and g have a nonconstant common factor then det M must be zero because $M\alpha = 0$ has then a nontrivial solution. Conversely, if det M = 0 then $M\alpha = 0$ has a nontrivial solution in Q^{n+m} , i.e., there exist $\phi, \psi \in Q[x]$ such that $\psi f = \phi g$. However, multiplying by a suitable $r \in R$ allows to find $\phi, \psi \in R[x]$ such that $\psi f = \phi g$ and hence f and g have a non-constant common factor.

The element det $M \in R$ is called the eliminant or resultant of f and g.

We proved the following

Theorem. Let R be a UFD. The polynomials $f, g \in R[x]$ have a non-constant common factor if and only if the eliminant of f and g vanishes.

3.4.19 Discriminant. If R is a UFD then the eliminant of f and f' is called the discriminant of f. Since (fg)' = fg' + gf' the discriminant of f is zero if and only if f has a non-constant repeated factor.

3.4.20 Homogeneous polynomials. Let R be an integral domain with infinitely many elements.

A polynomial of the form $ax_1^{j_1}...x_n^{j_n} \in R[x_1,...,x_n]$ is called a monomial of degree $j_1 + ... + j_n$. A polynomial $f \in R[x_1,...,x_n]$ is called homogeneous of degree m if it is a sum of monomials of degree m. Alternatively one may define that $f \in R[x_1,...,x_n]$ is homogeneous of degree m if the relationship $\hat{f}(tr_1,...,tr_n) = t^m \hat{f}(r_1,...,r_n)$ holds for the polynomial function \hat{f} for all $t, r_1, ..., r_n \in R$.

With each homogeneous polynomial F in $R[x_1, ..., x_{n+1}]$ of degree m which is not divisible by x_{n+1} there is associated a unique (nonhomogeneous) polynomial f in $R[x_1, ..., x_n]$ of degree, namely

$$F = \sum_{j=0}^{m} a_j x_{n+1}^j \mapsto f = \sum_{j=0}^{m} a_j$$

where $a_j \in R[x_1, ..., x_n]$.

Conversely with each nonhomogeneous polynomial f in $R[x_1, ..., x_n]$ of degree m there is associated a homogeneous polynomial F in $R[x_1, ..., x_{n+1}]$ of degree m for which x_{n+1} is not a factor: suppose $f = \sum_{j=0}^{m} f_j$ where each f_j is homogeneous of degree j. Then $F = \sum_{j=0}^{m} f_j x_{n+1}^{m-j}$ is homogeneous of degree m.

3.4.21 Euler's theorem. Suppose R has infinitely many elements and $F \in R[x_1, ..., x_n]$ is homogeneous of degree m. Then

$$\sum_{j=1}^{n} x_j D_j F = mF.$$

3.4.22 Factorization and elimination. Let R be an integral domain. The following statements hold:

- (1) If f and F are associated then any factor of f is associated with a factor of F and vice versa.
- (2) Any factor of a homogeneous polynomial is homogeneous and the product of homogeneous polynomials is homogeneous.
- (3) F is irreducible if and only if f is.
- (4) The homogeneous polynomials

$$F = a_n x_1^n + \dots + a_0 x_2^n, G = b_m x_1^m + \dots + b_0 x_2^m$$

have a common non-constant factor if and only if the eliminant of f and g is zero. The eliminant of f and g is also called the eliminant of F and G.

(5) If $F \in \mathbb{C}[x_1, x_2]$ is homogeneous of degree m then

$$F = \gamma \prod_{j=1}^{m} (a_j x_1 - b_j x_2)$$

- (6) Let $F, G \in R[x_1, ..., x_{n+1}]$ are homogeneous of degree n and m. Consider them as elements of $(R[x_1, ..., x_n])[x_{n+1}]$. Then there resultant is in $R[x_1, ..., x_n]$. It is either zero or else homogeneous of degree nm.
- (7) The eliminant of $\prod_{j=1}^{n} (x y_j)$ and $\prod_{j=1}^{m} (x z_j)$ is

$$a\prod_{j=1}^{n}\prod_{k=1}^{m}(y_j-z_k)$$

for some $a \in R$.

3.5. Algebraic Geometry

3.5.1 Algebraic varieties. Let $R = \mathbb{C}[x_1, ..., x_n]$ be the ring of polynomials in n indeterminates over the field of complex numbers. An *(affine) algebraic variety* is a subset V of \mathbb{C}^n such that there exists an ideal $J = \langle f_1, ..., f_k \rangle$ in R for which

$$V = \{x \in \mathbb{C}^n : (\forall f \in J : f(x) = 0)\} = \{x \in \mathbb{C}^n : f_1(x) = \dots = f_k(x) = 0\}.$$

(Since the ring of polynomials over \mathbb{C} in n indeterminates is isomorphic to the ring of polynomial functions on \mathbb{C}^n we will not distinguish polynomials and polynomial functions in this section. In fact, the word "polynomial" will always mean "polynomial function" subsequently.)

V is called an *(affine) algebraic hypersurface* if J is principal, i.e., if k = 1. An affine algebraic hypersurface in \mathbb{C}^2 , i.e., when n = 2, is called an *(affine) algebraic curve*.

Algebraic geometry is the mathematical discipline which studies algebraic varieties.

Examples: If n = 1 then R is a PID and every variety is a hypersurface. In fact, in this case, a variety is just the set of zeros of some polynomial and hence a finite isolated set. If $J = \langle f \rangle$, the ideal generating the variety V is maximal, i.e., if f is irreducible (cf. Theorem 3.3.5) then $f = x - \alpha$ for some $\alpha \in \mathbb{C}$ by Corollary 3.4.15. Hence $V = \{\alpha\}$. If J is any ideal distinct from C[x] itself the variety associated with it is not empty by the fundamental theorem of algebra.

Next let n = 2 and k = 1. In this case it is customary to denote the indeterminates by x and y. According to the above definition a variety V is called an algebraic curve. For instance, if f = ax + by + c where $a, b, c \in \mathbb{R}$ then $V = \{(x, y) : ax + by + c = 0\}$ is a subset of $\mathbb{R}^2 \subset \mathbb{C}^2$. In fact, in \mathbb{R}^2 the variety V is represented by a line. Therefore one calls V a line even though it is really a two-dimensional subset of \mathbb{C}^2 (given a reasonable definition of the term "two-dimensional"). A *conic* is an algebraic curve given by a quadratic polynomial $f = a_1x^2 + a_2xy + a_3y^2 + b_1x + b_2y + c$. These include, of course, ellipses, parabolas, and hyperbolas.

3.5.2 The correspondence between ideals and affine sets. Let $R = \mathbb{C}[x_1, ..., x_n]$. For every ideal J in R we have defined above a subset V(J) of \mathbb{C}^n , the algebraic variety of J. Note that $V(\{0\}) = \mathbb{C}^n$, $V(R) = \{\}$, and that $V(J_1) \subset V(J_2)$ if $J_2 \subset J_1$.

Conversely, for any subset A of \mathbb{C}^n (called an *affine set*) we define

$$I(A) = \{ f \in \mathbb{C}[x_1, ..., x_n] : (\forall a \in A : f(a) = 0) \}.$$

Then I(A) is an ideal. If $A \subset B \subset \mathbb{C}^n$ then $I(B) \subset I(A)$. For any $A \in \mathbb{C}^n$ we have $A \subset V(I(A))$ and A = V(I(A)) if and only if A is an algebraic variety. Also, for any ideal J in R we have $J \subset I(V(J))$.

Next note that $\langle x_1 - a_1, ..., x_n - a_n \rangle \subset I(\{a\})$. Conversely, since $x_j = (x_j - a_j) + a_j$ the binomial theorem shows that any $f \in R$ can be written as

$$f = c_{0,...,0} + \sum_{|\alpha|>0} c_{\alpha}(x-a)^{\alpha}$$

using a multiindex $\alpha = (\alpha_1, ..., \alpha_n)$. Now, if $f \in I(\{a\})$ then $c_{0,...,0} = 0$. Every other summand in the above expression for f, however, is, for some $j \in \{1, ..., n\}$, a multiple of $x_j - a_j$. Hence we have the following result:

Theorem. If $a = (a_1, ..., a_n) \in \mathbb{C}^n$ then $\langle x_1 - a_1, ..., x_n - a_n \rangle = I(\{a\}).$

3.5.3 Radicals. Let J be an ideal in R, a commutative ring with identity. Then the set

$$\sqrt{J} = \{r \in R : (\exists n \in \mathbb{N} : r^n \in J)\}$$

is called the *radical* of J. The set \sqrt{J} is an ideal containing J itself.

Let J be an ideal in R such that $\sqrt{J} = R$. Then $1 \in \sqrt{J}$ and hence there exists a natural number n such that $1 = 1^n \in J$ and therefore J = R.

3.5.4 Intersections of algebraic curves. Let $f, g \in \mathbb{C}[x, y]$ be two polynomials with no common divisors other than constants and let F be the fraction field of $\mathbb{C}[x]$. The polynomials f and g are polynomials in y with coefficients in $\mathbb{C}[x] \subset F$, i.e., f and g can be viewed as elements of F[y]. By Gauss's lemma f and g factor over F if and only if they factor over $\mathbb{C}[x]$ (cf. Theorem 3.4.14). Hence, even as elements of F[y] they have no common divisors other than constants. This shows that $\langle f, g \rangle = F[y]$ and thus there exist elements $r, s \in F[y]$ such that 1 = rf + sg. Note that $r = \tilde{r}/p$ and $s = \tilde{s}/p$ where $\tilde{r}, \tilde{s} \in \mathbb{C}[x, y]$ and $p \in \mathbb{C}[x]$. Therefore $p = \tilde{r}f + \tilde{s}g$. Now suppose that (α, β) is a zero of f and g. Then α must be a zero of p. Since p has only finitely many distinct zeros α can assume only finitely many values. Reversing the roles played by x and y shows that β also can assume only finitely many values. Hence f and g intersect in at most finitely many points. Thus we have shown the

Theorem. If $f, g \in \mathbb{C}[x, y]$ have no common factor other than units then the variety $\{(x, y) \in \mathbb{C}^2 : f(x, y) = g(x, y) = 0\}$ consists of finitely many points.

An interesting question is now to find the number of intersection points. Bezout's theorem states that this number equals the product of the degrees of f and g (provided they do not intersect at infinity).

On the other hand varieties in \mathbb{C}^2 generated by ideals which are not principal are not very interesting as point sets. Hence the most interesting varieties in \mathbb{C}^2 are the algebraic curves.

3.5.5 Hilbert's Nullstellensatz. Nullstelle is the German word for a zero of a function. The theorem is as follows:

Theorem Hilbert's Nullstellensatz. Let J be an ideal in $\mathbb{C}[x_1, ..., x_n]$. The algebraic variety V(J) is empty if and only if $J = \mathbb{C}[x_1, ..., x_n]$.

This theorem has the following consequences:

Corollary. Let $R = \mathbb{C}[x_1, ..., x_n]$. Then each of the following two statements is equivalent to Hilbert's Nullstellensatz.

1. Let J be an ideal in R. Then $I(V(J)) = \sqrt{J}$, i.e., $g|_{V(J)} = 0$ if and only if $g^k \in J$ for some $k \in \mathbb{N}$.

2. The ideal J in R is maximal if and only if there exists a point $a = (a_1, ..., a_n) \in \mathbb{C}^n$ such that

$$J = \langle x_1 - a_1, ..., x_n - a_n \rangle = \{ f \in R : f(a) = 0 \} = I(\{a\}).$$

Sketch of proof: First note that always $V(R) = \{\}$ since $1 \in R$ and $\sqrt{J} \subset I(V(J))$.

HN \Leftrightarrow 1.: Let $J = \langle f_1, ..., f_k \rangle$ and let $0 \neq g \in I(V(J))$, i.e., $g|_{V(J)} = 0$. Introduce another indeterminate y and consider $J' = \langle f_1 y^0, ..., f_k y^0, gy - 1 \rangle \subset R[y]$. If $(a, b) \in \mathbb{C}^{n+1}$ is in V(J') then $f_j(a)b^0 = 0$ for j = 1, ..., k and g(a)b - 1 = 0. Hence $a \in V(J)$ but $g(a) \neq 0$. Since this is impossible we have $V(J') = \{\}$ and J' = R[y]. Therefore

$$1 = \sum_{j=1}^{k} h_j f_j y^0 + h_0 (gy - 1)$$

for some $h_0, h_1, ..., h_k \in R[y]$. Let $N = \max\{\deg(h_0), ..., \deg(h_k)\}$. Then $h_j = \sum_{l=0}^N \eta_{j,l} y^l$ and

$$g^{N} = \sum_{j=1}^{k} f_{j} \sum_{l=0}^{N} \eta_{j,l} g^{N-l} (gy)^{l} + (gy-1) \sum_{l=0}^{N} \eta_{0,l} g^{N-l} (gy)^{l}$$

Consider this as an identity in F[y] where F is the fraction field of R. It then holds for all $y \in F$. Choose y = 1/g to obtain

$$g^{N} = \sum_{j=1}^{k} f_{j} \sum_{l=0}^{N} \eta_{j,l} g^{N-l},$$

i.e., $g^N \in J$ and $g \in \sqrt{J}$. This proves $\text{HN} \Rightarrow 1$.

Now assume $V(J) = \{\}$. Then $\sqrt{J} = I(V(J)) = I(\{\}) = R$ and hence J = R.

1. \Rightarrow 2.: Let *J* be a maximal ideal. Then $J = \sqrt{J}$ and hence I(V(J)) = J. Since $I(\{\}) = R$ we have that V(J) is not empty. Suppose $a \in V(J)$ then $J = I(V(J)) \subset I(\{a\})$. Since $I(\{a\}) \neq R$ and since *J* is maximal we get $J = I(\{a\})$. Conversely, assume that *J* is a proper ideal containing $I(\{a\})$. Then $\{\} \neq V(J) \subset V(I(\{a\})) = \{a\}$. Hence $J = I(\{a\})$, i.e., $I(\{a\})$ is maximal.

2. \Rightarrow HN: Let J be a proper ideal in R. Then, since R is Noetherian, there exists a maximal ideal M which contains J. Note that $M = I(\{a\})$ for some $a \in \mathbb{C}^n$. Since $\{a\} = V(M) \subset V(J)$ we have that V(J) is not empty. \Box

CHAPTER 4

Fields

Let P be a ring and R a unitary subring of P. Given a polynomial $f \in R[x_1, ..., x_n]$ we will denote, as before, the associated polynomial function on P^n by \hat{f} . The quotient field of $R[x_1, ..., x_n]$ is denoted by $R(x_1, ..., x_n)$.

4.1. Field Extensions

4.1.1 Fields and subfields. A commutative ring with identity is called a field if all its nonzero elements are units. If F is a field and K a subset of F which is field with respect to the binary operations of F, then K is called a *subfield* of F and F is called a *field* extension over K.

A subset K of F is a subfield if and only if it contains a nonzero element and if $a, b \in K$ and $b \neq 0$ imply that a - b and ab^{-1} are also in K.

4.1.2 Prime fields. A prime field is a field which does not contain a proper subfield.

Note that the intersection of all subfields of a field is a subfield itself. It does not properly contain another subfield, hence it is a prime field. Obviously a given field contains a unique prime field.

Theorem. If Π is a prime field then it is either isomorphic to \mathbb{Q} or else it is isomorphic to \mathbb{Z}_p where p is some prime number.

Sketch of proof: Π contains the identity elements 0 and 1 and, necessarily, the element n1 whenever n is an integer. Note that n1 + m1 = (n + m)1 and that (n1)(m1) = (nm)1, i.e., the set $\mathcal{P} = \{n1 : n \in \mathbb{Z}\}$ is the homomorphic image of \mathbb{Z} under the map $\phi : n \mapsto n1$. By the fundamental isomorphism theorem for rings \mathcal{P} is isomorphic to $\mathbb{Z}/\ker(\phi)$.

If $\ker(\phi) = \{0\}$ then \mathcal{P} is isomorphic to \mathbb{Z} . But the smallest field containing \mathbb{Z} is \mathbb{Q} . If $\ker(\phi)$ is nontrivial then $\ker(\phi) = \langle p \rangle$ for some natural number p. We then have firstly that p > 1. Secondly, since \mathcal{P} (a subset of a field) does not contain zero divisors p must be prime. Now note that $\mathbb{Z}/\langle p \rangle = \mathbb{Z}_p$ is a field if p is prime. \Box

4.1.3 The characteristic of field. Let F be a field and Π its prime field. If Π is isomorphic to \mathbb{Q} we say that Π and F have *characteristic* zero. If Π is isomorphic to \mathbb{Z}_p for some positive prime number p we say that Π and F have characteristic p.

4.1.4 Field extensions. Let F be a field, K a subfield of F and S a subset of F. Then we denote by K(S) the smallest field which contains $K \cup S$, i.e., the intersection of all subfields of F which contain $K \cup S$. We then have

$$K(S) = \{\frac{f(s_1, ..., s_n)}{\hat{g}(s_1, ..., s_n)} : n \in \mathbb{N}_0, s_1, ..., s_n \in S, f, g \in K[x_1, ..., x_n], \hat{g}(s_1, ..., s_n) \neq 0\}.$$

If S has n elements $s_1, ..., s_n$ we will write $K(S) = K(s_1, ..., s_n)$ and we will say that K(S) is the field generated over K by $s_1, ..., s_n$ or the field obtained by adjoining to K the elements $s_1, ..., s_n$.

If $S' \subset S$ then K(S') is a subfield of K(S). If S any subset of F define $C = \{T \subset S : T \text{ is finite}\}$. Then $K(S) = \bigcup_{T \in C} K(T)$.

4.1.5 Field extensions as vector spaces. If F is a field extension over K then F is also a vector space over K where the scalar multiplication is regular multiplication in F. The dimension of F as a vector space over K is called the degree of F over K. It is denoted by [F:K]. If [F:K] is finite then F is called a *finite field extension* over K.

4.1.6 *K*-isomorphisms. Let *F* and *F'* be two field extensions over a field *K*. If there exists an isomorphism $\phi : F \to F'$ such that $\phi|_K$ is the identity we call ϕ a *K*-isomorphism and we say that *F* and *F'* are *K*-isomorphic. The group of *K*-automorphisms of *F* is denoted by $\operatorname{Aut}_K(F)$.

For example the fields $\mathbb{Q}(\sqrt{2})$ and $\mathbb{Q}(-\sqrt{2})$ are \mathbb{Q} -isomorphic.

4.1.7 Simple field extensions. Let F be a field and K a subfield of F. If $s \in F$ define $S = \{\hat{f}(s) : f \in K[x]\}$. Then S is a ring which is isomorphic to K[x]/J for some ideal $J \subset K[x]$. Since K[x] is a PID we have that $J = \{0\}$ or that $J = \langle \varphi \rangle$ for some polynomial φ .

In the first case, where $J = \{0\}$, the element s is called *transcendental* over K and K(s) is called a simple transcendental extension over K. In this case K(s) is isomorphic to K(x), the quotient field of K[x]. On the other hand, since K can be considered a subfield of K(x), K(x) itself is a simple transcendental extension over K.

In the second case, where $J = \langle \varphi \rangle$, note first that φ is uniquely determined up to multiplication by nonzero elements of K. $\langle \varphi \rangle$ is a proper maximal ideal since $S \subset F$ and hence $K[x]/\langle \varphi \rangle$ do not have zero divisors. Therefore $K[x]/\langle \varphi \rangle$ is a field so that K(s), $K[x]/\langle \varphi \rangle$, and S are all isomorphic. Since $\hat{\varphi}(s) = 0$ the polynomial φ has the factor x - sand hence its degree is at least equal to one. We call *s algebraic* over *K* and K(s) a simple algebraic extension over *K*. The polynomial φ is called the minimal polynomial of *s* over *K*.

Now suppose a field K and an irreducible polynomial $\varphi \in K[x]$ are given. Consider the field $H = K[x]/\langle \varphi \rangle$ and the canonical homomorphism $\Psi : K[x] \to H$. The subset $\tilde{K} = \{ax^0 + \langle \varphi \rangle : a \in K\}$ is the isomorphic image of K and hence we consider K as a subset of H by identifying K and \tilde{K} . Note that H = K(s) where $s = \Psi(x) = x + \langle \varphi \rangle$ is an algebraic element of H, in fact $\hat{\varphi}(s) = 0$. Thus the existence of simple algebraic extensions over K can be obtained without an a priori knowledge of a field F containing the element s.

Example: Let $K = \mathbb{R}$ and $\varphi = x^2 + 1$ then $\mathbb{C} \cong \mathbb{R}[x]/\langle x^2 + 1 \rangle \cong \{\hat{f}(i) : f \in \mathbb{R}[x]\} \cong \mathbb{R}(i)$ where $i = x + \langle x^2 + 1 \rangle$.

4.1.8 Conjugate elements. Two elements s and s' of an extension F over K are called *conjugate* if they are algebraic over K and have the same minimal polynomial over K.

4.1.9 The degree of a simple extension. Let K be a field, φ an irreducible polynomial in K[x] of degree n, and $s = x + \langle \varphi \rangle$. Then [K(s) : K] = n, i.e., the degree of the simple algebraic extension K(s) over K is equal to n. The degree of a simple transcendental extension over K is infinite.

Sketch of proof: In the algebraic case one shows, using the division theorem, that $\{x^j + \langle \varphi \rangle : j = 0, ..., n - 1\}$ is a basis of H. Conversely, if [K(s) : K] is finite then K(s) must be algebraic.

For example, if $i = x + \langle x^2 + 1 \rangle$ then 1 and *i* are linearly independent in $\mathbb{R}[x]/\langle x^2 + 1 \rangle$ but $i^2 + 1 = 0$. Hence $[\mathbb{R}[x]/\langle x^2 + 1 \rangle : \mathbb{R}] = 2$ and $\mathbb{R}[x]/\langle x^2 + 1 \rangle \cong \{a + bi : a, b \in \mathbb{R}\}.$

4.1.10 Algebraic field extension. F is called an algebraic (field) extension over K if every element of F is algebraic over K. Any other extension is called a transcendental (field) extension.

Theorem. If F is a field extension over K and if [F : K] = n is finite, then F is an algebraic extension over K which is obtained by the adjunction of finitely many algebraic elements. Every element of F is then a zero of polynomial function over K of degree at most n. Conversely, any extension over a field K which is obtained by adjoining finitely many algebraic elements is a finite extension and hence an algebraic extension.

Sketch of proof: Let r be an element of F. Then $1, r, ..., r^n$ are linearly dependent vectors, i.e., there exists $k_0, ..., k_n \in K$, not all zero, such that $k_0 + k_1r + ... + k_nr^n = 0$. Let $r_1, ..., r_n$ be a basis of F then, obviously, $F = K(r_1, ..., r_n)$. For the converse note that, by 4.1.9, adjoining a single algebraic element to K produces a finite extension and so does adjoining finitely many algebraic elements.

4.1.11 Splitting fields. Let f be a polynomial in K[x] of degree n. Suppose there is a field extension F over K such that f, considered as an element of F[x], factors into n linear factors (polynomials of degree 1). Then one says that F contains all the roots of \hat{f} . The smallest field which contains K and all the roots of \hat{f} is called a *splitting field* of f over K.

Theorem. Let K be a field and $f \in K[x]$. Then there exists a splitting field F over K of f. Moreover, if deg(f) = n then there exist algebraic elements $s_1, ..., s_n$ of F such that $F = K(s_1, ..., s_n)$. In particular, a splitting field over K is algebraic over K.

Sketch of proof: Let φ be an irreducible factor of f and define $K(s_1) = K[x]/\langle \varphi \rangle$ as in 4.1.7. Then $\hat{\varphi}(s_1) = 0$ and hence $\hat{f}(s_1) = 0$. The polynomial f has a factor $x - s_1$ when it is considered as an element of $K(s_1)[x]$, i.e., $f = (x - s_1)g$ for some polynomial $g \in K(s_1)[x]$. This procedure has to be repeated at most n times if $n = \deg(f)$.

4.1.12 Theorem of the primitive element. If K has characteristic zero then every finite extension over K is simple.

Sketch of proof: Let F be the extension under consideration. By 4.1.10 we have that F is generated over K by finitely many, say n, elements. The induction principle gives that the theorem is proved once it is proved for n = 2.

Let F = K(s,t) and suppose that the minimal polynomials of s and t over K are f and g, respectively. Let H be a field containing F and the splitting fields of f and g. Suppose that in H we have $f = \prod_{j=1}^{k} (x - s_j)$ and $g = \prod_{j=1}^{\ell} (x - t_j)$ where $s_1 = s$ and $t_1 = t$. Next consider the polynomial $\varphi_{m,j} = (s_m - s_1)x + t_j - t_1$. We will show below that $s_m \neq s_1$ if $m \neq 1$ and the characteristic of K is zero. Hence $\hat{\varphi}_{m,j}$ has at most one root in K. Since K has infinitely many elements one may choose $c \in K$ such $\hat{\varphi}_{m,j}(c) \neq 0$ whenever $m \in \{2, ..., k\}$ and $j \in \{1, ..., \ell\}$. Define r = cs + t. Then K(s, t) = K(r) as will be shown presently.

Obviously $K(r) \subset K(s,t)$. We need to show that $s, t \in K(r)$. Assume that s is not in K(r) and let $h = \hat{g} \circ (r - cx) \in K(r)[x]$. Note that $\hat{f}(s_m) = 0$ for all $m \in \{1, ..., k\}$ but that $\hat{h}(s_m) = \hat{g}(r - cs_m)) \neq 0$ for all $m \in \{2, ..., k\}$. Let $d \in K(r)[x]$ be a common divisor of f and h. Since f and h split into linear factors when considered as elements of H[x] the same is true for d assuming that $\deg(d) > 0$. This would imply that d is an associate of x - s which is impossible, since $s \notin K(r)$. Therefore d must have degree zero and, by the GCD identity, there are polynomials α and β in K(r)[x] such that $\alpha f + \beta h = d$. This in turn gives that $d = \hat{\alpha}(s)\hat{f}(s) + \hat{\beta}(s)\hat{h}(s) = 0$. Since this is impossible we have $s \in K(r)$. But then we have also $t = r - cs \in K(r)$.

4. FIELDS

It remains to show that the roots of \hat{f} are simple. Suppose \hat{f} has a root s of multiplicity at least two. Since f is irreducible as an element of K[x] the ideal $\langle f \rangle$ is a proper maximal ideal in K[x]. Since the degree of f' is smaller than the degree of f we have f' = 0 or $\langle f, f' \rangle = K[x]$. The second possibility is ruled out since there would then be $\alpha, \beta \in K[x]$ such that $\alpha f + \beta f' = 1$ which contradicts the fact that $\hat{f}(s) = \hat{f}'(s) = 0$. If $f = \sum_{j=0}^{k} a_j x^j$ then $f' = \sum_{j=1}^{k} j a_j x^j$ so that f' = 0 implies $ja_j = 0$ for all $j \in \{1, ..., k\}$. Since K has characteristic zero this shows that $a_1 = ... = a_k = 0$ and hence that $f = a_0 x^0$ which is not irreducible. Therefore any root of \hat{f} is simple.

4.1.13 Extensions of isomorphisms. Let K and \tilde{K} be isomorphic fields and $\tau : K \to \tilde{K}$ the corresponding isomorphism. If $f = \sum_{j=0}^{n} \alpha_j x^j$ is a polynomial in K[x] we say that $\tilde{f} = \sum_{j=0}^{n} \tau(\alpha_j) x^j \in \tilde{K}[x]$ is the polynomial corresponding to f under τ .

Proposition. Let K and \tilde{K} be isomorphic fields and $\tau : K \to \tilde{K}$ the corresponding isomorphism. Suppose f is an irreducible polynomial in K[x] and that \tilde{f} is the polynomial corresponding to f under τ . Let F be a field extension over K containing a root r of f and let \tilde{F} be a field extension over \tilde{K} containing a root s of \tilde{f} . Then there is a unique isomorphism $\rho : K(r) \to \tilde{K}(s)$ such that $\rho|_K = \tau$ and $\rho(r) = s$.

Sketch of proof: First note that f is the minimal polynomial of r and that \tilde{f} is irreducible so that it is the minimal polynomial of s. Hence

$$K(r) = \{\sum_{j=0}^{n-1} a_j r^j : a_0, ..., a_{n-1} \in K\} \text{ and } \tilde{K}(s) = \{\sum_{j=0}^{n-1} b_j s^j : b_0, ..., b_{n-1} \in \tilde{K}\}$$

Define $\rho : \sum_{j=0}^{n-1} a_j r^j \mapsto \sum_{j=0}^{n-1} \tau(a_j) s^j$. One checks easily that ρ is an isomorphism with the desired properties. Uniqueness is trivial.

Theorem. Let K and \tilde{K} be isomorphic fields and $\tau : K \to \tilde{K}$ the corresponding isomorphism. Suppose f is a polynomial in K[x] and that \tilde{f} is the polynomial corresponding to f under τ . Let F and \tilde{F} be splitting fields of f over K and of \tilde{f} over \tilde{K} , respectively. Then the isomorphism τ maybe extended to an isomorphism $\rho : F \to \tilde{F}$ such that $\rho(r)$ is a root of \tilde{f} for every root r of f.

Sketch of proof: The proof will be by induction over n, the degree of f. Note that there is nothing to prove for n = 1 since F = K and $\tilde{F} = \tilde{K}$. Assume now the truth of the theorem in the case of polynomials of degree n and that f and \tilde{f} are polynomials of degree n + 1. Let φ be an irreducible factor of f and $r \in F$ a root of φ . Then $\tilde{\varphi}$, the polynomial corresponding to φ under τ , is an irreducible factor of \tilde{f} which has a root $s \in \tilde{F}$. By the previous proposition τ can extended to an isomorphism τ_1 from K(r) to $\tilde{K}(s)$ such that $\tau_1(r) = s$. Then $f = (x - r)f_1$ and $\tilde{f} = (x - s)\tilde{f}_1$ in K(r) and $\tilde{K}(s)$, respectively. Note that \tilde{f}_1 corresponds to f_1 under τ_1 and that F and \tilde{F} are splitting fields of f_1 and \tilde{f}_1 , respectively. Since f_1 and \tilde{f}_1 have degree n the isomorphism τ_1 may now be extended to an isomorphism ρ from F to \tilde{F} by the induction hypothesis.

Corollary. All splitting fields of f over K are K-isomorphic.

Sketch of proof: Choose $K = \tilde{K}$ and τ to be the identity. \Box **4.1.14 Normal field extensions.** A finite field extension F over K is called *normal* if F is the splitting field of some polynomial in K[x].

Theorem. If F is a normal extension over K and if the irreducible polynomial $\varphi \in K[x]$ has a root in F then F contains all the roots of φ .

Sketch of proof: Let F be the splitting field of f over K and let \tilde{F} be the splitting field of φ over F. Let s be the root of φ which is in F and let t be any other root of φ . By Proposition 4.1.13 there exists an isomorphism $\tau : K(s) \to K(t)$. Note that F is also the splitting field of f over K(s) and that F(t) is the splitting field of f over K(t). By Theorem 4.1.13 the isomorphism τ extends to an isomorphism from F to F(t). Hence F and F(t) have the same dimension as vector space over K and thus F(t) has dimension one as a vector space over F. But this implies $t \in F$.

4.2. Some more concepts from group theory

4.2.1 Composition series. Let $G = G_0$ be a group and suppose that there are subgroups $G_1, ..., G_{\ell-1}, G_{\ell} = \{1\}$ such that, for $j = 1, ..., \ell$, the subgroup G_j is a normal subgroup of G_{j-1} which is maximal among the proper normal subgroups of G_{j-1} . Then the sequence $(G_0, ..., G_\ell)$ is called a *composition series* of $G = G_0$. The quotient groups G_{j-1}/G_j are called *composition factors* of the series. The number ℓ is called the length of the series.

Theorem Jordan-Hölder theorem. Let S_1 and S_2 be two composition series of a group G. Then S_1 and S_2 have the same length and for each composition factor of S_1 there is an isomorphic composition factor of S_2 .

Example: $(\mathbb{Z}_6, \mathbb{Z}_3, \{0\})$ and $\mathbb{Z}_6, \mathbb{Z}_2, \{0\})$ are composition series.

4.2.2 Solvable Groups. A group G is called solvable if it has a composition series for which each composition factor is abelian.

Cyclic groups of prime order are solvable.

4.2.3 (Un)solvability of the symmetric groups. We denote the symmetric group on n letters by S_n and the associated alternating subgroup by A_n .

 S_2 is a cyclic group of order two and hence solvable.

 A_3 is a cyclic group of order three and hence solvable. A_3 is a normal subgroup of S_3 which is maximal among the proper subgroups of S_3 . The factor S_3/A_3 is isomorphic to \mathbb{Z}_2 and the factor $A_3/\{(1)\}$ is isomorphic to \mathbb{Z}_3 . Hence S_3 is solvable.

Let $V_4 = \{(1), (1, 2)(3, 4), (1, 3)(2, 4), (1, 4)(2, 3)\}$ which is isomorphic to Klein's four group $S_2 \times S_2$ and let $Z_2 = \{(1), (1, 2)(3, 4)\}$ which is isomorphic to S_2 . Then

$$(S_4, A_4, V_4, Z_2, \{(1)\})$$

is a composition series of S_4 all of whose composition factors are abelian.

Next suppose that n > 4 and that S_n is solvable. Then there exists a composition series $(G_0, ..., G_\ell)$ where $G_0 = S_n$ and $G_\ell = \{(1)\}$. We will show by induction that G_s contains all cycles of length three for all $s = 0, ..., \ell$. Since this is obviously true for G_0 assume that it is true for G_s . Let a, b, c, d, e be distinct numbers in $\{1, ..., n\}$. Define $\sigma_1 = (a, b, c)$ and $\sigma_2 = (c, d, e)$ and let ϕ be the canonical homomorphism from G_s to the abelian group G_s/G_{s+1} . Then

$$\phi(\sigma_1^{-1}\sigma_2^{-1}\sigma_1\sigma_2) = 1$$

and hence $(b, e, c) = \sigma_1^{-1} \sigma_2^{-1} \sigma_1 \sigma_2 \in G_{s+1}$. Since b, e, c were arbitrary the induction is completed. Hence we have that the trivial subgroup $\{(1)\}$ contains all cycles of length three which is impossible.

We have proved the following theorem:

Theorem. The symmetric groups S_2 , S_3 , and S_4 are solvable. The symmetric groups S_n are not solvable when n > 4.

4. FIELDS

4.3. Galois Theory

Throughout this section we assume that all fields under consideration have characteristic zero.

4.3.1 Groups of automorphisms. Let F be a field extension of E. Then we define the set

$$E^* = \{ \sigma \in \operatorname{Aut}(F) : \forall u \in E : \sigma(u) = u \}.$$

 E^* is a subgroup of the automorphism group of F. In fact $E^* = \operatorname{Aut}_E(F)$, the group of E-automorphisms of F.

Theorem. Suppose F is a finite (and hence simple) extension of E, i.e., F = E(s). Let φ be the minimal polynomial of s. Then the order of the group E^* equals the number of roots of φ contained in F. In particular, if F is a normal extension of E then $\operatorname{ord}(E^*) = \operatorname{deg}(\varphi) = [F : E]$.

Sketch of proof: Let $\varphi = \sum_{j=0}^{n} \alpha_j x^j$ and $F = \{a_0 + a_1 s + \ldots + a_{n-1} s^{n-1} : a_0, \ldots a_{n-1} \in E\}$. Denote the number of roots of φ in F by m. If $\sigma \in E^*$ then

$$0 = \sigma(\hat{\varphi}(s)) = \sigma(\sum_{j=0}^{n} \alpha_j s^j) = \sum_{j=0}^{n} \alpha_j \sigma(s)^j = \hat{\varphi}(\sigma(s)).$$

Hence $\sigma(s)$ is a root of φ in F. Since the image of s under σ determines σ completely we have $\operatorname{ord}(E^*) \leq m$. Conversely, let t be a root of φ in F. By Proposition 4.1.13 there is an E-isomorphism $\sigma_t : F = E(s) \to F = E(t)$, i.e., an element $\sigma_t \in E^*$, such that $\sigma_t(s) = t$. If t and r are distinct roots of φ then $\sigma_t \neq \sigma_r$. Hence $\operatorname{ord}(E^*) \geq m$.

4.3.2 Fixed fields. Let Σ be a subgroup of Aut(F) and define

$$\Sigma^* = \{ u \in F : \forall \sigma \in \Sigma : \sigma(u) = u \}.$$

Then Σ^* is a subfield of F. It is called the *fixed field* of Σ .

4.3.3 Duality properties. Let F be a field extension of K. Let $\Gamma = \operatorname{Aut}_K(F)$, the group of K-isomorphisms of F. Suppose that L, M are fields such that $K \subset L, M \subset F$ and that Σ and Π are subgroups of Γ . We denote the identity element of Γ by ι . Then the following properties hold:

- (1) $F^* = \{\iota\}, \{\iota\}^* = F$, and $K^* = \Gamma$.
- (2) If $L \subset M$ then $M^* \subset L^*$.
- (3) If $\Sigma \subset \Pi$ then $\Pi^* \subset \Sigma^*$.
- (4) $L \subset L^{**}$ and $\Sigma \subset \Sigma^{**}$.
- (5) $L^* = L^{***}$ and $\Sigma^* = \Sigma^{***}$.

Properties (1) through (4) follow immediately from the definitions. Property (5) follows from (2), (3), and (4). Note that Γ^* is not necessarily equal to K and that the inclusions in (4) can be strict. For example, when $K = \mathbb{Q}$ and $F = \mathbb{Q}(\sqrt[3]{2})$ then Γ contains only the identity. Hence $K^{**} = \Gamma^* = F$.

4.3.4 Closed fields and closed subgroups. Let K, L, F as well as Γ and Σ as before. Then L is called closed if $L = L^{**}$ and Σ is called closed if $\Sigma = \Sigma^{**}$.

Let F be an extension over K. Let \mathcal{E} be a set of all closed fields E such that $K \subset E \subset F$ and let \mathcal{G} be the set of all closed subgroups of $\Gamma = \operatorname{Aut}_K(F)$. Then we define the map $\psi : \mathcal{E} \to \mathcal{G} : E \mapsto E^*$.

Theorem. The map ψ is bijective.

58

4.3.5 The fundamental theorem of Galois theory. If F is a normal field extension over K then all fields E for which $K \subset E \subset F$ are closed and all subgroups Σ of $\operatorname{Aut}_K(F)$ are closed and hence there is a one-to-one correspondence between the intermediate fields E and the subgroups Σ .

Sketch of proof: Let E be an intermediate field and assume the inclusion $E \subset E^{**}$ is strict. Choose $s \in E^{**} - E$ and denote the minimal polynomial of s over E by φ . Since $s \notin E$ we know that $\deg(\varphi) > 1$. Since F is normal over K it is normal over E and hence F contains all the roots of φ . Let t be a root of φ different from s. Then, by Proposition 4.1.13, the identity map from E to E extends to an isomorphism between E(s) and E(t) and this isomorphism extends to an element σ of E^* . Hence $\sigma(s) = t$. However, since $s \in E^{**}$ we also have $\sigma(s) = s$ and hence s = t, a contradiction. Hence $E = E^{**}$.

Now let $\Sigma = \{\sigma_1, ..., \sigma_h\}$ be a subgroup of $\operatorname{Aut}_K(F) = K^*$ of order h. There is an element $u \in F$ such that F = K(u) and hence $F = \Sigma^*(u)$. Denote the minimal polynomial of u over Σ^* by φ and define the polynomial $f = (x - \sigma_1(u))...(x - \sigma_h(u)) = \sum_{j=0}^h a_j x^j$. The coefficients a_j are symmetric polynomials of the roots $\sigma_1(u), ..., \sigma_h(u)$, in particular, $a_h = 1, a_{h-1} = -\sum_{j=0}^h \sigma_j(u)$, and $a_0 = (-1)^h \prod_{j=0}^h \sigma_j(u)$. This implies that $\sigma_k(a_j) = a_j$ whence $f \in \Sigma^*[x]$. Since u is a root of f we have $\deg(\varphi) \leq \deg(f)$. Since F is normal over Σ^* Theorem 4.3.1 gives that $\operatorname{ord}(\Sigma^{**}) = [F : \Sigma^*] = \deg(\varphi)$ and hence $\operatorname{ord}(\Sigma^{**}) \leq \operatorname{ord}(\Sigma)$. This shows finally that $\Sigma = \Sigma^{**}$.

4.3.6 The Galois group of a field extension. If F is a normal field extension over K the set $\operatorname{Aut}_K(F)$ of all K-automorphism of F is called the *Galois group* of F over K and it is denoted by $\operatorname{Gal}(F/K)$.

4.3.7 Normal extension and normal subgroups. Suppose F is a normal field extension over K and E is an intermediate field. Then E is a normal extension over K if and only if E^* is a normal subgroup of $K^* = \operatorname{Aut}_K(F)$. Moreover, if E is a normal extension over K, then the Galois group $\operatorname{Gal}(E/K)$ is isomorphic to the quotient group K^*/E^* .

Sketch of proof: Assume that E = K(u) and let f be the minimal polynomial of u over K.

Suppose E is a normal extension over K. If $\gamma \in K^*$ then $\gamma(u)$ is a root of f and hence in E. Therefore, if $\sigma \in E^*$ and $s \in E$, then $s = \sum_{i=0}^{n-1} k_j u^j$ with $k_j \in K$ and hence

$$(\gamma \circ \sigma \circ \gamma^{-1})(s) = \sum_{j=0}^{n-1} k_j (\gamma \circ \sigma \circ \gamma^{-1})(u)^j = \sum_{j=0}^{n-1} k_j u^j = s.$$

This shows that $\gamma \circ \sigma \circ \gamma^{-1} \in E^*$, i.e., that E^* is a normal subgroup of K^* .

Now assume that E^* is a normal subgroup of K^* . Then $\sigma(\gamma(u)) = (\sigma \circ \gamma)(u) = \gamma(u)$ for all $\gamma \in K^*$ and all $\sigma \in E^*$. This implies $\gamma(u) \in E^{**} = E$ and hence that all roots of f are contained in E, i.e., E the splitting field of f.

Finally, define the homomorphism $\Xi : K^* \to \text{Gal}(E/K) : \sigma \mapsto \sigma|_E$. This homomorphism is surjective by Theorem 4.1.13 and its kernel is E^* . The last statement of the theorem follows now from the fundamental isomorphism theorem for commutative rings.

4.4. Radical Extensions

4.4.1 Primitive roots of unity. Let K be a field and let n be a natural number which is not a multiple of the characteristic p of K (any natural number if K has characteristic zero). Let f be the polynomial $x^n - 1$ in K[x]. If \hat{f} has a root $\zeta \in K$ then ζ is called an n-th root of unity. Note that \hat{f} has at most n roots and let F be the splitting field of f. The derivative of f is nx^{n-1} and is different from zero. In this case \hat{f} has therefore n distinct

n-th roots of unity in F. The set of all n-th roots of unity forms an abelian group under multiplication. One can in fact show that this group is cyclic, i.e., it is generated by one of its elements. Any such element is called a primitive n-th root of unity.

4.4.2 Radical extension. A simple field extension K(s) over K is called a *simple radical* extension over K if there exists a natural number n > 1 such that s^n in K, or, in other words, if the minimal polynomial of s over K is $x^n - r$ for some $r \in K$. The adjunction of s to K is then called *radical*.

A field extension F over K is called a *radical extension* over K if there exist elements t_1 , ..., t_m in F such that $F = K(t_1)(t_2)...(t_m)$ and each of the successive adjunctions is radical.

Theorem. For any radical extension E over K there is an extension F over E and an ascending tower of fields

$$K = F_0 \subset F_1 \subset \ldots \subset F_k = F$$

with the following properties:

- (1) For each $j \in \{1, ..., k\}$ there exists a prime number p_j and an element $s_j \in F_j$ such that $r_j = s_j^{p_j} \in F_{j-1}$ and $F_j = F_{j-1}(s_j)$.
- (2) F_j is a normal extension over F_{j-1} .

Sketch of proof: By the definition of a radical extension we have a tower of radical extension fields

$$K = F_0 \subset F_1 \subset \ldots \subset F_k = E$$

where $F_j = F_{j-1}(s_j)$. If $s_j^n \in F_{j-1}$ and *n* has the prime factorization $n = p_1...p_{\nu}$ let $t_{j,\ell} = s_j^{n/(p_1...p_{\ell})}$. Then replace the string $F_{j-1} \subset F_j$ in the above tower of fields by

 $F_{j-1} \subset F_{j-1}(t_{j,1}) \subset \ldots \subset F_{j-1}(t_{j,1})...(t_{j,\nu}) = F_j.$

Doing this for all j we obtain a tower of radical extensions having the first of the above properties (with F = E).

Suppose F_{j-1} is a normal extension over K. If F_{j-1} contains ζ_j , a primitive p_j -th root of unity, then F_j is a normal extension over F_{j-1} and hence a normal extension over K. If F_{j-1} does not contain ζ_j , a primitive p_j -th root of unity, replace the string $F_{j-1} \subset F_j$ in the tower by $F_{j-1} \subset F_{j-1}(\zeta_j) \subset F_j$. Then $F_{j-1}(\zeta_j)$ is the splitting field of the polynomial $x^{p_j}-1$ and hence normal. Also F_j contains, together with s_j , all other roots $s_j\zeta_j^k$, $k = 1, ..., p_j$ of the polynomial $x^{p_j} - r_j$ and hence is a splitting field. Repeating this for all j results in a tower having both of the required properties. Note that at step j the root ζ_j may not be an element of E so that the resulting field F may be an extension of E.

Example: Let $K = \mathbb{Q}$ and $E = \mathbb{Q}(\sqrt[3]{2})$. Then $F_1 = \mathbb{Q}(\omega)$ where $\omega = (-1 + i\sqrt{3})/2$ is a third root of unity. Also $F = F_2 = \mathbb{Q}(\omega, \sqrt[3]{2})$ is the splitting field of $x^3 - 2$.

4.4.3 Radical extension and solvability. If F is a normal radical extension K then the Galois group of F over K is solvable.

Sketch of proof: There exists a tower of fields

$$K = F_0 \subset F_1 \subset \ldots \subset F_k = F$$

with the properties of Theorem 4.4.2. By the fundamental theorem of Galois theory there is a tower of subgroups

$$\{\iota\}=F^*\subset F^*_{k-1}\subset\ldots\subset F^*_1\subset K^*=\mathrm{Gal}(F/K).$$

By Theorem 4.3.7 we have that F_j^* is a normal subgroup of F_{j-1}^* . The subgroup F_j^* is maximal among the normal subgroups of F_{j-1}^* since its index in F_{j-1}^* equals $[F_j:F_{j-1}] = p_j$

which is prime. Hence F_j^* is a maximal normal subgroup of F_{j-1}^* and our theorem is proved once we show that $F_{j-1}^*/F_j^* = \operatorname{Gal}(F_j/F_{j-1})$ is abelian. Suppose $\sigma_1, \sigma_2 \in F_{j-1}^*$. Since $F_j = F_{j-1}(s_j)$ we have that σ_1 and σ_2 are uniquely

Suppose $\sigma_1, \sigma_2 \in F_{j-1}^*$. Since $F_j = F_{j-1}(s_j)$ we have that σ_1 and σ_2 are uniquely determined respectively by the values $\sigma_1(s_j) = \zeta_j^{k_1} s_j$ and $\sigma_2(s_j) = \zeta_j^{k_2} s_j$ where ζ_j is a primitive p_j -th root of unity and k_1 and k_2 are appropriate integers. Hence

$$(\sigma_1 \circ \sigma_2)(s_j) = \sigma_1(\zeta_j^{k_2} s_j) = \sigma_1(\zeta_j)^{k_2} \zeta_j^{k_1} s_j \quad \text{and} \quad (\sigma_2 \circ \sigma_1)(s_j) = \sigma_2(\zeta_j^{k_1} s_j) = \sigma_2(\zeta_j)^{k_1} \zeta_j^{k_2} s_j.$$
 If $s_j^{p_j} = 1$ we may assume $s_j = \zeta_j$ whence

$$(\sigma_1 \circ \sigma_2)(s_j) = s_i^{(k_1+1)k_2+k_1+1} = (\sigma_2 \circ \sigma_1)(s_j).$$

If $s_j^{p_j} \neq 1$ we have $\zeta_j \in F_{j-1}$ and hence $\sigma_1(\zeta_j) = \sigma_2(\zeta_j) = \zeta_j$. Thus

$$(\sigma_1 \circ \sigma_2)(s_j) = \zeta_j^{k_1 + k_2} s_j = (\sigma_2 \circ \sigma_1)(s_j)$$

4.5. The Theorem of Ruffini and Abel

4.5.1 Historical facts. The quadratic formula was known (in a sense) to the Babylonians perhaps 5000 years ago. The cubic equation and the quartic equation were solved by radicals in the 1500s by the Italian mathematicians del Ferro, Tartaglia and Cardan and Ferrari.

By 1799 over 250 years had passed without anyone being able to solve the quintic equation by radicals even though attempts hade been made by many mathematicians including very famous people like Euler, Bézout, Vandermonde, and Lagrange.

Then in 1799 Ruffini proved (perhaps not entirely correctly) that this task was in fact impossible but his assertion did not enter the consciousness of the mathematical community of the time. The reason was perhaps that nobody really believed that it was impossible. A quarter century later Abel gave another proof of this fact and this time the message stuck. **4.5.2 The general polynomial equation of degree** n. Let k be a field and $u_0, ..., u_{n-1}$ indeterminates. Then

$$f = x^{n} + u_{n-1}x^{n-1} + \dots + u_0$$

is a polynomial over the field $K = k(u_0, ..., u_{n-1})$. The equation

$$\hat{f}(z) = 0$$

is called the general polynomial equation of degree n.

4.5.3 The concept of solvability by radicals. To solve a polynomial equation by radicals means finding a formula for its roots in terms of the coefficients so that the formula only involves the operations of addition, subtraction, multiplication, division and taking roots, each a finite number of times.

Definition. Let $f \in K[x]$. The polynomial equation $\hat{f}(z) = 0$ is called solvable by radicals if the splitting field of f is a radical extension over K.

For instance, the general quadratic equation in $K = \mathbb{C}(u_0, u_1)$ is $z^2 + u_1 z + u_0 = 0$ and the corresponding equation is solved by the formula

$$z_{1,2} = -\frac{u_1}{2} \pm \frac{1}{2}\sqrt{u_1^2 - 4u_0}$$

in the extension field $K(\sqrt{\Delta})$ where $\Delta = u_1^2 - 4u_0$. The purpose of the formula is, of course, that $z_{1,2}$ are the roots of $z^2 + u_1 z + u_0$ (in $K(\sqrt{\Delta})$) for any choice of u_1 and u_0 in \mathbb{C} or even K.

4. FIELDS

However, if $K = \mathbb{C}$, i.e., $f \in \mathbb{C}[x]$, then the splitting field of f is also \mathbb{C} and a radical extension of itself. So, by our definition, $\hat{f}(z) = 0$ is solvable, but this does not mean there are effective means to actually do it. For this reason we consider u_0, \ldots, u_{n-1} as indeterminates.

4.5.4 The theorem of Ruffini and Abel. Let K be a field of characteristic zero and let $f = x^n + u_{n-1}x^{n-1} + ... + u_0$ be a polynomial in $K(u_0, ..., u_{n-1})[x]$. Then $\hat{f}(z) = 0$ is not solvable by radicals if n > 4.

Sketch of proof: Let $x_1, ..., x_n$ be *n* indeterminates and let $v_0, ..., v_{n-1}$ be the elementary symmetric polynomials of these indeterminates, i.e.,

$$v_0 = (-1)^n x_1 \dots x_n \quad \dots \quad v_{n-1} = -x_1 - \dots - x_n.$$

Then $K(x_1, ..., x_n)$ is the splitting field of $g = x^n + v_{n-1}x^{n-1}... + v_0$ over $K(v_0, ..., v_{n-1})$. The Galois group of $K(x_1, ..., x_n)$ over $K(v_0, ..., v_{n-1})$ is the group of permutations of the indeterminates $x_1, ..., x_n$.

Now suppose that f(z) = 0 is solvable by radicals and let F be the splitting field of f which, by assumption, is a radical extension of $K(u_0, ..., u_{n-1})$. One may show that the rings $K[u_0, ..., u_{n-1}]$ and $K[v_0, ..., v_{n-1}]$ are isomorphic and this isomorphism extends to an isomorphism between the fraction fields $K(u_0, ..., u_{n-1})$ and $K(v_0, ..., v_{n-1})$ and, by Theorem 4.1.13, to an isomorphism between F and $K(x_1, ..., x_n)$ and therefore the Galois group of F over $K(u_0, ..., u_{n-1})$ is isomorphic to the permutation group on n letters which is not solvable by Theorem 4.2.3.

CHAPTER 5

Vector Spaces

5.1. Fundamentals

5.1.1 Vector spaces. Let V be a set and let K be a field. Suppose there is a binary operation on V (denoted by +) and a function σ from $K \times V$ to V (denoted by juxtaposition) such that the following properties are satisfied:

(a) (V, +) is an abelian group,

(b) (rs)x = r(sx) for all $r, s \in K$ and all $x \in V$,

(c) (r+s)x = rx + sx for all $r, s \in K$ and all $x \in V$,

(d) r(x+y) = rx + ry for all $r \in K$ and all $x, y \in V$,

(e) 1x = x for all $x \in V$.

Then $(V, K, +, \sigma)$ (or just V) is called a vector space over K. If $K = \mathbb{R}$ we call V a real vector space and if $K = \mathbb{C}$ we call V a complex vector space. The elements of V are called vectors and the elements of K are called scalars. The map $(r, x) \mapsto rx$ is called scalar multiplication. The identity element of the group (V, +), called the zero vector, is denoted by 0. A confusion with the scalar 0 can not arise.

For all $x \in V$ and all $r \in K$ we have rx = 0 if and only if r = 0 or x = 0.

5.1.2 Examples. $\mathbb{R}^2 = \mathbb{R} \times \mathbb{R}$, \mathbb{R}^3 and, more generally, \mathbb{R}^n are real vector spaces under the usual componentwise addition and scalar multiplication. \mathbb{C} and, more generally, \mathbb{C}^n can be thought of as a real or a complex vector space. The set $\mathbb{R}[x]$ of all polynomials in xwith real coefficients and the set $\mathbb{R}[x]_n$ of all polynomials of degree at most n in x with real coefficients are real vector spaces. Likewise $\mathbb{C}[x]$ and $\mathbb{C}[x]_n$ are complex vector spaces. The set of all functions from an interval $I \subset \mathbb{R}$ to \mathbb{R} which are k times continuously differentiable (denoted by $C^k(I)$) is a vector space. The set $\{0\}$ over any field is a vector space, the null space or trivial vector space.

5.1.3 Linear combinations. If $x_1, ..., x_n$ are elements of a vector space V and if $\alpha_1, ..., \alpha_n$ are scalars then the vector

$$\alpha_1 x_1 + \ldots + \alpha_n x_n$$

is called a *linear combination* of $x_1, ..., x_n$.

5.1.4 Subspaces and spans. A subset of a vector space V which is itself a vector space (with respect to the operations in V) is called a *subspace* of V.

A nonempty subset S of a vector space V is a subspace of V if and only if $\alpha x + \beta y \in S$ whenever $x, y \in S$ and $\alpha, \beta \in K$, the scalar field.

The intersection of a nonempty collection of subspaces of V is again a subspace of V.

Let A be a subset of V. Let C be the collection of all subspaces of V which include A. Then the set $\langle A \rangle = \bigcap_{U \in C} U$ is a subspace of V called the span of A. We also say that a vector space W is spanned by A or that A spans (or that the elements of A span) W if $W = \langle A \rangle$.

If $A = \{\}$ then $\langle A \rangle = \{0\}$. Otherwise $\langle A \rangle$ is the set of all linear combinations of elements of A.

Example: The vectors (3,5) and (0,-2) span \mathbb{R}^2 .

5.1.5 Dimension of a vector space. A vector space V is called *finite-dimensional* if it is equal to $\{0\}$ or if there exists a finite subset of V which spans V. Otherwise V is called *infinite-dimensional*.

Let $V \neq \{0\}$ be a finite-dimensional vector space and S the (nonempty) set of natural numbers n for which there is a subset of V with n elements that spans V. Then S contains a smallest number which is called the *dimension* of V. The dimension of V is denoted by dim V.

The null space $\{0\}$ is said to have dimension zero.

5.1.6 Linear independence. Let V be a vector space. The vectors $x_1, ..., x_n \in V$ are called *linearly independent* if $\alpha_1 x_1 + ... + \alpha_n x_n = 0$ implies that $\alpha_1 = ... = \alpha_n = 0$. Otherwise they are called linearly dependent. A set $M \subset V$ is called linearly independent if any finite number of distinct elements of M are linearly independent. Otherwise M is called linearly dependent. In particular, the empty set is linearly independent and a set consisting of precisely one element is linearly independent if and only if that element is not the zero element. Moreover, any set containing the zero element is linearly dependent. If $A \subset B$ and B is linearly independent then so is A.

The vectors $x_1, ..., x_n$ are linearly dependent if and only if one of them can be expressed as a linear combination of the others.

5.1.7 Bases. A set $B \subset V$ is called a *basis* of the nontrivial vector space V if it is linearly independent and spans V. This is equivalent to the statement that every $x \in V$ can be expressed uniquely as a linear combination of the elements of B. The empty set is the basis of the trivial vector space.

Theorem. Let M be a linearly independent subset of a nontrivial vector space V. Then there exists a basis B of V such that $M \subset B$. In particular, every nontrivial vector space has a basis.

Sketch of proof: Let Σ be the set $\Sigma = \{A \subset V : M \subset A \land A \text{ linearly independent}\}$. This is a nonempty set which is partially ordered by set inclusion. For any nonempty totally ordered subset Γ of Σ define $A' = \bigcup_{T \in \Gamma} T$. We will show that $A' \in \Sigma$ and hence that Γ has an upper bound in Σ . Therefore, by Zorn's Lemma, Σ has a maximal element, which we call B. If we also show that B spans V, the theorem is proven.

It is obvious that $M \subset A'$. Let $a_1, ..., a_n$ be distinct elements of A'. Then there exist $T_1, ..., T_n \in \Gamma$ such that, for j = 1, ..., n we have that $a_j \in T_j$. Since Γ is totally ordered there is, in fact, a $k \in \{1, ..., n\}$ such that $T_1, ..., T_n \subset T_k$ and hence $a_1, ..., a_n \in T_k$. Since T_k is a linearly independent set the vectors $a_1, ..., a_n$ are linearly independent. Since they were arbitrary, A' is linearly independent and hence $A' \in \Sigma$.

Finally, let v be any element of V and suppose that v is not in the span of B. Then $B \cup \{v\}$ is a linearly independent subset of V which contains M, i.e., $B \cup \{v\} \in \Sigma$. But then B is not maximal which contradicts its definition. Hence no such v can exist, i.e., B spans V.

Obviously, the basis of an infinite-dimensional vector space must have infinitely many elements.

In the case of a finite-dimensional vector space 5.3.1 and 5.3.2 can be combined to give another (constructive) proof of the above theorem. It will also be shown in Corollary 5.3.2 that the number of elements of any basis of a finite-dimensional vector space is equal to the dimension of that vector space.

5.1.8 Direct sums. Let V and W be two vector spaces over the field K. We define the direct sum $V \oplus W$ of V and W to be the vector space $(V \times W, K, +, \sigma)$ where + is defined

by $(v_1, w_1) + (v_2, w_2) = (v_1 + v_2, w_1 + w_2)$ and $\sigma : K \times (V \times W) \to V \times W$ is defined by $(r, (v, w)) \mapsto r(v, w) = (rv, rw)$.

Theorem. The dimension of a direct sum satisfies $\dim(V \oplus W) = \dim V + \dim W$.

Sketch of proof: Let B_1 be a basis of V and B_2 a basis of W. Define $A_1 = \{(v,0) : v \in B_1\}$, $A_2 = \{(0,w) : w \in B_2\}$ and $B = A_1 \cup A_2$. Then B is a basis of $V \oplus W$. \Box Example: Let $V = \{(a,a) : a \in \mathbb{R}\}$ and $W = \{(a,b) : a,b \in \mathbb{R}\}$. Then $V \oplus W = \{((a,a), (b,c)) : a, b, c \in \mathbb{R}\}$. A basis is $\{((1,1), (0,0)), ((0,0), (1,0)), ((0,0), (0,1))\}$.

5.1.9 Internal sums and internal direct sums. Let X and Y be two subspaces of a vector space V. The union of X and Y is not necessarily a subspace of V. We define the sum of X and Y to be the subspace generated by their union, i.e., $X + Y = \langle X \cup Y \rangle$. It turns out that $X + Y = \{x + y : x \in X, y \in Y\}$. The dimension of X + Y is infinite if at least one of X and Y has infinite dimension. Otherwise it is given by

$$\dim(X+Y) = \dim X + \dim Y - \dim(X \cap Y).$$

If $X \cap Y = \{0\}$ we denote X + Y by $X \oplus Y$ and call it the *internal direct sum* of X and Y.¹

Theorem. Let X be a subspace of a vector space V. Then there exists a subspace Y such that $X \cap Y = \{0\}$ and $X \oplus Y = V$.

Sketch of proof: Let A be a basis of X. By Theorem 5.1.7 there is a basis C of V such that $A \subset C$. Define B = C - A and $Y = \langle B \rangle$. Then Y is a subspace of V. Assume $x \in X \cap Y$. Then x can be written as a linear combination of elements of A and as a linear combination of elements of B. Subtracting these two expressions gives that a linear combination of elements of C equals zero. Hence all coefficients are zero and x is equal to zero. Thus $X \cap Y = \{0\}$. Since $A \cup B$ spans V we get that X + Y = V.

5.2. Linear Transformations

5.2.1 Linear transformations.

Let V and W be two vector spaces over the same field K. A function $f: V \to W$ is called a *linear transformation* or a vector space homomorphism if

$$f(\alpha x + \beta y) = \alpha f(x) + \beta f(y)$$

for all $\alpha, \beta \in K$ and all $x, y \in V$.

A bijective linear transformation from V to W is called a (vector space) *isomorphism*. If W = V it is called a (vector space) *automorphism*. Two vector spaces V and W are called *isomorphic* if there exists an isomorphism from V to W.

The following properties hold for linear transformations $f: V \to W$:

- f(0) = 0, f(-x) = -f(x), since a linear transformation is a group homomorphism from (V, +) to (W, +).
- The composition of linear transformations is a linear transformation.
- The relation "is isomorphic to" is an equivalence relation on the set of all vector spaces over K.
- -f(V) is a subspace of W.
- $\ker f = \{x \in V : f(x) = 0\} \text{ is a subspace of } V.$
- f is injective if and only if ker $f = \{0\}$.

¹The internal direct sum of X and Y is isomorphic (see 5.2.1) to their (external) direct sum. This justifies using the same notation.

5. VECTOR SPACES

- The automorphisms of a vector space V form a group under composition called the linear group of V and denoted by GL(V).
- The functions from V or W to $V \oplus W$ which maps $v \in V$ to (v, 0) or $w \in W$ to (0, w) are injective linear transformations called embeddings.
- The functions from $V \oplus W$ to V or W which maps (v, w) to v or w are surjective linear transformations called projections.

Theorem. Suppose V and W are vector spaces and B is a basis of V. Then any function from B to W extends uniquely to a linear transformation from V to W. In particular, any linear transformation is uniquely determined by the images of a basis of the domain of the transformation.

Sketch of proof: Denote the given function from B to W by f. Define g on V by $g(\alpha_1 x_1 + ... + \alpha_n x_n) = \alpha_1 f(x_1) + ... + \alpha_n f(x_n)$ where $\{x_1, ..., x_n\} \subset B$. Then g is a linear transformation from V to W. It is the only linear transformation whose restriction to B is equal to f. \Box

Examples: Let A be an $n \times m$ matrix. Then $x \mapsto Ax$ is a linear transformation from \mathbb{R}^m to \mathbb{R}^n . If $q : \mathbb{R} \to \mathbb{R}$ is a continuous function then $-d^2/dx^2 + q$ which maps $y \in C^2(\mathbb{R})$ to $-y'' + qy \in C^0(\mathbb{R})$ is a linear transformation.

5.2.2 Dimensions of images and kernels. The dimensions of the kernel and the image of a linear transformation are not independent as the following theorem shows. This theorem is sometimes called the fundamental theorem of linear algebra.

Theorem. Let $f: V \to W$ be a linear transformation. Then dim f(V) + dim ker $f = \dim V$.

Sketch of proof: If dim ker $f = \infty$ the theorem becomes trivial. Hence assume that dim ker f = k and that $B = \{b_1, ..., b_k\}$ is a basis of ker f. By Theorem 5.1.9 there exists a subspace Y of V such that ker $f \oplus Y = V$. Let C be a basis of Y. We show below that f(C) is a basis of f(V). If V is infinite-dimensional then C and f(C) are infinite sets and hence dim f(V) is infinite. If dim $V = n < \infty$ then C has n - k elements and so does f(C). Hence dim f(V) = n - k proving the theorem, provided f(C) is a linearly independent spanning subset of f(V).

In the following c_k denotes an element of C whenever $k \in \mathbb{N}$ and $w_k = f(c_k)$. Consider the equation $\alpha_1 w_1 + \ldots + \alpha_j w_j = 0$. Then $x = \alpha_1 c_1 + \ldots + \alpha_j c_j \in \ker f$, i.e., $x \in Y \cap \ker f$ and hence x = 0. This shows that all coefficients α_j are zero and hence that $\{w_1, \ldots, w_j\}$ is a linearly independent set.

Let $w \in f(V)$. Then w = f(v) for some $v \in V$. Hence v = x + y where $x \in \ker f$ and $y \in Y$. But this shows that f(v) = f(y) since f(x) = 0. Thus f(Y) = f(V) but f(C) spans f(Y).

Corollary. Let f be a linear transformation between vector spaces V and W and B a basis of V. Then

- (a) dim $f(V) \leq \dim V$ and dim $f(V) \leq \dim W$,
- (b) f is injective if and only if $f|_B$ is injective and f(B) is linearly independent,
- (c) f is surjective if and only if f(B) spans W.
- If V and W are finite-dimensional we have
- (d) f is injective if and only if dim $f(V) = \dim V$,
- (e) f is surjective if and only if dim $f(V) = \dim W$.

The dimension of f(V), i.e., dim f(V) is called the rank of f denoted by rank f.

5.2.3 A few minor facts. Throughout this section V is a finite-dimensional vector space. We will prove a few statements which will be needed later:

66

A. Let A be a linear transformations from V to V. If ker $A^2 = \ker A$ then $V = \ker A \oplus A(V)$.

Sketch of proof: Note that $A(V) \cap \ker A = \{Ax : x \in \ker A^2\}$. The latter set, however, equals $\{0\}$ since $\ker A^2 = \ker A$. By 5.1.9 and the fundamental theorem of linear algebra we obtain then that $\dim(A(V) \oplus \ker A) = \dim V$. Since V is finite-dimensional this implies that $A(V) \cup \ker(A)$ spans V.

B. Suppose A_1 and A_2 are commutative linear transformations from V to V, that $\ker A_1 \cap \ker A_2$ is trivial, and that $\ker A_1^2 = \ker A_1$. Then $\ker(A_1A_2) = \ker A_1 \oplus \ker A_2$.

Sketch of proof: Assume x = u + v where $u \in \ker A_1$ and $v \in \ker A_2$. Then $A_1A_2x = A_2A_1u + A_1A_2v = 0$, i.e., $x \in \ker(A_1A_2)$. Next suppose that $x \in \ker(A_1A_2)$. By part A we know that $x = u + A_1v$ for some $v \in V$ and some $u \in \ker A_1$. Hence $A_1^2A_2v = 0$. Since $\ker A_1^2 = \ker A_1$ we have also $0 = A_1A_2v = A_2A_1v$. Hence $A_1v \in \ker A_2$.

C. Assume that $A_1, ..., A_m$ are pairwise commutative linear transformations from V to V such that ker $A_j^2 = \ker A_j$ for j = 1, ..., m and ker $A_j \cap \ker A_\ell = \{0\}$ for $j \neq \ell$. Then $\ker(A_1...A_m) = \ker A_1 \oplus ... \oplus \ker A_m$.

Sketch of proof: First show (by induction) that $\ker A_1 \cap \ker(A_{m+1-k}...A_m) = \{0\}$ for k = 1, ..., m-1. Hence, by part B, $\ker(A_1...A_m) = \ker A_1 \oplus \ker(A_2...A_m)$. Another induction completes the proof.

5.2.4 Quotient spaces. Let $(V, K, +, \sigma)$ be a vector space and U a subspace of V. Recall that (V/U, +) is a commutative group. One may define a scalar multiplication $K \times V/U \rightarrow V/U$ by $(\alpha, x + U) \mapsto \alpha x + U$. This scalar multiplication (also denoted by σ is well defined and turns $(V/U, K, +, \sigma)$ into a vector space, called the quotient space of V with respect to U.

Let $\varphi: V \to V/U$ be the canonical group homomorphism which is indeed a vector space homomorphism, i.e., a linear map. Since $U = \ker(\varphi)$ and since φ is surjective we obtain from Theorem 5.2.2

$$\dim(V/U) = \dim V - \dim U$$

if U is finite-dimensional.

5.2.5 Invariant subspaces. Let $f: V \to V$ be a linear transformation on a vector space V and U a subspace of V. If $f(U) \subset U$ then U is called an *invariant* subspace with respect to f. For example, the subspace $C^{\infty}(\mathbb{R})$ of $C^{1}(\mathbb{R})$ is an invariant subspace with respect to the linear transformation $y \mapsto y'$.

5.3. Finite-dimensional vector spaces

5.3.1 Existence of a basis of a finite-dimensional vector space. Suppose that V is an *n*-dimensional vector space and that A is a subset of V with *n* elements. If A spans V, then it is a basis of V. In particular, any *n*-dimensional vector space has a basis consisting of *n* elements.

Sketch of proof: Let V be a vector space of dimension n. Then there exists a set $\{x_1, ..., x_n\}$ which spans V. Assume that this set is not linearly independent. Then one of its elements, say x_n , is a linear combination of the others. This fact is used to show that any $x \in V$ can be represented by a linear combination of $x_1, ..., x_{n-1}$. Hence V is in fact spanned by $\{x_1, ..., x_{n-1}\}$ but this is a contradiction to the fact that V has dimension n. \Box **5.3.2 Exchange of basis elements.** Suppose V is a vector space of dimension n > 0 and that $B = \{b_1, ..., b_n\}$ spans V (and hence is a basis of V). If A is a linearly independent subset of V with $k \leq n$ elements such that $A \cap B$ is empty, then there exists a basis of V which contains A and n - k elements of B.

5. VECTOR SPACES

Sketch of proof: If S is a set we will use #S to denote the number of elements of S. We say that A has property E if there exists a subset B' of B such that $A \cup B'$ is a basis of V and #A + #B' = n. Now define

$$\begin{split} M &= \{ j \in \mathbb{N} : j > n+1 \lor [\forall A \subset V : (A \text{ linearly independent} \land \#A = j-1) \\ &\to A \text{ has property } E] \}. \end{split}$$

Then $1 \in M$ since the empty set has property E (choose B' = B). Next suppose that $j \in M, j \leq n$, and $A = \{a_1, ..., a_j\}$. Then $\{a_1, ..., a_{j-1}\}$ has property E and hence, after a proper relabeling of the b_ℓ , the set $\{a_1, ..., a_{j-1}, b_j, ..., b_n\}$ is a basis of V. Hence there exist scalars $\alpha_1, ..., \alpha_{j-1}$ and $\beta_j, ..., \beta_n$ such that

(3)
$$a_j = \sum_{\ell=1}^{j-1} \alpha_\ell a_\ell + \sum_{\ell=j}^n \beta_\ell b_\ell$$

where at least one of the β_{ℓ} , say β_m , must be different from zero. Hence equation (3) can be solved for b_m . Now one may show that b_m can be replaced by a_j as a basis vector. But this means that A has property E and hence that $j + 1 \in M$. The induction principle proves now that $M = \mathbb{N}$ and thus the theorem.

Corollary. Suppose that V is an n-dimensional vector space.

(a) If A is a linearly independent subset of V with n elements, then it is a basis of V.

(b) Every basis of V has precisely n elements.

5.3.3 Arithmetic vector spaces. If K is a field then K^n is an n-dimensional vector space over K. It is customary to consider K^n as a space of columns with n components. Subsequently a column will be denoted by $(\alpha_1, ..., \alpha_n)^{\top}$. For j = 1, ..., n let e_j denote the column all of whose entries are zero except for the j-th entry which is one, i.e., $e_1 = (1, 0, ... 0)^{\top}$ etc. Then $\{e_1, ..., e_n\}$ is a basis of K^n called the canonical basis. Note that $(x_1, ..., x_n)^{\top} = \sum_{j=1}^n x_j e_j$.

5.3.4 Coordinates. Let V be a vector space of dimension n and $\{b_1, ..., b_n\}$ a basis of V. Then we call the tuple $B = (b_1, ..., b_n) \in V^n$ an ordered basis of V.

Now let x be an element of V. Then the (uniquely determined) coefficients $\alpha_1, ..., \alpha_n$ in the representation $x = \sum_{j=1}^n \alpha_j b_j$ are called *coordinates* of x with respect to the ordered basis B. The tuple $(\alpha_1, ..., \alpha_n)^{\top}$ of coordinates is denoted by x_B .

If $V = K^n$ we have to distinguish between an element x of V which is denoted by $(x_1, ..., x_n)^{\top}$ and the tuple of coefficients of x with respect to an ordered basis B, which we will denote by $[\beta_1, ..., \beta_n]_B^{\top}$. In particular, if B is the canonical basis, then $x = (x_1, ..., x_n)^{\top} = [x_1, ..., x_n]_B^{\top}$.

The concept of coordinates allows to exhibit a very close relationship between vector spaces of the same finite dimension:

Theorem. An *n*-dimensional vector space V over the field K is isomorphic to K^n .

Sketch of proof: Choose an ordered basis B for V. Then the map

$$g: V \to K^n: \sum_{j=1}^n \alpha_j b_j \mapsto (\alpha_1, ..., \alpha_n)^\top$$

is an isomorphism.

5.3.5 Matrix representations of a linear transformations. Let f be a linear transformation from K^n to K^m . Choose ordered bases $A = (a_1, ..., a_n)$ and $B = (b_1, ..., b_m)$ for K^n to K^m , respectively, (e.g., the canonical bases). Form an $m \times n$ matrix M (a matrix

68

with *m* rows and *n* columns) by letting the *j*-th column of *M* be the tuple of coefficients of the vector $f(a_j)$ with respect to the basis *B*, i.e., if $f(a_j) = \sum_{k=1}^{m} \beta_{k,j} b_k$ then

$$M = \begin{pmatrix} \beta_{1,1} & \dots & \beta_{1,n} \\ \vdots & & \vdots \\ \beta_{m,1} & \dots & \beta_{m,n} \end{pmatrix}.$$

Then computing f(x) is reduced to matrix multiplication, since we have for any $x \in K^n$

$$f(x) = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}_B = M \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_A.$$

Conversely, every $m \times n$ matrix represents a linear transformation from K^n to K^m after choosing ordered bases in both domain and range.

More generally, if V and W are vector spaces over K of dimensions n and m, respectively, then any linear transformation f from V to W is represented by a matrix with entries in K and every matrix of suitable dimensions represents a linear transformation after ordered bases $(v_1, ..., v_n)$ and $(w_1, ..., w_m)$ of V and W are chosen. Specifically, given these bases define isomorphisms $g: V \to K^n$ and $h: W \to K^m$ as in 5.3.4. Then, according to the above, the linear transformation $h \circ f \circ g^{-1}: K^n \to K^m$ may be represented by an $m \times n$ matrix M. Then the j-th column of M is $h(f(g^{-1}(e_j))) = h(f(v_j)) = [f(v_j)]_{(w_1,...,w_m)}$, i.e., the m-tuple of coefficients of $f(v_j)$ when written in terms of the ordered basis $(w_1, ..., w_m)$. **5.3.6 Systems of linear equations.** Consider a system of m linear equations in n unknowns $x_1, ..., x_n$, i.e,

$$a_{1,1}x_1 + \dots + a_{1,n}x_n = b_1,$$

$$a_{2,1}x_1 + \dots + a_{2,n}x_n = b_2,$$

$$\dots$$

$$a_{m,1}x_1 + \dots + a_{m,n}x_n = b_m.$$

Let A be the matrix with entries $A_{i,k} = a_{i,k} \in K$, $b = (b_1, ..., b_m)^\top \in K^m$, and $x = (x_1, ..., x_n)^\top \in K^n$. Then the above system can be concisely written as Ax = b and A can be considered as the matrix of a linear transformation from K^n to K^m (also denoted by A).

We will now consider the questions of existence and uniqueness of solutions of such a system of linear equations. The columns of A are elements of K^m . They span $A(K^n)$, the image of K^n under the linear transformation A. The number of linear independent columns of A is called the column rank of (the matrix) A. It is equal to the dimension of $A(K^n)$ and hence equal to the rank of the linear transformation A.

Let (A, b) be the $m \times (n+1)$ matrix obtained by adjoining b as a column to A. The vector b depends linearly on the columns of A if and only if the matrices A and (A, b) have the same (column) rank. This implies that Ax = b has a solution if and only if rank $(A, b) = \operatorname{rank} A$. Since $n = \operatorname{rank} A + \dim \ker A$ the linear transformation A is injective and hence a solution of Ax = b is unique if and only if rank A = n. The solutions of a homogeneous equation Ax = 0 are given as the elements of ker A. If x and x_p are both solutions of Ax = b then $x - x_p \in \ker A$. Hence the solutions of Ax = b are given as the elements of the coset $x_p + \ker A = \{a + x_p : a \in \ker A\}$. Thus we have proven the following

Theorem. Let A be an $m \times n$ matrix and $b \in K^m$. Then existence and uniqueness of solutions of the system Ax = b is given by the following table

	Ax = b = 0	$Ax = b \neq 0$
$\operatorname{rank}(A, b) > \operatorname{rank} A$	does not happen	system is unsolvable
$\operatorname{rank}(A, b) = \operatorname{rank} A = n$	x = 0 is the unique solution	system is uniquely solvable
$\operatorname{rank}(A, b) = \operatorname{rank} A < n$	system has nontrivial solutions	system has many solutions

Moreover, the set of solutions of Ax = 0 is a subspace of K^n whose dimension is given by $n - \operatorname{rank} A$. If x_p is some solution of Ax = b then any solution x of Ax = b can be expressed by $x = x_p + x_h$ where x_h is a solution of Ax = 0.

5.3.7 Gaussian elimination. Two systems Ax = b and A'x = b' of *m* linear equations in *n* unknowns are called equivalent if they have precisely the same solutions. A system Ax = b is transformed into an equivalent system by any of the following *elementary row operations*:

(a) Multiply an equation by a nonzero scalar.

(b) Add a multiple of one equation to another one.

(c) Interchange any two of the equations.

To these row operations among equations correspond similar row operations among the rows of the matrix (A, b). Two matrices are called row-equivalent if it is possible to transform one into the other by a finite sequence of elementary row-operations. Row equivalence is (of course) an equivalence relation. Of course, the matrices (A, b) and (A', b') are row-equivalent if and only if the systems Ax = b and A'x = b' are equivalent.

A matrix is said to be a *row-echelon matrix* if it satisfies

(a) all zero rows occur below any nonzero row,

(b) the first nonzero entry of each nonzero row occurs to the right of the first nonzero entry of any row above.

The first nonzero entry of a row is called a *pivot*.

Theorem. Every matrix A is row-equivalent to a row-echelon matrix R. The number of linearly independent rows of A (called the row rank of A) equals both the number of linearly independent columns of A (i.e., the column rank of A) and the number of nonzero rows of R.

Sketch of proof: The first claim follows by induction on the number of rows. The elementary row operations leave the number of linearly independent rows of a matrix invariant (as they do not change the space spanned). Since the systems Ax = 0 and Rx = 0 are equivalent, i.e., since ker $A = \ker R$, we obtain rank $A = n - \dim \ker A = n - \dim \ker R = \operatorname{rank} R$, assuming that A is an $m \times n$ matrix. Hence the elementary row operations leave also the column rank of a matrix invariant. The proof is completed by the observation that column rank and row rank of R are equal to the number of nonzero rows of R.

Therefore elementary row operations can be used to determine the rank of A and the rank of (A, b) and hence to answer the questions of existence and uniqueness of solutions of Ax = b. But they are also useful in actually computing the solution: Let (B, b') be a row-echelon matrix which is row equivalent to (A, b) and assume that rank $B = \operatorname{rank}(B, b')$. Then x solves Ax = b if and only if x solves Bx = b'. Let $x = (x_1, ..., x_n)^{\top}$. If column k of B contains the pivot of some row then x_k is called a *determined* and otherwise a *free* variable. Any choice of scalars for the free variables will lead to a solution. The determined variables can now be computed one by one starting with the one having the largest index.

Example: Consider

$$\begin{pmatrix} 1 & -1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1 \\ -5 \\ b \end{pmatrix}.$$

This is equivalent to

$$\begin{pmatrix} 1 & -1 & 1 \\ 0 & 3 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1 \\ -6 \\ b+3 \end{pmatrix}.$$

If $b \neq -3$ then there is no solution. If b = -3 then x_3 is a free variable and x_1 and x_2 are determined variables. The second equation gives $x_2 = -2$ and the first gives $x_1 = 1 + x_2 - x_3 = -1 - x_3$ for any choice of x_3 .

5.4. Eigenvalues and Eigenvectors

5.4.1 Eigenvalues and eigenvectors. Let W be a vector space, V a subspace of W, $T: V \to W$ a linear transformation, and λ a scalar. If there exists a nontrivial (nonzero) element $x \in V$ such that $Tx = \lambda x$ then λ is called an *eigenvalue* of T and x is called an *eigenvector* of T associated with λ . Thus λ is an eigenvalue if and only if $T - \lambda I$ is not injective.

If W = V and V is finite-dimensional then the set of all eigenvalues is called the *spectrum* of T and is denoted by $\sigma(T)$. (The spectrum of a general linear transformation is defined differently and may contain points which are not eigenvalues. Here the word spectrum will only be used in the case described.)

Theorem. Eigenvectors corresponding to distinct eigenvalues are linearly independent.

Sketch of proof: This will follow from the more general Theorem 5.4.2.

Note that for any $\lambda \in K$ the set ker $(T - \lambda I)$ is a subspace of V. If λ is not an eigenvalue then this space is equal to the trivial space $\{0\}$. If λ is an eigenvalue then this space is the set containing zero and all eigenvectors of T which are associated with λ . It is called the *(geometric) eigenspace* of T associated with λ . The dimension of this space is called the geometric multiplicity of λ . We will occasionally say that λ has geometric multiplicity zero if λ is not an eigenvalue.

5.4.2 Generalized eigenvectors. A nontrivial vector x in a subspace V of a vector space W is called a *generalized eigenvector* of the linear transformation $T: V \to W$ associated with λ if there exists $k \in \mathbb{N}$ such that $(T - \lambda I)^k x = 0$.

Theorem. Generalized eigenvectors corresponding to distinct eigenvalues are linearly independent.

Sketch of proof: Let $v_1, ..., v_m$ be generalized eigenvectors of T associated with the pairwise distinct eigenvalues $\lambda_1, ..., \lambda_m$ and define $k_r = \min\{\ell \in \mathbb{N} : (T - \lambda_r I)^\ell v_r = 0\}$. Suppose that $\alpha_1 v_1 + ... + \alpha_m v_m = 0$. Apply the operator

$$(T - \lambda_1 I)^{k_1 - 1} (T - \lambda_2 I)^{k_2} ... (T - \lambda_m I)^{k_m}$$

to both sides of the equation. Since

$$(T - \lambda_1)^{k_1 - 1} (T - \lambda_j)^{k_j} v_1 = \sum_{\ell=0}^{k_j} {k_j \choose \ell} (\lambda_1 - \lambda_j)^{\ell} (T - \lambda_1)^{k_1 - 1 + k_j - \ell} v_1$$
$$= (\lambda_1 - \lambda_j)^{k_j} (T - \lambda_1)^{k_1 - 1} v_1$$

is different from zero whenever $j \neq 1$ it follows that $\alpha_1 = 0$. Similarly, $\alpha_2 = \dots = \alpha_m = 0$.

The set $\bigcup_{k \in \mathbb{N}} \ker(T - \lambda I)^k$ is a subspace of V called the *algebraic eigenspace* of λ . Its dimension is called the algebraic multiplicity of λ . The algebraic eigenspace of λ includes

the geometric eigenspace of λ as a subspace. Hence the algebraic multiplicity is at least as large as the geometric multiplicity.

Theorem. The following three statements are equivalent:

- (a) λ is an eigenvalue.
- (b) The geometric multiplicity of λ is positive.
- (c) The algebraic multiplicity of λ is positive.

Sketch of proof: That (a) implies (b) and that (b) implies (c) is immediate. Next assume that λ is not an eigenvalue. Then $T - \lambda I$ is injective. Since the composition of injective functions is injective we get $(T - \lambda I)^k$ is injective for all $k \in \mathbb{N}$. Thus (c) implies (a). \Box **5.4.3 Index of an eigenvalue.** With respect to any linear transformation T we define the *index* of a scalar λ to be the smallest nonnegative integer ν , if one exists, such that $\ker(T - \lambda I)^{\nu}$ is equal to the algebraic eigenspace of λ . If no such integer exists for λ , we define the index to be infinity. In particular then, the index of λ is zero, if λ is not an eigenvalue at all, and one, if the geometric and the algebraic eigenspace of λ coincide.

5.5. Spectral Theory in Finite-dimensional Complex Vector Spaces

5.5.1 Polynomials of linear transformations. Let V be a finite-dimensional complex vector space, T a linear transformation from V to V, and $f = \sum_{j=0}^{N} \alpha_j z^j$ a polynomial in $\mathbb{C}[z]$. Then one may define the transformation

$$f(T): V \to V: x \mapsto \sum_{j=0}^{N} \alpha_j T^j x,$$

where, by definition, T^0 is the identity transformation. Note that, if f = gh, then f(T) = g(T)h(T) = h(T)g(T).

5.5.2 Existence of eigenvalues. Let $T : V \to V$ be a linear transformation from an *n*-dimensional complex vector space V to itself. If n > 0, then T has at least one eigenvalue and at most n.

Sketch of proof: Choose $0 \neq x \in V$. Then $x, Tx, ..., T^n x$ are a linearly dependent. Hence there exist numbers $\alpha_0, ..., \alpha_n$, at least one of which is not zero, such that $\alpha_0 x + \alpha_1 Tx + ... + \alpha_n T^n x = 0$. Let m be the largest index j such that $\alpha_j \neq 0$. Then $m \geq 1$ and, without loss of generality, $\alpha_m = 1$. Let f be the polynomial defined by $f(z) = z^m + \alpha_{m-1} z^{m-1} + ... + \alpha_0$. Then, by the fundamental theorem of algebra, there exist complex numbers $\lambda_1, ..., \lambda_m$ such that $f(z) = (z - \lambda_1)...(z - \lambda_m)$. Hence

$$0 = f(T)x = (T - \lambda_1 I)...(T - \lambda_m I)x.$$

This implies that at least one of the operators $T - \lambda_j I$ is not injective, i.e., at least one of the λ_j is an eigenvalue. There can be no more than n distinct eigenvalues by Theorem 5.4.1.

5.5.3 Algebraic multiplicities of eigenvalues. If T is a linear transformation in an n-dimensional vector space and λ is a complex number, then the index of λ is at most equal to the algebraic multiplicity of λ and the algebraic multiplicity of λ is at most n.

Sketch of proof: Let m denote the index of λ . There exists an x in the algebraic eigenspace of λ such that $(T - \lambda I)^{m-1}x \neq 0$. Assume that $y = \alpha_0 x + ... + \alpha_{m-1}(T - \lambda I)^{m-1}x = 0$. Applying the operator $(T - \lambda I)^{m-1}$ to y shows that $\alpha_0 = 0$. Next, applying the operator $(T - \lambda I)^{m-2}$ to y shows that $\alpha_1 = 0$, also. Proceeding in this manner shows that the m vectors x, $(T - \lambda I)x$, ..., $(T - \lambda I)^{m-1}x$ are linearly independent and hence that the algebraic eigenspace has dimension at least m. The second statement is obvious, since the algebraic eigenspace of λ is a subspace of V.

5.5.4 Equality of f(T) and g(T). Suppose f, g are polynomials over \mathbb{C} and T is a linear transformation in a finite-dimensional complex vector space. Let $\nu(\lambda)$ denote the index of λ . Then f(T) = g(T) if and only if $\operatorname{ord}_{\lambda}(f-g) \geq \nu(\lambda)$ for all $\lambda \in \mathbb{C}$.

Sketch of proof: We denote the vector space on which T is defined by V. As we might compare the polynomial f - g with the zero polynomial no harm is done in assuming that g = 0.

First assume f(T) = 0. Let λ_0 be an eigenvalue of T with index ν_0 . Then there exists a vector x in the algebraic eigenspace of λ_0 such that $0 \neq y = (T - \lambda_0)^{\nu_0 - 1} x$ is an eigenvector of T associated with λ_0 . Assume $m_0 = \operatorname{ord}_{\lambda_0}(f)$, i.e., there is a polynomial h such that $f(\lambda) = h(\lambda)(\lambda - \lambda_0)^{m_0}$ with $h(\lambda_0) \neq 0$. We have to show that $m_0 \geq \nu_0$ and hence we assume, on the contrary, that $m_0 < \nu_0$. Then

$$0 = f(T)(T - \lambda_0)^{\nu_0 - m_0 - 1} x = h(T)(T - \lambda_0)^{\nu_0 - 1} x = h(T)y = h(\lambda_0)y.$$

Hence $h(\lambda_0) = 0$ which is impossible.

Next assume that $\operatorname{ord}_{\lambda}(f) \geq \nu(\lambda)$ for all $\lambda \in \mathbb{C}$. Suppose the distinct eigenvalues of T are $\lambda_1, ..., \lambda_r$. Let $S = \prod_{j=1}^r (T - \lambda_j I)^{\nu(\lambda_j)}$. Note that our hypothesis implies that there exists a polynomial h such that f(T) = h(T)S. From the definition of the index ν it follows that $\ker S^2 = \ker S$ and hence, by part A of 5.2.3, that $V = \ker S \oplus S(V)$. Also note that S(V) is an invariant subspace for T. Let $T' = T|_{S(V)} : S(V) \to S(V)$. If S(V) has dimension n > 0, then T' has an eigenvalue μ . This implies that μ is also an eigenvalue of T with eigenvector in S(V) which is impossible. Hence $S(V) = \{0\}$ and $f(T)(V) = (h(T)S)(V) = \{0\}$. \Box **5.5.5 Jordan blocks and Jordan matrices.** Let V be a vector space and T a linear transformation from V to V. In the following a subspace W of V will be called T-cyclic, if there exists $x \in W$ and $m \in \mathbb{N}$ such that $B = \{x, Tx, ..., T^{m-1}x\}$ is a basis of W but $T^m x = 0$. In this case we will write $W = [x]_T$. Note that W is an invariant subspace for T.

The ordered basis $(T^{m-1}x, ..., x)$ is called a *Jordan chain*. Its first element is an eigenvector of T associated with the eigenvalue zero. Every other element is a generalized eigenvector of T associated with zero. With respect to this the transformation $T|_W : W \to W$ has the matrix

$$\begin{pmatrix} 0 & 1 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 1 \\ 0 & 0 & 0 & \dots & 0 & 0 \end{pmatrix}$$

Matrices of this form a particularly important. We therefore make the following definitions: A matrix of the form

$$\begin{pmatrix} \lambda & 1 & 0 & \dots & 0 & 0 \\ 0 & \lambda & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \lambda & 1 \\ 0 & 0 & 0 & \dots & 0 & \lambda \end{pmatrix}$$

is called a *Jordan block* with eigenvalue λ . A matrix of the form

$$\begin{pmatrix} J_1 & 0 & \dots & 0 \\ 0 & J_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & J_r \end{pmatrix}$$

where each matrix J_l is a Jordan block is called a *Jordan matrix*.

5.5.6 Nilpotent transformations. A linear transformation T from a finite-dimensional vector space V to itself is called *nilpotent* if there exists a natural number m such that T^m is the zero transformation.

Suppose that $T^m = 0$ and that λ is an eigenvalue of T with associated eigenvector x. Then $T^m x = \lambda^m x$ and this shows that $\lambda = 0$. Hence, a nilpotent transformation has only one eigenvalue, namely zero. That eigenvalue has algebraic multiplicity equal to the dimension of V.

Theorem. Suppose $T: V \to V$ is a nilpotent linear transformation on an *n*-dimensional complex vector space V. Then V is a direct sum of T-cyclic subspaces.

Sketch of proof: The proof is by induction on the dimension of V. Let M be the set of all natural numbers for which the theorem is true. Then $1 \in M$, because dim(V) = 1 and $0 \neq x \in V$ imply that Tx = 0 and hence $V = [x]_T$.

Next assume that $n \in M$ and that $\dim(V) = n+1$. Since T is nilpotent $\dim(T(V)) \leq n$. Therefore there is a subspace F of V of dimension n such that $T(V) \subset F$. Since $\dim(F) = n$ and F is invariant under T there exist vectors $x_1, ..., x_k$ such that $F = \bigoplus_{j=1}^k [x_j]_T$ where the indices are chosen such that $m_1 \leq ... \leq m_k$ when m_j denotes $\dim([x_j]_T)$. Now choose $g \in V - F$. Then there exists $h \in F$ and numbers $\alpha_1, ..., \alpha_k$ such that $Tg = Th + \sum_{j=1}^k \alpha_j x_j$. Indeed, since $Tg \in F$ we have

$$Tg = \sum_{j=1}^{k} \sum_{\ell=0}^{m_j - 1} \alpha_{j,\ell} T^\ell x_j$$

and can therefore choose

$$h = \sum_{j=1}^{k} \sum_{\ell=1}^{m_j - 1} \alpha_{j,\ell} T^{\ell - 1} x_j$$

and $\alpha_j = \alpha_{j,0}$ for j = 1, ..., k.

We now distinguish two cases. In the first case all of the numbers α_j are zero. Defining $x_{k+1} = g - h$ gives $Tx_{k+1} = 0$, $\langle x_{k+1} \rangle = [x_{k+1}]_T$, and $V = F \oplus [x_{k+1}]_T$.

In the second case we have a number p such that $\alpha_p \neq 0$ but $\alpha_j = 0$ whenever j > p. In this case we define $\tilde{x}_p = (g-h)/\alpha_p$ and $F' = \bigoplus_{j \neq p} [x_j]_T$. Then

$$T^{\ell+1}\tilde{x}_p = T^\ell x_p + \sum_{j=1}^{p-1} \frac{\alpha_j}{\alpha_p} T^\ell x_j.$$

This implies that $T^{m_p+1}\tilde{x}_p = 0$ and that $\{\tilde{x}_p, T\tilde{x}_p, ..., T^{m_p}\tilde{x}_p\}$ is linearly independent. Hence $[\tilde{x}_p]_T$ is a *T*-cyclic subspace of dimension $m_p + 1$. Since $[\tilde{x}_p]_T \cap F' = \{0\}$ we obtain that the sum $[\tilde{x}_p]_T + F'$ is direct and has dimension n + 1. Thus $V = [\tilde{x}_p]_T \oplus F'$. This shows that $n + 1 \in M$. The theorem follows now from the induction principle.

The above proof is due to Gohberg and Goldberg (A simple proof of the Jordan decomposition theorem for matrices, American Mathematical Monthly 103 (1996), p. 157 – 159).

Corollary. V has an ordered basis such that the matrix associated with T is a Jordan matrix all of whose Jordan blocks have eigenvalue zero. Conversely, every such matrix is nilpotent.

5.5.7 The structure theorem. The theorem below shows that every linear transformation on a finite-dimensional complex vector space is built up from a number of simple transformation on certain invariant subspaces.

Theorem. Let T be a linear transformation from a finite-dimensional complex vector space V to itself and suppose that $\lambda_1, ..., \lambda_m$ are the distinct eigenvalues of T with respective indices $\nu_1, ..., \nu_m$. Denote by $U_j = \ker(T - \lambda_j I)^{\nu_j}$ the algebraic eigenspaces of λ_j . Then the following statements are true.

(a) T maps each U_j to itself, i.e., each U_j is an invariant subspace with respect to T.

(b) $V = U_1 \oplus \ldots \oplus U_m$.

(c) The sum of the algebraic multiplicities of all eigenvalues of T is equal to the dimension of V.

(d) Each $(T - \lambda_j I)|_{U_j}$ is nilpotent.

(e) Each $T|_{U_i}$ has exactly one eigenvalue, namely λ_j .

Sketch of proof: Suppose $x \in U_j$ and let $y = (T - \lambda_j I)x$. Then $y \in U_j$ and hence $Tx = y + \lambda_j x \in U_j$. This proves (a). Next let $A_j = (T - \lambda_j I)^{\nu_j}$ and $S = A_1...A_m$. By Theorem 5.5.4 we have S = 0, i.e., ker S = V. Claim (b) follows now from part C of 5.2.3. The proofs of (c), (d), and (e) are trivial.

Corollary. Let $T: V \to V$ be a linear transformation on an *n*-dimensional complex vector space V. Let $\lambda_1, ..., \lambda_m$ be the distinct eigenvalues of T with respective algebraic multiplicities $k_1, ..., k_m$. Then there exists a basis of V with respect to which the matrix associated with T is of the form

$$\begin{pmatrix} A_1 & 0 & \dots & 0 \\ 0 & A_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & A_m \end{pmatrix}$$

where, for l = 1, ..., m, A_l is a $k_l \times k_l$ Jordan matrix all of whose Jordan blocks have eigenvalue λ_l .

The matrix described in this corollary is called the *Jordan normal form* of the linear transformation T.

5.5.8 Functional Calculus. Let V be a finite-dimensional complex vector space, T a linear transformation from V to V. Let $\mathcal{F}(T)$ be the set of all functions f for which there is a (not necessarily connected) open set $\Omega(f)$ such that $\Omega(f)$ contains all eigenvalues of T and f is analytic $\Omega(f)$. For any $f \in \mathcal{F}(T)$ there exists a polynomial P such that

$$f^{(m)}(\lambda) = P^{(m)}(\lambda)$$

for every eigenvalue λ and all $m \in \{0, ..., \nu(\lambda) - 1\}$. We then define f(T) = P(T). Note that f(T) is well-defined for, if Q is another polynomial satisfying that condition, then P(T) = Q(T) by theorem 5.5.4.

Theorem. Suppose $f, g \in \mathcal{F}(T)$ and $\alpha, \beta \in \mathbb{C}$. Then

- (1) $\alpha f + \beta g \in \mathcal{F}(T)$ and $(\alpha f + \beta g)(T) = \alpha f(T) + \beta g(T)$.
- (2) $fg \in \mathcal{F}(T)$ and (fg)(T) = f(T)g(T) = g(T)f(T).
- (3) f(T) = 0 if and only if $f^{(m)}(\lambda) = 0$ for every eigenvalue of T and for all $m \in \{0, ..., \nu(\lambda) 1\}$.

5.5.9 Spectral projections. An operator E such that $E^2 = E$ is called *idempotent*. Since $E(V) \cap (1-E)(V) = \{0\}$ an idempotent transformation is a projection in the sense of 5.2.1.

For any $\lambda \in \mathbb{C}$ let $U(\lambda)$ be a neighborhood of λ whose closure contains no eigenvalue other than possibly λ itself. Define

$$e_{\lambda}(\mu) = \begin{cases} 1 & \text{if } \mu \text{ is in } U(\lambda), \\ 0 & \text{if } \mu \text{ is in the interior of } U(\lambda)^c \end{cases}$$

Now let $T: V \to V$ be a linear transformation on a finite-dimensional complex vector space V. Then $e_{\lambda} \in \mathcal{F}(T)$ and the transformation $E_{\lambda} = e_{\lambda}(T)$, called the *eigenprojection* of λ , has the following properties:

(a)
$$E_{\lambda} \neq 0$$
 if and only if $\lambda \in \sigma(T)$.

(b)
$$E_{\lambda}^2 = E_{\lambda}$$
.

- (c) If $\lambda_1 \neq \lambda_2$ then $E_{\lambda_1} E_{\lambda_2} = 0$.
- (d) $\sum_{\lambda \in \sigma(T)} E_{\lambda} = I.$
- (e) $E_{\lambda}(V)$ is the algebraic eigenspace of λ .

In view of 5.5.8 only the last statement needs proof. Let $f: \mu \mapsto (\mu - \lambda)^{\nu(\lambda)} e_{\lambda}(\mu)$. Then f(T) = 0 and hence $E_{\lambda}(V)$ is a subset of $\ker(T - \lambda)^{\nu(\lambda)}$, the algebraic eigenspace of λ . To prove the converse note that we have from (d) that $x = \sum_{\mu \in \sigma(T)} E_{\mu}(x)$. Since $E_{\mu}(x)$ is in the algebraic eigenspace of μ and since the sum of the algebraic eigenspaces is direct we get $x = E_{\lambda}(x)$ if $x \in \ker(T - \lambda)^{\nu(\lambda)}$.

5.5.10 The spectral theorem. Let T be a linear transformation on a finite-dimensional vector space V and $f \in \mathcal{F}(T)$. Then the function

$$g: \mu \mapsto \sum_{\lambda \in \sigma(T)} \sum_{m=0}^{\nu(\lambda)-1} \frac{(\mu - \lambda)^m}{m!} f^{(m)}(\lambda) e_{\lambda}(\mu)$$

satisfies $g^{(m)}(\lambda) = f^{(m)}(\lambda)$ for all eigenvalues λ and all $m \in \{0, ..., \nu(\lambda - 1)\}$. This gives rise to the so called spectral theorem:

Theorem. Let T be a linear transformation on a finite-dimensional vector space V and $f \in \mathcal{F}(T)$. Then

$$f(T) = \sum_{\lambda \in \sigma(T)} \sum_{m=0}^{\nu(\lambda)-1} \frac{(T-\lambda I)^m}{m!} f^{(m)}(\lambda) E_{\lambda}.$$

In particular, if $\nu(\lambda) = 1$ for every eigenvalue λ then

$$T = \sum_{\lambda \in \sigma(T)} \lambda E_{\lambda}.$$

5.5.11 Representation by Cauchy's integral formula. Suppose $f \in \mathcal{F}(T)$ and that the domain $\Omega(f)$ of f has the following properties: there is an open set U such that $\sigma(T) \subset U$, $\overline{U} \subset \Omega(f)$, and the boundary Γ of U consists of finitely many positively oriented, closed, rectifiable Jordan curves. Then

$$f(T) = \frac{1}{2\pi i} \int_{\Gamma} f(z)(zI - T)^{-1} dz.$$

Sketch of proof: If z is in Γ , the boundary of U, then $r_z : U \to \mathbb{C} : \mu \mapsto (z - \mu)^{-1}$ is in $\mathcal{F}(T)$. Since

$$r_z^{(m)}(\mu) = \frac{m!}{(z-\mu)^{m+1}}$$

the spectral theorem implies

$$(zI-T)^{-1} = \sum_{\lambda \in \sigma(T)} \sum_{m=0}^{\nu(\lambda)-1} \frac{(T-\lambda I)^m}{(z-\lambda)^{m+1}} E_{\lambda}.$$

Then use Cauchy's integral formula and again the spectral theorem. \Box **5.5.12 The minimal polynomial and Cayley's theorem.** Let $T: V \to V$ be a linear transformation on a finite-dimensional complex vector space. Suppose that $\lambda_1, ..., \lambda_m$ are its distinct eigenvalues whose respective indices are $\nu_1, ..., \nu_m$. Then the polynomial

$$(\lambda - \lambda_1)^{\nu_1} \dots (\lambda - \lambda_m)^{\nu_m}$$

is called the *minimal polynomial* of T.

In view of Theorem 5.5.4 we obtain immediately the validity of the following theorem of Cayley.

Theorem. If q is (a multiple of) the minimal polynomial of T then q(T) = 0.

Among all polynomials f with leading coefficient one satisfying f(T) = 0 the minimal polynomial of T is the one of lowest degree.

5.5.13 The characteristic polynomial, trace and determinant. Let $T: V \to V$ be a linear transformation on a finite-dimensional complex vector space. Suppose that $\lambda_1, ..., \lambda_m$ are its distinct eigenvalues whose respective algebraic multiplicities are $k_1, ..., k_m$. Then the polynomial

$$(\lambda - \lambda_1)^{k_1} \dots (\lambda - \lambda_m)^{k_m} = \lambda^n - a_1 \lambda^{n-1} + \dots + (-1)^n a_n$$

is called the *characteristic polynomial* of T. The number $a_1 = k_1\lambda_1 + ... + k_m\lambda_m$ is called the *trace* of T and is denoted by tr T. The number $a_n = \lambda_1^{k_1} ... \lambda_m^{k_m}$ is called the *determinant* of T and is denoted by det T. Since $\lambda I - T$ has an eigenvalue $\lambda - \lambda_j$ of algebraic multiplicity k_j if and only if T has an eigenvalue λ_j of the same multiplicity it follows that the characteristic polynomial of T equals det $(\lambda I - T)$.

Theorem. Suppose $T: V \to V$ is a linear transformation on a complex finite-dimensional vector space. Then T is invertible if and only if det $T \neq 0$. In this case T^{-1} is also a linear transformation from V to V.

5.6. Multilinear Algebra

5.6.1 Multilinear forms. Suppose V is a vector space over the field K. To any function $f: V^k \to K$ and any element $(x, y) \in V^{j-1} \times V^{k-j}$ we may associate the function $f_{x,y}: V \to K: t \mapsto f(x, t, y)$. The function f is called a *multilinear form* or a *(covariant) tensor* of rank k over V if for all $j \in \{1, ..., k\}$ and all $(x, y) \in V^{j-1} \times V^{k-j}$ the function $f_{x,y}$ is linear.

If V has dimension n the set T_k of all tensors of rank k over V is a vector space of dimension n^k .

5.6.2 Antisymmetric tensors. Suppose f is a tensor of rank k over the vector space V. Then f is called *antisymmetric* if for every transposition τ the relationship $f(x_1, ..., x_k) = -f(x_{\tau 1}, ..., x_{\tau k})$ holds. Equivalently, f is called antisymmetric when $j \neq l$ and $x_j = x_l$ imply that $f(x_1, ..., x_k) = 0$.

The antisymmetric tensors of rank k form a subspace of T_k .

5.6.3 Determinant forms. Let V be an n-dimensional vector space and $\{b_1, ..., b_n\}$ a basis of V. An antisymmetric tensor f of rank n is called a determinant form. It is uniquely

determined when $f(b_1, ..., b_n)$ is given. In particular, the antisymmetric tensors of rank n form a one-dimensional vector space. To see this let $x_j = \sum_{k=1}^n \alpha_{k,j} b_k$. Then

$$f(x_1, ..., x_n) = \sum_{k_1, ..., k_n} \alpha_{k_1, 1} ... \alpha_{k_n, n} f(b_{k_1}, ..., b_{k_n}).$$

The only nonzero terms in this sum are those for which $(k_1, ..., k_n) = \pi(1, ..., n)$, where π is a permutation on n letters. Let $(-1)^{\pi}$ denote +1 or -1 depending on whether π is an even or an odd permutation. Then $f(b_{k_1}, ..., b_{k_n}) = (-1)^{\pi} f(b_1, ..., b_n)$ and hence

$$f(x_1, ..., x_n) = f(b_1, ..., b_n) \sum_{\pi \in S_n} (-1)^{\pi} \alpha_{\pi 1, 1} ... \alpha_{\pi n, n}.$$

In particular, after having chosen an *n*-tuple $(b_1, ..., b_n)$ of basis vectors there is one and only one determinant form D such that $D(b_1, ..., b_n) = 1$.

5.6.4 Determinants. Let V be an n-dimensional vector space. Suppose $T: V \to V$ is a linear transformation. Given any nonzero determinant form $f: V^n \to K$ define $f_T: V^n \to K: (x_1, ..., x_n) \mapsto f(Tx_1, ..., Tx_n)$ which is also a determinant form.

Theorem. Let $\{b_1, ..., b_n\}$ be any basis of V. Then

$$\det T = \frac{f_T(b_1, \dots, b_n)}{f(b_1, \dots, b_n)} = \frac{f(Tb_1, \dots, Tb_n)}{f(b_1, \dots, b_n)}.$$

Sketch of proof: Since the space of determinant forms is one-dimensional f_T is a multiple of f (f is a basis vector). This shows that $f(Tb_1, ..., Tb_n)/f(b_1, ..., b_n)$ is independent of the basis chosen. Recall from 5.5.7 that V has a basis $\{b_1, ..., b_n\}$ such that

$$Tb_j = \lambda_j b_j + \alpha_j b_{j-1}$$

where $\alpha_j \in \{0, 1\}$. Among the numbers $\lambda_1, ..., \lambda_n$ each eigenvalue appears as often as indicated by its multiplicity.

5.6.5 The multiplicative property of determinants. Let V be an n-dimensional vector space and suppose that T and S are linear transformations from $V \to V$. If either S or T is not injective, then neither is TS. Hence, by Theorem 5.5.12, $\det(TS) = \det(T) \det(S) = 0$. Now suppose T, S, and TS are injective and $\{b_1, ..., b_n\}$ is a basis of V. Then $\{Sb_1, ..., Sb_n\}$ is likewise a basis of V and hence

$$\det(TS) = \frac{f(TSb_1, ..., TSb_n)}{f(b_1, ..., b_n)} = \frac{f(TSb_1, ..., TSb_n)}{f(Sb_1, ..., Sb_n)} \frac{f(Sb_1, ..., Sb_n)}{f(b_1, ..., b_n)} = \det(T) \det(S).$$

In particular, $\det(T^{-1}) = \det(T)^{-1}$.

5.6.6 Determinants of matrices. In K^n choose the canonical basis $(e_1, ..., e_n)$. Since the linear transformations from K^n to K^n are in a one-to-one correspondence to $n \times n$ -matrices the determinant of a matrix A is defined as the determinant of the associated linear transformation. Equivalently,

$$\det A = \sum_{\pi \in S_n} (-1)^{\pi} A_{\pi 1,1} \dots A_{\pi n,n}$$

if the entries of A are denoted by $A_{j,k}$.

We list the most important properties of the determinant of an $n \times n$ -matrix A:

- (1) det $A = \det A^{\top}$, where A^{\top} denotes the transpose of A.
- (2) A transposition of two columns or two rows of A changes the sign of the determinant.
- (3) Adding to a column (or row) of A a linear combination of the remaining columns (or rows) does not change the determinant.

- (4) Multiplying a column or row by a constant results in a multiplication of det A by that constant.
- (5) If the columns or rows of A are linearly dependent (in particular, if two columns or rows are equal) then $\det A = 0$.
- (6) $\det(AB) = \det(A) \det(B)$, in particular, $\det(A^{-1}) = \det(A)^{-1}$ and $\det I = 1$ when I is the identity matrix.

5.6.7 Laplace's expansion theorem.

Theorem. Let A be an $n \times n$ -matrix and choose a partition (J, K) of $\{1, ..., n\}$ such that $j_1 < ... < j_p$ and $k_1 < ... < k_q$ when $J = \{j_1, ..., j_p\}$ and $K = \{k_1, ..., k_q\}$. Let Γ be the set of those permutations γ in S_n such that

$$\gamma j_1 < \dots < \gamma j_p \text{ and } \gamma k_1 < \dots < \gamma k_q$$

Furthermore let U_{γ} be the submatrix of A consisting only of the columns numbered $j_1, ..., j_p$ and the rows numbered $\gamma j_1, ..., \gamma j_p$ and V_{γ} the submatrix of A consisting only of the columns numbered $k_1, ..., k_q$ and the rows numbered $\gamma k_1, ..., \gamma k_q$. Then

$$\det A = \sum_{\gamma \in \Gamma} (-1)^{\gamma} \det U_{\gamma} \det V_{\gamma}$$

Sketch of proof: Let A_{γ} be the set of all permutations in S_n whose support is contained in $\{\gamma k_1, ..., \gamma k_q\}$ and B_{γ} the set of all permutations in S_n whose support is contained in $\{\gamma j_1, ..., \gamma j_q\}$. Then every permutation in S_n can be represented uniquely as $\alpha \circ \beta \circ \gamma$ where $\gamma \in \Gamma, \alpha \in A_{\gamma}, \text{ and } \beta \in B_{\gamma}.$ \square

By choosing p = 1 and $j_1 = j$ we expand the determinant along row j. Since the determinant of A equals the determinant of the transpose we can as well expand along a column.

Corollary. Suppose A is an $n \times n$ -matrix. Denote the matrix obtained by deleting row j and column k by M(j,k). Then

$$\det A = \sum_{k=1}^{n} (-1)^{j+k} A_{j,k} \det M(j,k) = \sum_{j=1}^{n} (-1)^{j+k} A_{j,k} \det M(j,k).$$

5.7. Normed Spaces and Inner product Spaces

5.7.1 Norms. Let V be a vector space over either the real or the complex numbers, i.e., $K = \mathbb{R}$ or $K = \mathbb{C}$.

A function $\|\cdot\|: V \to \mathbb{R}$ is called a norm on V if it satisfies the following conditions: (a) $||x|| \ge 0$ for all $x \in V$,

- (b) ||x|| = 0 if and only if x = 0,
- (c) $\|\alpha x\| = |\alpha| \|x\|$ for all $\alpha \in K$ and all $x \in V$,
- (d) $||x + y|| \le ||x|| + ||y||$ for all $x, y \in V$. (This inequality is called the triangle inequality.) If there exists a norm $\|\cdot\|$ on V then $(V, \|\cdot\|)$ is called a normed vector space.

Examples: \mathbb{R}^n and \mathbb{C}^n can be considered as normed vector spaces. Denote the components of $x \in \mathbb{C}^n$ by $x_1, ..., x_n$. Then the following are norms:

- (a) $||x||_{\infty} = \max\{|x_1|, ..., |x_n|\},\$

(b) $||x||_p = \left(\sum_{k=1}^n |x_k|^p\right)^{1/p}$ if $p \ge 1$. The set $C^0([a, b])$, i.e., the set of continuous complex-valued or real-valued functions on [a, b], can be considered as a normed vector space when [a, b] is a closed interval in \mathbb{R} . The following are norms on $C^0([a, b])$:

- (a) $||f||_{\infty} = \max\{|f(x)| : x \in [a, b]\},\$
- (b) $||f||_p = \left(\int_a^b |f(x)|^p dx\right)^{1/p}$ if $p \ge 1$.

5.7.2 Inner products. Again let V be a vector space over either the real or the complex numbers. A function $(\cdot, \cdot) : V \times V \to K$ is called an *inner product* or a *scalar product* if it satisfies the following conditions:

- (a) $(x, x) \ge 0$ for all $x \in V$,
- (b) (x, x) = 0 if and only if x = 0,
- (c) $(\alpha x + \beta y, z) = \alpha(x, z) + \beta(y, z)$, for all $\alpha, \beta \in K$ and all $x, y, z \in V$,
- (d) $(x, y) = \overline{(y, x)}$ for all $x, y \in V$.

Here $\overline{\alpha}$ equals α or its complex conjugate depending on whether $K = \mathbb{R}$ or $K = \mathbb{C}$.

Note that y = 0 if and only if (x, y) = 0 for all $x \in V$.

If $K = \mathbb{R}$ the inner product is *bilinear* (linear in both of its arguments). If $K = \mathbb{C}$ the inner product is linear in its first argument but *antilinear* in its second: $(x, \alpha y + \beta z) = \overline{\alpha}(x, y) + \overline{\beta}(x, z)$.

If there exists an inner product (\cdot, \cdot) on V then $(V, (\cdot, \cdot))$ is called an *inner product space*. Examples: \mathbb{R}^n and \mathbb{C}^n can be considered as inner product spaces. An inner product on \mathbb{C}^n is given by $(x, y) = \sum_{k=1}^n x_k \overline{y_k}$. Also $C^0([a, b])$ can be considered as an inner product space when [a, b] is a closed interval in \mathbb{R} : an inner product is $(f, g) = \int_a^b f \overline{g} dx$.

Theorem. 1. An inner product satisfies Schwarz's inequality, i.e.,

$$|(x,y)| \le (x,x)^{1/2} (y,y)^{1/2}$$

2. Every inner product space is a normed vector space under the norm $x \mapsto ||x|| = (x, x)^{1/2}$.

Sketch of proof: 1. Assume $(x, y) \neq 0$. Let $\alpha = |(x, y)|/(y, x)$. For any real r

$$0 \le (x - r\alpha y, x - r\alpha y) = (x, x) - 2r|(x, y)| + r^2(y, y).$$

Now $(y, y) \neq 0$ since otherwise (x, y) would be zero. Schwarz's inequality follows now from choosing r = |(x, y)|/(y, y).

2. The triangle inequality follows from Schwarz's inequality.

5.7.3 Orthogonality. Suppose (\cdot, \cdot) is an inner product on a vector space V. If (x, y) = 0 we say that x and y are *orthogonal* and denote this by $x \perp y$.

Let M be a subset of V and define $M^{\perp} = \{x \in V : (\forall y \in M : (x, y) = 0)\}$. If $x_1, x_2 \in M^{\perp}$ then so is $\alpha x_1 + \beta x_2$, i.e., M^{\perp} is a subspace of V. M^{\perp} is called the orthogonal complement of M.

Theorem. Let X be a subset of an inner product space. If the elements of X are nonzero and pairwise orthogonal then X is linearly independent.

Sketch of proof: Take the inner product of a linear combination of the vectors with each of the vectors themselves. $\hfill \Box$

5.7.4 Orthonormal subsets. A set X whose elements have norm one and are pairwise orthogonal is called *orthonormal*.

Theorem. Let X be a linearly independent countable (e.g., finite) subset of an inner product space. Then $\langle X \rangle$, the span of X, has an orthonormal basis.

Sketch of proof: The basis can be produced from X using the so called Gram-Schmidt process: define $z_1 = x_1/||x_1||$. Then z_1 has norm one and spans $\langle \{x_1\} \rangle$. Next assume that,

for some k < n, the set $\{z_1, ..., z_k\}$ is orthonormal and spans $\langle \{x_1, ..., x_k\} \rangle$. Define

$$y_{k+1} = x_{k+1} - \sum_{j=1}^{k} (x_{k+1}, z_j) z_j.$$

Then y_{k+1} is different from zero and orthogonal to each of the z_j . Finally, define $z_{k+1} = y_{k+1}/||y_{k+1}||$ to obtain an orthonormal set $\{z_1, ..., z_k, z_{k+1}\}$ which spans $\langle \{x_1, ..., x_{k+1}\} \rangle$. Induction shows therefore the existence of an orthonormal and hence linearly independent set $\{z_1, z_2, ...\}$ which spans $\langle X \rangle$.

Let x be a linear combination of elements of an orthonormal set, i.e., $x = \alpha_1 x_1 + ... + \alpha_n x_n$ when $x_1, ..., x_n$ are distinct elements of some orthonormal set. Then the coefficients α_j are given by $\alpha_j = (x, x_j)$.

5.7.5 Finite-dimensional inner product spaces. We restrict our attention again to finite-dimensional spaces since the appropriate generalizations to infinite-dimensional spaces require some topology (closed subsets, continuity, limits).

Let M be a subspace of the finite-dimensional inner product space V. Then $M \cap M^{\perp} = \{0\}$ and $V = M \oplus M^{\perp}$. In particular, $(M^{\perp})^{\perp} = M$. For each $x \in V$ there exists a unique decomposition x = Px + Qx where $Px \in M$ and $Qx \in M^{\perp}$. As in 5.2.1 the mappings $P: V \to M$ and $Q: V \to M^{\perp}$ are projections, which, in this case, are called orthogonal projections.

Theorem Pythagorean theorem. With the above notation we have

$$\|x\|^2 = \|Px\|^2 + \|Qx\|^2.$$

5.7.6 Dual spaces and the representation theorem. Let $L: V \to K$ be a linear transformation on the vector space V. Then L is called a *linear functional*. The set of all such linear functionals is a vector space (define $(\alpha L_1 + \beta L_2)(x) = \alpha L_1 x + \beta L_2 x$). This space of functionals is called the dual space of V and is denoted by V^* . The elements of V^* are precisely the tensors of rank 1 over V.

Examples: A $1 \times n$ matrix gives rise to a linear functional on K^n . In fact these are all linear functionals on K^n . The set of all $1 \times n$ matrices is again a vector space which is in a one-to-one correspondence with K^n .

Let $g: [a, b] \to \mathbb{C}$ be a continuous function. Then $G: C^0([a, b]) \to \mathbb{C}$ given by $Gf = \int_a^b f(x)g(x)dx$ is a linear functional. However, the continuous functions g give not rise to all linear functionals on $C^0([a, b])$. In particular the map $f \mapsto f(x_0)$, where x_0 is a fixed point in [a, b], is a linear functional but can not be represented by any of the functionals G. Thus $C^0([a, b])$ can be considered as a proper subset of its own dual space.

The following theorem determines the dual space of any finite-dimensional inner product space.

Theorem Representation theorem. Let V be a real or complex finite-dimensional inner product space and L a linear functional on V. Then there exists a unique $y \in V$ such that Lx = (x, y) for all $x \in V$. Conversely, for every $y \in V$ the function $x \mapsto (x, y)$ is a linear functional on V. In particular, there exists a bijection from V^* to V. This bijection is linear if the scalar field is \mathbb{R} and antilinear if the scalar field is \mathbb{C} .

5.7.7 Adjoints. Let V be a finite-dimensional inner product space and $T: V \to V$ a linear transformation. Choose $y \in V$. Then $x \mapsto (Tx, y)$ is a linear functional. Hence, by Theorem 5.7.6, for each y there is a unique $z \in V$ such that (Tx, y) = (x, z) and this relationship defines a function $T^*: V \to V$, i.e., $T^*y = z$. Hence

$$S = T^* \quad \Leftrightarrow \quad \forall x, y \in V : (Tx, y) = (x, Sy).$$

 T^* is a linear transformation. It is called the *adjoint* of T. Note that $T^{**} = T$ since

$$(T^*x,y) = \overline{(y,T^*x)} = \overline{(Ty,x)} = (x,Ty).$$

Let $\{x_1, ..., x_n\}$ be an orthonormal basis of V. Suppose that, with respect to that basis, T is represented by the matrix A while T^* is represented by the matrix B, i.e., $Tx_j = \sum_{l=1}^n A_{l,j}x_l$ and $T^*x_j = \sum_{l=1}^n B_{l,j}x_l$. Observe that

$$A_{l,j} = (Tx_j, x_l) = (x_j, T^*x_l) = (x_j, \sum_{k=1}^n B_{k,l}x_k) = \overline{B_{j,l}}.$$

Hence we found that the matrix representing T^* is the complex conjugate of the transpose of the matrix representing T (or just the transpose if the scalar field is \mathbb{R}). We will write $B = A^*$.

Theorem. Let $T: V \to V$ be a linear transformation on a finite-dimensional inner product space. Then ker $T = T^*(V)^{\perp}$ and ker $T^* = T(V)^{\perp}$. In particular, T and T^* have the same rank.

Sketch of proof: The following sequence of equivalences shows the first statement.

$$x \in \ker T \Leftrightarrow Tx = 0 \Leftrightarrow \forall y : (Tx, y) = 0 \Leftrightarrow \forall y : (x, T^*y) = 0 \Leftrightarrow x \in (T^*(V))^{\perp}. \quad \Box$$

5.7.8 Normal linear transformations. Suppose V is a finite-dimensional vector space. A linear transformation $T: V \to V$ on an inner product space V is called *normal* if it commutes with its adjoint, i.e., if $TT^* = T^*T$.

If T is normal and $x \in V$ then $(Tx, Tx) = (T^*x, T^*x)$. This implies that ker $T = \ker T^*$ for any normal linear transformation.

If T is normal and $x \in \ker T^2$, then $Tx \in \ker T = \ker T^*$. Therefore $0 = (T^*Tx, x) = (Tx, Tx)$ and hence $\ker T = \ker T^2$. Since T is normal if and only if $T - \lambda I$ is normal, this implies that the algebraic and geometric multiplicities of any eigenvalue of a normal linear transformation coincide. In particular, if V, the domain of T, is a complex vector space then V has a basis of eigenvectors of T.

Suppose λ and λ' are distinct eigenvalues of a normal linear transformation T and that x and x' are the associated eigenvectors. Then

$$(\lambda - \lambda')(x, x') = (Tx, x') - (x, T^*x') = 0$$

which implies that (x, x') = 0, i.e., that x and x' are orthogonal. In particular, if V, the domain of T, is a complex vector space then V has an orthonormal basis of eigenvectors of T. In summary we have the

Theorem Spectral theorem. Let T be a normal linear transformation on a finite-dimensional complex vector space V with distinct eigenvalues $\lambda_1, ..., \lambda_m$. Then there exist pairwise orthogonal subspaces $U_1, ..., U_m$ such that $V = U_1 \oplus ... \oplus U_m$ and

$$T = \sum_{j=1}^{m} \lambda_j P_j$$

where the P_i are orthogonal projections onto U_i .

5.7.9 Self-adjoint linear transformations. A linear transformation on an inner product space is called *self-adjoint* if $T = T^*$. Note that every self-adjoint linear transformation is normal.

Example: A linear transformation P is an orthogonal projection if and only if it is idempotent and self-adjoint. i.e., $P = P^2 = P^*$.

Theorem. Every eigenvalue of a self-adjoint linear transformation is real.

Sketch of proof: Suppose T is self-adjoint and λ is an eigenvalue of T with associated eigenvector x. Then

$$\lambda(x,x) = (Tx,x) = (x,Tx) = \lambda(x,x). \quad \Box$$

Index

K-isomorphism, 54 adjoint, 82algebra associative, 14 algebraic, 54 antilinear, 80antisymmetric, 77 associative, 13automorphism of vector spaces, 65axiom, 7basis, 64 bijective, 11 bilinear, 80binary operation, 13cardinality, 16 Cartesian product, 10center, 27 centralizer, 27 characteristic, 53characteristic polynomial, 77 commutative, 13 complement of a set, 8composition factors, 57composition series, 57congruent, 36conjugate, 54connective, 1 constant, 4contrapositive, 1converse, 1coordinate, 68 countable, 17countably infinite, 17 cyclic, 73 de Morgan's laws in logic, 2in set theory, 9 definition, 7 determinant, 77determinant form, 77

difference of sets, 8dimension, 64direct sum internal, 65 distributive, 13 left, 13 right, 13 domain of a relation, 10domain of discourse, 5 eigenprojection, $\mathbf{76}$ eigenspace algebraic, 71 geometric, 71 eigenvalue, 71 eigenvector, 71 generalized, 71 element, 5 elementary row operations, 70 empty set, 8equality of sets, 8extension, 12field, 14field extension, 53finite, 54 finite set, 17 fixed field, 58formula, 4free abelian group, $\frac{31}{2}$ free group, 31function, 11 bijective, 11 injective, 11 surjective, 11Galois group, 59 greatest lower bound, 12 ${\rm group},\, 14$ idempotent, 75 identity, 13 image of a relation, 10

INDEX

of a set under a relation, 10index of an eigenvalue, 72 infimum, 12 infinite, 17 countably, 17 injective, 11 inner product, 80inner product space, 80intersection of sets, 8invariant, 67 inverse element, 13irreducible, 19 isomorphic vector spaces, 65isomorphism of vector spaces, 65Jordan block, 73 Jordan chain, 73 Jordan matrix, 74 Jordan normal form, 75 least upper bound, 12linear combination, 63linear functional, 81 linear independence, 64 logical equivalence, 2 logical implication, 1 lower bound, 12maximal element, 12maximum, 12minimal element, 12minimal polynomial, 77 minimum, 12 multilinear form, 77 multiplicity of a zero, 44nilpotent, 74 normal, 82 normal field extension, 56one-to-one, 11 onto, 11 order of a group, 23 of a group element, 27ordered pair, 10ordering partial, 12total, 12well. 12orthogonal, 80orthonormal, 80pairwise disjoint, 8parity, 26 partition, 11 permutation, 25 pivot, 70

polynomial function, 42power set, 8predicate, 4preimage of a set under a relation, 10premise, 2 prime number, 19 proof by contradiction, 4quantifier existential, 5 universal, 5 R-module, 14 radical adjunction, 60radical extension, 60simple, 60range, 11 rank, 66 relation, 10 antisymmetric, 12composite, 10equivalence, 11 inverse, 10 reflexive, 11 symmetric, 11 transitive, 11 residue class, 36restriction, 12 ring, 14 root, 44row-echelon matrix, 70 scalar multiplication, 14 scalar product, 80scope of a quantifier, 5self-adjoint, 82sentence, 1 singular, 4set, 5spectrum, 71 splitting field, 55statement, 1subfield, 53 subset, 8proper, 8 subspace, 63successor, 14successor function, 14 supremum, 12 surjective, 11 tautology, 1tensor, 77 theorem, 7 trace, 77 transcendental, 54 transformation linear, 65

INDEX

triangle inequality, 22, 79

uncountable, 17 union of sets, 8 unitary subring, 35 upper bound, 12 variable, 4 bound, 5 determined, 70 free, 5, 70 vector space, 14 finite-dimensional, 64 infinite-dimensional, 64

well-defined, 17

 $m zero, \, {44}$