

Discrete Structures (2IT25)

Rob Hoogerwoord & Hans Zantema

July 11, 2011

Contents

1	Relations	6
1.1	Binary relations	6
1.2	Equivalence relations	10
1.3	Operations on Relations	14
1.3.1	Set operations	14
1.3.2	Transposition	15
1.3.3	Composition	16
1.3.4	Closures	18
1.4	Warshall's Algorithm	23
1.4.1	Introduction	23
1.4.2	Preliminaries	24
1.4.3	The algorithm	25
1.5	Exercises	27
2	Graphs	32
2.1	Directed Graphs	32
2.2	Undirected Graphs	33
2.3	A more compact notation for undirected graphs	35
2.4	Additional notions and some properties	35
2.5	Connectivity	37
2.5.1	Paths	37
2.5.2	Path concatenation	39
2.5.3	The triangular inequality	40
2.5.4	A lemma and its proof	41
2.5.5	Connected components	42
2.6	Cycles	43
2.6.1	Directed cycles	43
2.6.2	Undirected cycles	44
2.7	Euler and Hamiltonian cycles	45
2.7.1	Euler cycles	45
2.7.2	Hamiltonian cycles	47
2.7.3	A theorem on Hamiltonian cycles	47
2.7.4	A proof by contradiction	48
2.7.5	A more explicit proof	49
2.7.6	Proof of the Core Property	49
2.8	Ramsey's theorem	52
2.8.1	Introduction	52
2.8.2	Ramsey's theorem	53
2.8.3	A few applications	56
2.9	Trees	57
2.9.1	Undirected trees	57
2.9.2	Rooted trees	59
2.10	Exercises	61

3	Functions	63
3.1	Functions	63
3.2	Equality of functions	64
3.3	Monotonicity of function types	64
3.4	Function composition	65
3.5	Lifting a function	66
3.6	Surjective, injective, and bijective functions	68
3.7	Some counting arguments	71
3.8	Of finite and infinite	75
3.9	A Useful Classification	76
3.9.1	Several kinds of relations	76
3.9.2	Several kinds of functions	77
3.9.3	Applications	77
3.9.4	A more abstract and more algebraic approach	79
3.9.5	Afterthoughts	83
3.10	Exercises	83
4	Posets and lattices	86
4.1	Partial orders	86
4.2	Indirect Equality	89
4.3	Extreme elements	89
4.4	Upper and lower bounds	92
4.5	Lattices	95
4.5.1	Definition	95
4.5.2	Algebraic properties	97
4.5.3	Distributive lattices	99
4.5.4	Complete lattices	100
4.6	Exercises	102
4.7	Monotonic and continuous functions	103
4.8	Fixed point theorems	105
4.8.1	Least solutions	105
4.8.2	Fixed points and prefix points	106
4.8.3	Existence of fixed points	108
4.9	Complete Partial Orders	110
4.10	Boolean algebras	111
4.11	Applications	114
4.11.1	Recursively defined sets	114
4.11.2	The natural numbers and Mathematical Induction	115
4.11.3	Finite lists	117
4.11.4	Closures of relations	119
4.11.5	Grammars and languages	119
4.12	Exercises	120

5	Monoids and Groups	122
5.1	Operators and their properties	122
5.2	Semigroups and monoids	124
5.3	Groups	126
5.4	Subgroups	129
5.5	Cosets and Lagrange's Theorem	130
5.6	Permutation Groups	131
5.6.1	Function restriction and extension	131
5.6.2	Continued Compositions	132
5.6.3	Bijections	133
5.6.4	Permutations	133
5.6.5	Swaps	134
5.6.6	Neighbour swaps	136
5.6.7	Even and odd permutations	137
5.7	Exercises	139
6	Combinatorics: the art of counting	142
6.1	Introduction	142
6.2	Recurrence Relations	149
6.2.1	An example	149
6.2.2	The characteristic equation	150
6.2.3	A strange, but beautiful, phenomenon	152
6.2.4	Linear recurrence relations	152
6.2.5	Keeping simple things simple	155
6.2.6	Slightly more complicated cases	156
6.2.7	A more computational approach	158
6.3	Binomial Coefficients	162
6.3.1	Factorials	162
6.3.2	Binomial coefficients	163
6.3.3	The Shepherd's Principle	165
6.3.4	Newton's binomial formula	166
6.4	A few examples	167
6.5	Exercises	169
7	Number Theory	173
7.1	Introduction	173
7.2	Divisibility	173
7.3	Greatest common divisors	176
7.4	Euclid's algorithm and its extension	179
7.5	Equations and their solutions	181
7.6	The prime numbers	183
7.7	Modular Arithmetic	187
7.7.1	Congruence relations	187
7.7.2	An application: the nine and eleven tests	190
7.7.3	Linear congruence equations	191

7.7.4	An example	192
7.7.5	Multiple linear congruences: an example	193
7.7.6	Two linear congruences: the general case	195
7.7.7	The Chinese Remainder Theorem	197
7.8	Fermat's little theorem	198
7.9	Cryptography: the RSA algorithm	199
7.10	Exercises	201

1 Relations

1.1 Binary relations

A (binary) *relation* R from set U to set V is a subset of the cartesian product $U \times V$. If $(u, v) \in R$, we say that u is in relation R to v . We usually denote this by uRv . Set U is called the *domain* of the relation and V its *range* (or: *codomain*). If $U = V$ we call R an (*endo*)*relation on* U .

1.1 Examples.

- (a) “Is the mother of” is a relation from the set of all women to the set of all people. It consists of all pairs $(person1, person2)$ where $person1$ is the mother of $person2$. Of course, this relation also is an (endo)relation on the set of people.
- (b) “There is a train connection between” is a relation on the set of cities in the Netherlands.
- (c) The equality relation “=” is a relation on every set. This relation is often denoted by I (and also called the “identity” relation). Because, however, every set has its “own” identity relation we sometimes use subscription to distinguish all these different identity relations. That is, for every set U we define I_U by:

$$I_U = \{ (u, u) \mid u \in U \} .$$

Whenever no confusion is possible and it is clear which set is intended, we drop the subscript and write just I instead of I_U , and in ordinary mathematical language we use “=”, as always. So, for any set U and for all $u, v \in U$, we have: $uIv \Leftrightarrow u = v$.

- (d) Integer n divides integer m , notation $n|m$, if there is an integer $q \in \mathbb{Z}$ such that $q * n = m$. Divides $|$ is the relation on \mathbb{Z} that consists of all pairs $(n, m) \in \mathbb{Z} \times \mathbb{Z}$ with $(\exists q : q \in \mathbb{Z} : q * n = m)$.
- (e) “Less than” ($<$) and “greater than” ($>$) are relations on \mathbb{R} , and on \mathbb{Q} , \mathbb{Z} , and \mathbb{N} as well, and so are “at most” (\leq) and “at least” (\geq).
- (f) The set $\{(a, p), (b, p), (b, q), (c, q)\}$ is a relation from $\{a, b, c\}$ to $\{p, q\}$.
- (g) The set $\{(x, y) \in \mathbb{R}^2 \mid y = x^2\}$ is a relation on \mathbb{R} .
- (h) Let Ω be a set, then “is a subset of” (\subseteq) is a relation on the set of all subsets of Ω .

Besides binary relations we can also consider n -ary relations for any $n \geq 0$. An n -ary relation on sets U_0, \dots, U_{n-1} is a subset of the cartesian product $U_0 \times \dots \times U_{n-1}$. Unless stated otherwise, in this text relations are binary.

Let R be a relation from set U to set V . Then for each element $u \in U$ we define $[u]_R$ as a subset of V , as follows:

$$[u]_R = \{v \in V \mid uRv\} .$$

(Sometimes $[u]_R$ is also denoted by $R(u)$.) This set is called the (R -)image of u . Similarly, for $v \in V$ a subset of U called ${}_R[v]$ is defined by:

$${}_R[v] = \{u \in U \mid uRv\} ,$$

which is called the (R -)pre-image of v .

1.2 Definition. If R is a relation from finite set U to finite set V , then R can be represented by means of a so-called *adjacency matrix*; sometimes this is convenient because it allows computations with (finite) relations to be carried out in terms of matrix calculations. We will see examples of this later.

With m for the *size* – the number of elements – of U and with n for the size of V , sets U and V can be represented by finite sequences, by numbering their elements. That is, we assume $U = \{u_1, \dots, u_m\}$ and we assume $V = \{v_1, \dots, v_n\}$. The adjacency matrix of relation R then is an $m \times n$ matrix A_R , say, the elements of which are 0 or 1 only, and defined by, for all $i, j: 1 \leq i \leq m \wedge 1 \leq j \leq n$:

$$A_R[i, j] = 1 \Leftrightarrow u_i R v_j .$$

Here $A_R[i, j]$ denotes the element of matrix A_R at row i and column j . Note that this definition is equivalent to stating that $A_R[i, j] = 0$ if and only if $\neg(u_i R v_j)$, for all i, j . Actually, adjacency matrices are *boolean* matrices in which, for the sake of conciseness, true is encoded as 1 and false as 0; thus, we might as well state that: $A_R[i, j] \Leftrightarrow u_i R v_j$.

□

Notice that this representation is not *unique*: the elements of finite sets can be assigned numbers in very many ways, and the distribution of 0's and 1's over the matrix depends crucially on how the elements of the two sets are numbered. For instance, if U has m elements it can be represented by $m!$ different sequences of length m ; thus, a relation between sets of sizes m and n admits as many as $m! * n!$ (potentially different) adjacency matrices for its representation. Not surprisingly, if $U = V$ it is *good practice* to use one and the same element numbering for the two U 's (in $U \times U$). If $1 \leq i \leq m$ then the set $[u_i]_R$ is represented by the row with index i in the adjacency matrix, that is:

$$[u_i]_R = \{v_j \mid 1 \leq j \leq n \wedge A_R[i, j] = 1\} .$$

Similarly, for $1 \leq j \leq n$ we have:

$${}_R[v_j] = \{u_i \mid 1 \leq i \leq m \wedge A_R[i, j] = 1\} .$$

1.3 Examples.

- (a) An adjacency matrix for the relation $\{(a, p), (b, p), (b, q), (c, q)\}$ from $\{a, b, c\}$ to $\{p, q\}$ is:

$$\begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{pmatrix} .$$

- (b) Another adjacency matrix for the same relation and the same sets is obtained by reversing the order of the elements in one set: if we take (c, b, a) instead of (a, b, c) and if we keep (p, q) (as above), then the adjacency matrix becomes:

$$\begin{pmatrix} 0 & 1 \\ 1 & 1 \\ 1 & 0 \end{pmatrix} .$$

Note that standard set notation is *overspecific*, as the order of the elements in an expression like $\{a, b, c\}$ is irrelevant: $\{a, b, c\}$ and $\{c, b, a\}$ are *the same* set! Therefore, when we decide to represent a relation by an adjacency matrix we need not take the order of the set's elements for granted: we really have quite some freedom here.

- (c) An adjacency matrix for the identity relation on a set of size n is the $n \times n$ *identity matrix* I_n :

$$I_n = \begin{pmatrix} 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1 \end{pmatrix} .$$

This matrix is unique, that is, independent of how the elements of the set are ordered, provided we stick to the convention of good practice, that both occurrences of the same set are ordered in the same way.

- (d) An adjacency matrix of relation \leq on the set $\{1, 2, 3, 4, 5\}$ is the upper triangular matrix

$$\begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} .$$

* * *

Some relations have special properties, which deserve to be named.

1.4 Definition. Let R be a relation on a set U . Then R is called:

- *reflexive*, if for all $x \in U$ we have: xRx ;
- *irreflexive*, if for all $x \in U$ we have: $\neg(xRx)$;
- *symmetric*, if for all $x, y \in U$ we have: $xRy \Leftrightarrow yRx$;
- *antisymmetric*, if for all $x, y \in U$ we have: $xRy \wedge yRx \Rightarrow x = y$;
- *transitive*, if for all $x, y, z \in U$ we have: $xRy \wedge yRz \Rightarrow xRz$.

1.5 Examples. We consider some of the examples given earlier:

- (a) “Is the mother of” is a relation on the set of all people. It is irreflexive, antisymmetric, and not transitive.
- (b) “There is a train connection between” is symmetric and transitive. If one is willing to accept travelling over a zero distance as a train connection, then this relation also is reflexive.
- (c) On every set relation “equals” ($=$) is reflexive, symmetric, and transitive.
- (d) Relation “divides” ($|$) is reflexive, antisymmetric, and transitive.
- (e) “Less than” ($<$) and “greater than” ($>$) on \mathbb{R} are irreflexive, antisymmetric, and transitive, whereas “at most” (\leq) and “at least” (\geq) are reflexive, antisymmetric, and transitive.
- (f) The relation $\{(x, y) \in \mathbb{R}^2 \mid y = x^2\}$ is neither reflexive nor irreflexive.

□

For any relation R the proposition $(\forall x, y : x, y \in U : xRy \Leftrightarrow yRx)$ is (logically) equivalent to the proposition $(\forall x, y : x, y \in U : xRy \Rightarrow yRx)$, which is (formally) weaker. Hence, relation R is symmetric if $xRy \Rightarrow yRx$, for all $x, y \in U$. To *prove* that R is symmetric, therefore, it suffices to prove the latter, weaker, version of the proposition, whereas to *use* (in other proofs) that R is symmetric we may use the stronger version.

1.6 Lemma. Every reflexive relation R on set U satisfies: $u \in [u]_R$, for all $u \in U$.

Proof. By calculation:

$$\begin{aligned}
 & u \in [u]_R \\
 \Leftrightarrow & \quad \{ \text{definition of } [u]_R \} \\
 & uRu \\
 \Leftrightarrow & \quad \{ R \text{ is reflexive} \} \\
 & \text{true}
 \end{aligned}$$

□

1.7 Lemma. Every symmetric relation R on set U satisfies: $v \in [u]_R \Leftrightarrow u \in [v]_R$, for all $u, v \in U$.

Proof. By calculation:

$$\begin{aligned} & v \in [u]_R \\ \Leftrightarrow & \quad \{ \text{definition of } [u]_R \} \\ & u R v \\ \Leftrightarrow & \quad \{ R \text{ is symmetric} \} \\ & v R u \\ \Leftrightarrow & \quad \{ \text{definition of } [v]_R \} \\ & u \in [v]_R \end{aligned}$$

□

* * *

If R is a relation on a finite set S , then special properties like reflexivity, symmetry and transitivity can be read off from the adjacency matrix. For example, R is reflexive if and only if the main diagonal of R 's adjacency matrix contains 1's only, that is if $A_R[i, i] = 1$ for all (relevant) i .

Relation R is symmetric if and only if the transposed matrix A_R^T equals A_R . The *transposed matrix* M^T of an $m \times n$ matrix M is the $n \times m$ matrix defined by, for all i, j :

$$M^T[j, i] = M[i, j] .$$

1.2 Equivalence relations

Relations that are reflexive, symmetric, and transitive deserve some special attention: they are called *equivalence relations*.

1.8 Definition. A relation R is an *equivalence relation* if and only if it is reflexive, symmetric, and transitive.

□

If elements u and v are related by an equivalence relation R , that is, if $u R v$, then u and v are also called “equivalent (under R)”.

1.9 Example. On every set relation “equals” ($=$) is an equivalence relation.

1.10 Example. Consider the plane \mathbb{R}^2 and in it the set S of straight lines. We call two lines in S parallel if and only if they are equal or do not intersect. Hence, two lines in S are parallel if and only if their slopes are equal. Being parallel is an equivalence relation on the set S .

1.11 Example. We consider a fixed $d \in \mathbb{Z}$, $d > 0$, and we define a relation R on \mathbb{Z} by: mRn if and only if $m-n$ is divisible by d . The latter can be formulated as $(m-n) \bmod d = 0$, and a more traditional mathematical rendering of this is $m = n \pmod{d}$. Thus defined, R is an equivalence relation.
□

Actually, the last two examples are instances of the following, rather general, theorem.

1.12 Theorem. We consider a (fixed) function f of type $U \rightarrow V$, for some sets U and V . Then the relation, on U , “having the same function value” is an equivalence relation. This is the relation \sim defined by:

$$x \sim y \Leftrightarrow f(x) = f(y) \text{ , for all } x, y \in U \text{ .}$$

□

1.13 Example. We reconsider Example 1.11. The predicate $(m-n) \bmod d = 0$ is equivalent to $m \bmod d = n \bmod d$, so with \mathbb{Z} both for set U and for set V , function f , defined by $f(m) = m \bmod d$, for all $m \in \mathbb{Z}$, does the job.
□

If R is an equivalence relation on set U , then, for every $u \in U$ the set $[u]_R$ is called the *equivalence class of u* . Because equivalence relations are reflexive we have, as we have seen in lemma 1.6: $u \in [u]_R$, for all $u \in U$. From this it follows immediately that equivalence classes are nonempty. Equivalence classes have several other interesting properties. For example, the equivalence classes of two elements are equal if and only if these elements are equivalent:

1.14 Lemma. Every equivalence relation R on set U satisfies, for all $u, v \in U$:

$$[u]_R = [v]_R \Leftrightarrow u R v \text{ .}$$

Proof. The left-hand side of this equivalence contains the function $[\cdot]_R$, whereas the right-hand side does not. To eliminate $[\cdot]_R$ we rewrite the left-hand side first:

$$\begin{aligned} & [u]_R = [v]_R \\ \Leftrightarrow & \quad \{ \text{set equality} \} \\ & (\forall x : x \in U : x \in [u]_R \Leftrightarrow x \in [v]_R) \\ \Leftrightarrow & \quad \{ \text{definition of } [\cdot]_R \} \\ & (\forall x : x \in U : u R x \Leftrightarrow v R x) \text{ ,} \end{aligned}$$

hence, the lemma is equivalent to:

$$(\forall x : x \in U : u R x \Leftrightarrow v R x) \Leftrightarrow u R v \text{ .}$$

This we prove by mutual implication.

$$\begin{aligned}
\text{"}\Rightarrow\text{"}: & \quad (\forall x : x \in U : uRx \Leftrightarrow vRx) \\
& \Rightarrow \quad \{ \text{instantiation } x := v \} \\
& \quad uRv \Leftrightarrow vRv \\
& \Leftrightarrow \quad \{ R \text{ is an equivalence relation, so it is reflexive} \} \\
& \quad uRv \text{ .}
\end{aligned}$$

" \Leftarrow ": Assuming uRv and for any $x \in U$ we prove $uRx \Leftrightarrow vRx$, again by mutual implication:

$$\begin{aligned}
& uRx \\
& \Leftrightarrow \quad \{ \text{assumption} \} \\
& \quad uRv \wedge uRx \\
& \Leftrightarrow \quad \{ R \text{ is an equivalence relation, so it is symmetric} \} \\
& \quad vRu \wedge uRx \\
& \Rightarrow \quad \{ R \text{ is an equivalence relation, so it is transitive} \} \\
& \quad vRx \text{ ,}
\end{aligned}$$

and:

$$\begin{aligned}
& vRx \\
& \Leftrightarrow \quad \{ \text{assumption} \} \\
& \quad uRv \wedge vRx \\
& \Rightarrow \quad \{ R \text{ is an equivalence relation, so it is transitive} \} \\
& \quad uRx \text{ ,}
\end{aligned}$$

which concludes the proof of this lemma.

□

Furthermore, equivalence classes are either disjoint or equal:

1.15 Lemma. Every equivalence relation R on set U satisfies, for all $u, v \in U$:

$$[u]_R \cap [v]_R = \emptyset \vee [u]_R = [v]_R \text{ .}$$

Proof. This proposition is equivalent to:

$$[u]_R \cap [v]_R \neq \emptyset \Rightarrow [u]_R = [v]_R \text{ ,}$$

which we prove as follows:

$$\begin{aligned}
& [u]_R \cap [v]_R \neq \emptyset \\
& \Leftrightarrow \quad \{ \text{definition of } \emptyset \text{ and } \cap \} \\
& \quad (\exists x : x \in U : x \in [u]_R \wedge x \in [v]_R) \\
& \Leftrightarrow \quad \{ \text{definition of } [\cdot]_R \}
\end{aligned}$$

$$\begin{aligned}
& (\exists x : x \in U : u R x \wedge v R x) \\
\Rightarrow & \quad \{ R \text{ is symmetric and transitive} \} \\
& (\exists x : x \in U : u R v) \\
\Rightarrow & \quad \{ \text{predicate calculus} \} \\
& u R v \\
\Leftrightarrow & \quad \{ \text{lemma 1.14} \} \\
& [u]_R = [v]_R
\end{aligned}$$

□

The equivalence classes of an equivalence relation “cover” the set:

1.16 Lemma. Every equivalence relation R on set U satisfies: $(\bigcup_{u:u \in U} [u]_R) = U$.

Proof. By mutual set inclusion. On the one hand, every equivalence class is a subset of U , that is: $[u]_R \subseteq U$, for all $u \in U$; hence, their union, $(\bigcup_{u:u \in U} [u]_R)$, is a subset of U as well. On the other hand, we have for every $v \in U$ that $v \in [v]_R$, so, also $v \in (\bigcup_{u:u \in U} [u]_R)$. Hence, U is a subset of $(\bigcup_{u:u \in U} [u]_R)$ too.

□

These lemmata show that the equivalence classes of an equivalence relation form a, so-called, *partition* of set U .

1.17 Definition. A partition of set U is a set Π of nonempty and disjoint subsets of U , the union of which equals U . Formally, that set Π is a partition of U means the conjunction of:

- (a) $(\forall X : X \in \Pi : X \subseteq U \wedge X \neq \emptyset)$
- (b) $(\forall X, Y : X, Y \in \Pi : X \cap Y = \emptyset \vee X = Y)$
- (c) $(\bigcup_{X: X \in \Pi} X) = U$

□

Clause (a) in this definition expresses that the elements of a partition of U are nonempty subsets of U , clause (b) expresses that the sets in a partition are disjoint, whereas clause (c) expresses that the sets in a partition together “cover the whole” U . Phrased differently, clause (b) and (c) together express that every element of U is an element of *exactly one* of the sets in the partition.

Conversely, every partition also represents an equivalence relation. Every element of set U is element of exactly one of the subsets in the partition. “Being in the same subset” (in the partition) is an equivalence relation.

1.18 Theorem. Every partition Π of a set U represents an equivalence relation on U , the equivalence classes of which are the sets in Π .

Proof. Because Π is a partition, every element of U is an element of a unique subset

in Π . Now, the relation “being elements of the same subset in Π ” is an equivalence relation. Formally, we prove this by defining a function $\varphi: U \rightarrow \Pi$, as follows, for all $u \in U$ and $X \in \Pi$:

$$\varphi(u) = X \Leftrightarrow u \in X \text{ .}$$

Thus defined, φ is a function indeed, because for every $u \in U$ one and only one $X \in \Pi$ exists satisfying $u \in X$. Now relation \sim on U , defined by, for all $u, v \in U$:

$$u \sim v \Leftrightarrow \varphi(u) = \varphi(v) \text{ ,}$$

is an equivalence relation – Theorem 1.12!–. Furthermore, by its very construction φ satisfies $u \in \varphi(u)$ and, hence, $\varphi(u)$ is the equivalence class of u , for all $u \in U$.

□

1.3 Operations on Relations

Relations between two sets are subsets of the Cartesian Product of these two sets. Hence, all usual set operations can be applied to relations as well. In addition, relations admit of some dedicated operations that happen to have nice algebraic properties. It is even possible to develop a viable Relational Calculus, but this falls outside the scope of this text.

These relational operations play an important role in the mathematical study of programming constructs, such as recursion and data structures. They are also useful in some theorems about graphs. We will see applications of this later.

1.3.1 Set operations

- For sets U and V , the *extreme relations* from U to V are the *empty relation* \emptyset and the *full relation* $U \times V$. For the sake of brevity and symmetry, we denote these two relations by \perp (“bottom”) and \top (“top”), respectively; element wise, they satisfy, for all $u \in U$ and $v \in V$:

$$\neg(u \perp v) \wedge u \top v \text{ .}$$

For example, every relation R satisfies: $\perp \subseteq R$ and $R \subseteq \top$, which is why we call \perp and \top the extreme relations.

- If R and S are relations, with the same domain and with the same range, then $R \cup S$, and $R \cap S$, and $R \setminus S$ are relations too, between the same sets as R and S , and with the obvious meaning. The complement R^C of relation R is $\top \setminus R$.
- These operations have their usual algebraic properties. In particular, \top and \perp are the identity elements of \cup and \cap , respectively: $R \cup \perp = R$ and $R \cap \top = R$. They are zero elements as well, that is: $R \cup \top = \top$ and $R \cap \perp = \perp$.

1.3.2 Transposition

With every relation R from set U to set V a corresponding relation exists from V to U that contains (v, u) if and only if $(u, v) \in R$. This relation is called the *transposition* of R and is denoted by R^T . (Some mathematicians use R^{-1} , but this may be confusing: transposition is not the same as inversion, especially with functions.) Formally, transposition is defined as follows.

1.19 Definition. For every relation R from set U to set V , relation R^T from V to U is defined by, for all $v \in V$ and $u \in U$:

$$v R^T u \Leftrightarrow u R v .$$

1.20 Lemma. Transposition distributes over all set operations, that is:

$$\begin{aligned} \perp^T &= \perp \quad \text{and:} \quad \top^T = \top ; \\ (R \cup S)^T &= R^T \cup S^T ; \\ (R \cap S)^T &= R^T \cap S^T ; \\ (R \setminus S)^T &= R^T \setminus S^T ; \\ (R^C)^T &= (R^T)^C . \end{aligned}$$

1.21 Lemma. Transposition is its own *inverse*, that is, every relation R satisfies:

$$(R^T)^T = R .$$

□

For finite relations there is a direct connection between relation transposition and matrix transposition:

1.22 Lemma. If A_R is an adjacency matrix for relation R then $(A_R)^T$ is an adjacency matrix for R^T .

1.23 Examples. Properties of relations, like (ir)reflexivity and (anti)symmetry, can now be expressed concisely by means of relational operations; for R a relation on set U :

- “ R is reflexive” $\Leftrightarrow I_U \subseteq R$
- “ R is irreflexive” $\Leftrightarrow I_U \cap R = \perp$
- “ R is symmetric” $\Leftrightarrow R^T = R$
- “ R is antisymmetric” $\Leftrightarrow R \cap R^T \subseteq I_U$

□

Unfortunately, transitivity cannot be expressed so nicely in terms of the set operations. For this we need yet another operation on relations, which turns out to be quite useful for other purposes too.

1.3.3 Composition

Let R be a relation from U to V and let S be a relation from V to W . If uRv , for some $v \in V$ and if vSw , for that *same* v , then we say that u is related to w in the *composition* of R and S , written as $R;S$. So, the composition of R and S is a relation from U to W . Phrased differently, in this composition $u \in U$ is related to $w \in W$ if u and w are “connected via” some “intermediate” value in V . This is rendered formally as follows.

1.24 Definition. If R is a relation from U to V , and if S is a relation from V to W , then the composition $R;S$ is the relation from U to W defined by, for all $u \in U$ and $w \in W$:

$$u(R;S)w \Leftrightarrow (\exists v: v \in V: uRv \wedge vSw) .$$

1.25 Example. Let $R = \{(1, 2), (2, 3), (2, 4), (3, 1), (3, 3)\}$ be a relation from $\{1, 2, 3\}$ to $\{1, 2, 3, 4\}$ and let $S = \{(1, a), (2, c), (3, a), (3, d), (4, b)\}$ be a relation from $\{1, 2, 3, 4\}$ to $\{a, b, c, d\}$. Then the composition $R;S$ is the relation $\{(1, c), (2, a), (2, b), (2, d), (3, a), (3, d)\}$, from $\{1, 2, 3\}$ to $\{a, b, c, d\}$.

1.26 Lemma. Now we have, for endorelation R :

$$\text{“}R \text{ is transitive”} \Leftrightarrow (R;R) \subseteq R .$$

Proof. Left as an exercise.

□

1.27 Lemma. The identity relation is the identity of relation composition. More precisely, every relation R from set U to set V satisfies: $I_U;R = R$ and $R;I_V = R$.

Proof. Left as an exercise.

□

1.28 Lemma. Relation composition is associative, that is, all relations R, S, T satisfy: $(R;S);T = R;(S;T)$.

Proof. For all u, x we calculate:

$$\begin{aligned} & u((R;S);T)x \\ \Leftrightarrow & \quad \{ \text{definition of } ; \} \\ & (\exists w:: u(R;S)w \wedge wTx) \\ \Leftrightarrow & \quad \{ \text{definition of } ; \} \\ & (\exists w:: (\exists v:: uRv \wedge vSw) \wedge wTx) \\ \Leftrightarrow & \quad \{ \wedge \text{ over } \exists \} \\ & (\exists w:: (\exists v:: uRv \wedge vSw \wedge wTx)) \\ \Leftrightarrow & \quad \{ \text{swapping dummies} \} \\ & (\exists v:: (\exists w:: uRv \wedge vSw \wedge wTx)) \\ \Leftrightarrow & \quad \{ \text{(almost) the same steps as above, in reverse order} \} \\ & u(R;(S;T))x \end{aligned}$$

□

Remark: In other mathematical texts relation composition is sometimes called “(relational) product”, denoted by infix operator $*$. From a formal point of view, this is harmless, of course, but it is important to keep in mind that composition is *not commutative*: generally, $R;S$ differs from $S;R$! This is the reason why we prefer to use an asymmetric symbol, “;”, to denote composition: from a practical point of view the term “product” and the symbol “ $*$ ” may be misleading.

□

A very important property is that relation composition *distributes over arbitrary unions* of relations, both from the left and from the right:

1.29 Theorem. Every relation R and every collection Ω of relations satisfies:

$$R; \left(\bigcup_{X: X \in \Omega} X \right) = \left(\bigcup_{X: X \in \Omega} R; X \right) ,$$

and also:

$$\left(\bigcup_{X: X \in \Omega} X \right); R = \left(\bigcup_{X: X \in \Omega} X; R \right) .$$

Proof. Left as an exercise.

□

Corollary: Relation composition is *monotonic*, that is, for all relations R, S, T :

$$S \subseteq T \Rightarrow R; S \subseteq R; T , \text{ and also:}$$

$$R \subseteq S \Rightarrow R; T \subseteq S; T .$$

□

* * *

The n -fold composition of a relation R with itself is sometimes written as R^n , for natural n . More precisely, for all n , $0 \leq n$, we define (recursively):

$$R^0 = I \wedge R^{n+1} = R; R^n .$$

For example, the formula expressing transitivity of R , as in Lemma 1.26, can now also be written as: $R^2 \subseteq R$.

* * *

In the representation of relations by adjacency matrices, relation composition is represented by matrix multiplication. That is, if A_R is an adjacency matrix for relation R and if A_S is an adjacency matrix for relation S then the product matrix $A_R \times A_S$ is an adjacency matrix for the composition $R;S$. This matrix product is well-defined only if the number of columns of matrix A_R equals the number of rows of matrix

A_S . This is true because the number of columns of A_R equals the size of the range of relation R . As this range also is the domain of relation S – otherwise composition of R and S is impossible – this size also equals the number of rows of A_S .

Recall that adjacency matrices actually are boolean matrices; hence, the matrix multiplication must be performed with boolean operations, not integer operations, in such a way that addition and multiplication boil down to disjunction (“or”) and conjunction (“and”) respectively. So, a formula like $(\Sigma j :: A_R[i, j] * A_S[j, k])$ actually becomes: $(\exists j :: A_R[i, j] \wedge A_S[j, k])$.

1.30 Example. Let $R = \{(1, 2), (2, 3), (2, 4), (3, 1), (3, 3)\}$ be a relation from $\{1, 2, 3\}$ to $\{1, 2, 3, 4\}$ and let $S = \{(1, a), (2, c), (3, a), (3, d), (4, b)\}$ be a relation from $\{1, 2, 3, 4\}$ to $\{a, b, c, d\}$. Then adjacency matrices for R and S are:

$$\begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix}, \text{ and: } \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{pmatrix}.$$

The product of these matrices is an adjacency matrix for $R;S$:

$$\begin{pmatrix} 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{pmatrix}.$$

1.3.4 Closures

Some relations have properties, like reflexivity, symmetry, or transitivity, whereas other relations do not. For any such property, the *closure* of a relation with respect to that property is the *smallest extension* of the relation that does have the property. More precisely, if \mathcal{P} is a predicate on relations, then the \mathcal{P} -closure of relation R is the relation S , say, satisfying the following three requirements:

$$R \subseteq S,$$

which expresses that S is an extension of R , and:

$$\mathcal{P}(S),$$

which expresses that S has property \mathcal{P} , and:

$$R \subseteq X \wedge \mathcal{P}(X) \Rightarrow S \subseteq X,$$
 for all relations X ,

which expresses that S is contained in every relation X that is an extension of R and that has property \mathcal{P} ; this is what we mean with the *smallest extension* of R .

For any given property \mathcal{P} and relation R the \mathcal{P} -closure of R need not exist, but if it exists it is unique. Moreover, if relation R already has property \mathcal{P} , so $\mathcal{P}(R)$ holds, then, of course, the \mathcal{P} -closure of R is R itself.

remark: In this subsection we are studying properties of the general shape $\phi(X) \subseteq X$, where ϕ is a *monotonic* function from relations to relations, and where parameter X is a relation. In a later chapter, on Partial Orders, we will see that for monotonic functions ϕ , the smallest relation having such a property is the *intersection* of all relations having that property, that is: $(\bigcap_{X:\phi(X)\subseteq X} X)$. We will also see that, under some additional conditions, this smallest relation is equal to the union of all “approximations from below”, that is: $(\bigcup_{i:0\leq i} \phi^i(\perp))$. As these are rather general properties, which are not specific to closures of relations, we will not elaborate this here.

□

* * *

The simplest possible property of relations is reflexivity. The reflexive closure of an (endo)relation R now is the smallest extension of R that is reflexive. So, the reflexive closure of R is the smallest of all relations X satisfying:

$$R \subseteq X \wedge I \subseteq X .$$

This formula is logically equivalent to:

$$R \cup I \subseteq X ,$$

and, of course, the smallest relation X satisfying this is: $R \cup I$. So, the reflexive closure of relation R equals $R \cup I$. Notice that if R itself is reflexive, which means $I \subseteq R$, then $R \cup I$ equals R and we obtain R itself as its reflexive closure, as we should.

* * *

In the same vein we may ask for the *symmetric closure* of (endo)relation R . This time, it is defined as the smallest of all relations X satisfying:

$$R \subseteq X \wedge X^T \subseteq X ,$$

which, again, is logically equivalent to: $R \cup X^T \subseteq X$.

Although this equation is recursive, it is not very difficult to prove that its smallest solution is: $R \cup R^T$. (This is left to the exercises.)

* * *

The game becomes more interesting when we ask for the *transitive closure* of relation R . The transitive closure of R is the smallest of all relations X that contain R and that are transitive; that is, formally, it is the smallest of all relations X satisfying:

$$R \subseteq X \wedge X ; X \subseteq X .$$

1.31 Theorem. For every endorelation R its transitive closure is also the smallest of all relations X satisfying:

$$(0) \quad R \subseteq X \wedge R; X \subseteq X \text{ ,}$$

and it also is the smallest of all relations X satisfying:

$$(1) \quad R \subseteq X \wedge X; R \subseteq X \text{ .}$$

Proof. We prove (0), as follows; the proof of (1) is, mutatis mutandis, exactly the same. Let S be the transitive closure of R ; by definition, this means that S has the following three properties:

$$(2) \quad R \subseteq S \text{ , and:}$$

$$(3) \quad S; S \subseteq S \text{ , and:}$$

$$(4) \quad R \subseteq X \wedge X; X \subseteq X \Rightarrow S \subseteq X \text{ , for all relations } X \text{ .}$$

Also let T be the smallest of all predicates X satisfying (0); this means:

$$(5) \quad R \subseteq T \text{ , and:}$$

$$(6) \quad R; T \subseteq T \text{ , and:}$$

$$(7) \quad R \subseteq X \wedge R; X \subseteq X \Rightarrow T \subseteq X \text{ , for all relations } X \text{ .}$$

Now we must prove that $S = T$, which we prove by mutual set inclusion – what else, in view of the above formulae, can we do?–:

$$\begin{aligned} & T \subseteq S \\ \Leftrightarrow & \{ (7) \text{ , with } X := S \} \\ & R \subseteq S \wedge R; S \subseteq S \\ \Leftrightarrow & \{ (2) \} \\ & R; S \subseteq S \\ \Leftrightarrow & \{ \subseteq \text{ is transitive, to prepare for use of (3)} \} \\ & R; S \subseteq S; S \wedge S; S \subseteq S \\ \Leftrightarrow & \{ (3) \} \\ & R; S \subseteq S; S \\ \Leftrightarrow & \{ ; \text{ is monotonic} \} \\ & R \subseteq S \\ \Leftrightarrow & \{ (2) \} \\ & \text{true ,} \end{aligned}$$

and:

$$\begin{aligned}
& S \subseteq T \\
\Leftarrow & \quad \{ (4) , \text{ with } X := T \} \\
& R \subseteq T \wedge T ; T \subseteq T \\
\Leftrightarrow & \quad \{ (5) \} \\
& T ; T \subseteq T .
\end{aligned}$$

Thus far, the proof is entirely of the kind “nothing-else-one-can-do”, but now we are stuck. What remains to be proved is $T ; T \subseteq T$, and to do so in as simple a relational way as above we would need some (elementary) knowledge of (so-called) *Galois connections*. Fortunately, there is a simple way out. The next theorem provides an explicit formula for T , and using this we prove $T ; T \subseteq T$ as follows:

$$\begin{aligned}
& T ; T \\
= & \quad \{ \text{Theorem 1.32, see below} \} \\
& (\bigcup_{m:1 \leq m} R^m) ; (\bigcup_{n:1 \leq n} R^n) \\
= & \quad \{ \text{Theorem 1.29: ; distributes over arbitrary unions} \} \\
& (\bigcup_{m,n:1 \leq m \wedge 1 \leq n} R^{m+n}) \\
= & \quad \{ \text{dummy transformation and simplification} \} \\
& (\bigcup_{k:2 \leq k} R^k) \\
\subseteq & \quad \{ \text{domain extension} \} \\
& (\bigcup_{k:1 \leq k} R^k) \\
= & \quad \{ \text{Theorem 1.32} \} \\
& T .
\end{aligned}$$

□

1.32 Theorem. The smallest of all relations X satisfying (0) (in Theorem 1.31) is $(\bigcup_{n:1 \leq n} R^n)$.

Proof. We must prove that Z , with $Z = (\bigcup_{n:1 \leq n} R^n)$, is the smallest of all relations X satisfying $R \subseteq X \wedge R ; X \subseteq X$, or, equivalently, $R \cup R ; X \subseteq X$. This means that we must prove two things: firstly, that Z satisfies this equation, and, secondly, that $Z \subseteq X$ for every X satisfying this equation. As for the first, we calculate:

$$\begin{aligned}
& R \cup R ; Z \\
= & \quad \{ \text{definition of } Z \} \\
& R \cup R ; (\bigcup_{n:1 \leq n} R^n) \\
= & \quad \{ ; \text{ distributes over } \cup \} \\
& R \cup (\bigcup_{n:1 \leq n} R ; R^n) \\
= & \quad \{ \text{definition of } R^{n+1} \} \\
& R \cup (\bigcup_{n:1 \leq n} R^{n+1}) \\
= & \quad \{ \text{dummy transformation } n := n-1 \}
\end{aligned}$$

$$\begin{aligned}
& R \cup \left(\bigcup_{n:2 \leq n} R^n \right) \\
= & \quad \{ R = R^1; \text{ join } n=1 \} \\
& \left(\bigcup_{n:1 \leq n} R^n \right) \\
= & \quad \{ \text{definition of } Z \} \\
& Z .
\end{aligned}$$

Actually, we now have proved the stronger $R \cup R;Z = Z$, which implies the required $R \cup R;Z \subseteq Z$. The second proof obligation is, for all X :

$$R \cup R;X \subseteq X \Rightarrow Z \subseteq X ,$$

which we prove by assuming $R \cup R;X \subseteq X$ and, subsequently, by observing that $(\bigcup_{n:1 \leq n} R^n) \subseteq X$ is equivalent to: $(\forall n:1 \leq n: R^n \subseteq X)$. This lends itself to a proof by Mathematical Induction (on n). We recall that assumption $R \cup R;X \subseteq X$ is equivalent to the conjunction of $R \subseteq X$ and $R;X \subseteq X$: we will use whichever is convenient:

$$\begin{aligned}
& R^1 \\
= & \quad \{ \text{definition of } R^1 \} \\
& R \\
\subseteq & \quad \{ \text{assumption } R \subseteq X \} \\
& X
\end{aligned}$$

and, for $n, 1 \leq n$, and assuming $R^n \subseteq X$ –the Induction Hypothesis–:

$$\begin{aligned}
& R^{n+1} \\
= & \quad \{ \text{definition of } R^{n+1} \} \\
& R;R^n \\
\subseteq & \quad \{ \text{Induction Hypothesis, using monotonicity of } ; \} \\
& R;X \\
\subseteq & \quad \{ \text{assumption } R;X \subseteq X \} \\
& X
\end{aligned}$$

□

Corollary: We use the notation R^+ for the expression $(\bigcup_{n:1 \leq n} R^n)$. Theorem 1.32 then states that R^+ equals the smallest of all relations X satisfying (0) in Theorem 1.31, and this latter theorem states that, hence, R^+ also equals the transitive closure of R .

□

* * *

Finally, the reflexive-transitive closure of relation R is denoted by R^* ; it is the smallest extension of R into a relation that is both reflexive and transitive, so it is defined as the smallest of all relations X satisfying:

$$I \subseteq X \wedge R \subseteq X \wedge X;X \subseteq X \ .$$

For the reflexive-transitive closure similar theorems can be formulated as for the transitive closure.

1.33 Lemma. For any relation R we have:

$$R^* = I \cup R^+ \wedge R^+ = R;R^* \ .$$

Proof. Left to the exercises.

□

1.34 Theorem. For every endorelation R its reflexive-transitive closure R^* is also the smallest of all relations X satisfying:

$$I \cup R;X \subseteq X \ .$$

Similarly, R^* is also the smallest of all relations X satisfying:

$$I \cup X;R \subseteq X \ .$$

Proof. Use Lemma 1.33 and Theorem 1.31.

□

1.35 Theorem. Every endorelation R satisfies: $R^* = (\bigcup_{n:0 \leq n} R^n)$.

Proof. Use Lemma 1.33 and Theorem 1.32.

□

1.4 Warshall's Algorithm

1.4.1 Introduction

For relations on *finite* sets their transitive closure can be computed in a finite amount of time. Relations on finite sets are also known as *graphs* and Warshall's algorithm presented here generally is considered an algorithm on graphs.

Theorem 1.32 gives an explicit formula for the transitive closure of a relation R , namely:

$$R^+ = (\bigcup_{n:1 \leq n} R^n) \ .$$

If R is a relation on a finite set U with N , $1 \leq N$, elements, it can be proved that the range of the dummy in this quantification may be restricted to $n \leq N$, that is, for finite U of size N we have:

$$R^+ = (\bigcup_{n:1 \leq n \leq N} R^n) \ .$$

If relation R is represented by an adjacency matrix then an adjacency matrix for R^+ can be calculated with as few as $\mathcal{O}(\log(N))$ matrix multiplications. As every $N \times N$ matrix multiplication itself requires $\mathcal{O}(N^3)$ elementary operations, this way R^+ can be calculated in $\mathcal{O}(\log(N) * N^3)$ time.

Warshall's algorithm, however, does better: its time complexity is only $\mathcal{O}(N^3)$. This is achieved by replacing matrix multiplication by a simpler matrix operation, requiring only $\mathcal{O}(N^2)$ steps; this simpler operation is used (exactly) N times, thus giving rise to an algorithm with $\mathcal{O}(N^3)$ time complexity.

1.4.2 Preliminaries

Set U , with size N , being finite it can, without any loss of generality, be represented as the interval $[1..N]$. The crucial idea behind Warshall's algorithm is to *generalize* (relation) composition to a form which we call "restricted composition" here. It involves an additional parameter n , $0 \leq n \leq N$, say, which we add as a subscript to the symbol for composition, as follows, for all relations S, T on $[1..N]$ and for all $u, w \in [1..N]$:

$$u(S;_n T)w \Leftrightarrow (\exists v: 1 \leq v \leq n: uSv \wedge vTw) .$$

Restricted composition then has the following properties; the latter of these shows that this is a generalization indeed, as ";" appears as a special case:

$$S;_0 T = \perp , \text{ and:}$$

$$S;_N T = S;T .$$

It even is possible to derive a *recurrence relation* for ";" _{n} ; for this purpose we need a collection of auxiliary relations, which we denote by $\langle v \rangle$ ("via v "), and which we define, for all $u, v, w \in [1..N]$, by:

$$u(S\langle v \rangle T)w \Leftrightarrow uSv \wedge vTw .$$

In terms of these auxiliary relations restricted composition can now be (re)defined by:

$$u(S;_n T)w \Leftrightarrow (\exists v: 1 \leq v \leq n: u(S\langle v \rangle T)w) ,$$

which is equivalent to:

$$S;_n T = (\bigcup_{v: 1 \leq v \leq n} S\langle v \rangle T) .$$

For $n=0$ the range in this quantification is empty and, hence – as we have seen already –: $S;_0 T = \perp$. For $n \in [1..N]$ the term with $v=n$ can be split of, and we obtain:

$$S;_n T = S;_{n-1} T \cup S\langle n \rangle T .$$

1.4.3 The algorithm

We now consider a fixed relation R on $[1..N]$. We define a sequence of relations W_n , for $0 \leq n \leq N$, which serve as *approximations* of R^+ . Relation W_n is the transitive closure of R not based on ordinary composition, however, but on the restricted composition “ $;$ ” instead. That is, W_n is the smallest of all relations X satisfying:

$$R \cup X ;_n X \subseteq X .$$

Because $X ;_0 X = \perp$ we have: $W_0 = R$; because $X ;_N X = X ; X$ we have: $W_N = R^+$. Hence, if we are able to formulate a useful recurrence relation for W_n we have a basis for an algorithm.

1.36 Lemma. For all n , $1 \leq n \leq N$: $W_n = W_{n-1} \cup W_{n-1} \langle n \rangle W_{n-1}$.

Proof. Omitted

□

In terms of relational operations this recurrence relation can be implemented in an algorithm for the computation of R^+ in a straightforward way.

Warshall’s algorithm (abstract):

```

S := R
; { S = W0 }
for n := 1 to N
do { 1 ≤ n ≤ N ∧ S = Wn-1 }
    S := S ∪ S⟨n⟩S
    { 1 ≤ n ≤ N ∧ S = Wn }
od
{ n = N ∧ S = Wn , hence: S = R+ }

```

□

To implement this algorithm in terms of operations on the individual elements of $[1..N]$ we observe, for $u, w \in [1..N]$:

$$\begin{aligned}
& u(S \cup S \langle n \rangle S) w \\
\Leftrightarrow & \quad \{ \text{definition of } \cup \} \\
& uSw \vee u(S \langle n \rangle S) w \\
\Leftrightarrow & \quad \{ \text{definition of } \langle n \rangle \} \\
& uSw \vee (uSn \wedge nSw) .
\end{aligned}$$

Thus we obtain a more detailed implementation of the algorithm.

Warshall's algorithm (concrete):

```

 $S := R$ 
; {  $S = W_0$  }
for  $n := 1$  to  $N$ 
do {  $1 \leq n \leq N \wedge S = W_{n-1}$  }
   $T := S$ 
  ; {  $1 \leq n \leq N \wedge T = W_{n-1}$  }
  for all  $u, w \in [1..N]$ 
  do  $uSw := uTw \vee (uTn \wedge nTw)$ 
  od
  {  $1 \leq n \leq N \wedge S = W_n$  }
od
{  $S = R^+$  }

```

□

This algorithm has $\mathcal{O}(N^3)$ time complexity: the outermost for-construct takes (exactly) N steps, whereas the assignment $T := S$ and the innermost for-construct require (exactly) N^2 steps each. This innermost for-construct enumerates all pairs $u, w \in [1..N]$; notice that the order in which these pairs are enumerated is irrelevant.

In this algorithm a local variable, T , is used to buffer the value of W_{n-1} during the calculation of W_n in S . This local variable is unnecessary, though, because we have, for all u, w, n :

$$\begin{aligned}
& uW_n n \\
\Leftrightarrow & \quad \{ \text{Lemma 1.36, elementwise} \} \\
& uW_{n-1} n \vee (uW_{n-1} n \wedge nW_{n-1} n) \\
\Leftrightarrow & \quad \{ \text{absorption} \} \\
& uW_{n-1} n \quad ,
\end{aligned}$$

and, similarly:

$$\begin{aligned}
& nW_n w \\
\Leftrightarrow & \quad \{ \text{as above} \} \\
& nW_{n-1} w \quad .
\end{aligned}$$

This shows that both row n and column n in (the matrix for) W_n are equal to row n and column n , respectively, in W_{n-1} . Therefore, for $u=n$ or $n=w$, the assignment $uSw := uTw \vee (uTn \wedge nTw)$ boils down to **skip** (“no change of state”). For all other assignments, now T may be safely replaced by S as well; still, the order in which all pairs u, w are enumerated is irrelevant.

This yields the following, and final, encoding of the algorithm.

Warshall's algorithm (optimized):

```

 $S := R$ 
; {  $S = W_0$  }
for  $n := 1$  to  $N$ 
do {  $1 \leq n \leq N \wedge S = W_{n-1}$  }
  for all  $u, w \in [1..N]$ 
  do  $uSw := uSw \vee (uSn \wedge nSw)$ 
  od
  {  $1 \leq n \leq N \wedge S = W_n$  }
od
{  $S = R^+$  }

```

□

1.5 Exercises

1. Give an example of a relation that is:
 - (a) both reflexive and irreflexive;
 - (b) neither reflexive nor irreflexive;
 - (c) both symmetric and antisymmetric;
 - (d) neither symmetric nor antisymmetric.
2. For each of the following relations, investigate whether it is (ir)reflexive, (anti-)symmetric, and/or transitive:
 - (a) $R = \{(x, y) \in \mathbb{R}^2 \mid x+1 < y\}$
 - (b) $S = \{(x, y) \in \mathbb{R}^2 \mid x < y+1\}$
 - (c) $T = \{(x, y) \in \mathbb{Z}^2 \mid x < y+1\}$
3. Prove that each irreflexive and transitive relation is antisymmetric.
4. Let R be a relation on a set U . Prove that, if $[u]_R \neq \emptyset$, for all $u \in U$, and if R is symmetric and transitive, then R is reflexive.
5. Prove Theorem 1.12.
6. The natural numbers admit addition but not subtraction: if $a < b$ the difference $a - b$ is undefined, because it is not a natural number. To achieve a structure in which all differences are defined we need the “integer numbers”. These can be constructed from the naturals in the following way, a process called “definition by abstraction”.

We consider the set V of all pairs of natural numbers, so $V = \mathbb{N} \times \mathbb{N}$. On V we define a relation \sim , as follows, for all $a, b, c, d \in \mathbb{N}$:

$$(a, b) \sim (c, d) \Leftrightarrow a + d = c + b .$$

- (a) Prove that \sim is an equivalence relation.
- (b) Formulate in words what this equivalence relation expresses.
- (c) We investigate the equivalence classes of \sim . Obviously, there is a class containing the pair $(0, 0)$. Prove that, in addition, every *other* class contains *exactly one* pair of the shape either $(a, 0)$ or $(0, a)$, so not both in the same class, with $1 \leq a$.
- (d) We call the pairs $(0, 0)$, $(a, 0)$ and $(0, a)$, with $1 \leq a$, the “representants” of the equivalence classes. These classes can now be ordered in the following way, by means of their representants:

$$\dots, (0, 2), (0, 1), (0, 0), (1, 0), (2, 0), (3, 0), \dots .$$

We call these classes “integer numbers”; a more usual notation of the representants is:

$$\dots, -2, -1, 0, +1, +2, +3, \dots .$$

Thus, we obtain the integer numbers indeed. To illustrate this: define, on the set of representants, two binary operators **pls** and **min** that correspond with the usual “addition” and “subtraction”. Also define the “less than” relation on the set of representants.

7. Laat \mathcal{D} be the set of differentiable functions $f: \mathbb{R} \rightarrow \mathbb{R}$. On \mathcal{D} we define a relation \sim as follows, for all $f, g \in \mathcal{D}$: We definiëren op \mathcal{D} een relatie R door

$$f \sim g \Leftrightarrow \text{“function } f - g \text{ is constant”} .$$

Prove that \sim is an equivalence relation. How can relation \sim be defined in the way of Theorem 1.12?

8. We consider a linear vector space V and a (fixed) subspace W of V . On V we define a relation \sim by, for all $x, y \in V$:

$$x \sim y \Leftrightarrow x - y \in W .$$

Prove that \sim is an equivalence relation. Describe the equivalence classes for the special case that $V = \mathbb{R}^2$ W is the straight line given by the equation $x_1 + x_2 = 0$. Also characterise, for this special case, the equivalence relation in the way of Theorem 1.12.

9. An adjacency matrix for a relation R is: $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$. Investigate whether R is (ir)reflexive, (anti)symmetric, and/or transitive.
10. Prove that $(R; S)^T = S^T; R^T$.

11. Prove Lemma 1.26.
12. (a) Prove that, for all sets A, B, C : $A \subseteq C \wedge B \subseteq C \Leftrightarrow A \cup B \subseteq C$.
 (b) Prove that, for all sets A, B : $A \subseteq B \Leftrightarrow A \cup B = B$ and also:
 $A \subseteq B \Leftrightarrow A \cap B = A$.
 (c) Prove that relation composition distributes over union, that is:
 $R; (S \cup T) = (R; S) \cup (R; T)$ and: $(R \cup S); T = (R; T) \cup (S; T)$.
 (d) Using the previous result(s), prove that $;$ is *monotonic*, that is:
 $S \subseteq T \Rightarrow R; S \subseteq R; T$ and also: $R \subseteq S \Rightarrow R; T \subseteq S; T$.
13. Prove that, for endorelation R and for all natural m and n :
 $R^{m+n} = R^m; R^n$.
14. Prove that, indeed, $R \cup R^T$ is the smallest solution of equation, with unknown X : $R \subseteq X \wedge X^T \subseteq X$.
15. Prove that $R; \perp = \perp$, for every relation R .
16. Prove that \top is a solution of each of the equations (with unknown X) in Subsection 1.3.4.
17. We consider a relation R from U to V for which it is given that it is a function. Prove that R is surjective if and only if $I_V = R^T; R$.
18. We investigate some well-known relations on \mathbb{R} :
 - (a) What is the reflexive closure of $<$?
 - (b) What is the symmetric closure of $<$?
 - (c) What is the symmetric closure of \leq ?
 - (d) What is the reflexive closure of \neq ?
19. Relation R , on \mathbb{Z} , is defined by $mRn \Leftrightarrow m+1 = n$, for all $m, n \in \mathbb{Z}$. What is relation R^+ ?
20. For some given set Ω , a function ϕ , mapping subsets of Ω to subsets of Ω , is called *monotonic* if $X \subseteq Y \Rightarrow \phi(X) \subseteq \phi(Y)$, for all $X, Y \subseteq \Omega$.
 - (a) We consider the equations: $X: \phi(X) \subseteq X$ and: $X: \phi(X) = X$, and we *assume* they have smallest solutions; so, proving the *existence* of these smallest solutions is not the subject of this exercise. Prove that, if ϕ is monotonic then the smallest solutions of these equations are equal.
 - (b) For each of the closures in Subsection 1.3.4, define a function ϕ such that the corresponding equation is equivalent to $\phi(X) \subseteq X$; for each case, prove that ϕ is monotonic. What is Ω in these cases?
21. Prove that $R^* = I \cup R^+$ and that $R^+ = R; R^*$.

22. Prove that for every endorelation R : “ R is transitive” $\Leftrightarrow R^+ = R$.
23. We consider two endorelations R and S satisfying $R;S \subseteq S;R^+$. Prove that: $R^+;S \subseteq S;R^+$.
24. (a) Let R be an endorelation and let S be a transitive relation. Prove that:

$$R \subseteq S \Rightarrow R^+ \subseteq S \text{ .}$$

- (b) Apply this, by defining suitable relations R and S , to prove that every function f on \mathbb{N} satisfies:

$$(\forall i: 0 \leq i < n: f_i = f_{i+1}) \Rightarrow f_0 = f_n \text{ , for all } n \in \mathbb{N} \text{ .}$$

25. We call a relation on a set *inductive* if it admits proofs by Mathematical Induction. Formally, a relation R on a set V is inductive if, for every predicate P on V :

$$(\forall v: v \in V: (\forall u: u R v: P(u)) \Rightarrow P(v)) \Rightarrow (\forall v: v \in V: P(v)) \text{ .}$$

Prove that, for every relation R :

$$\text{“}R \text{ is inductive”} \Rightarrow \text{“}R^+ \text{ is inductive”} \text{ .}$$

Hint: To prove the right-hand side of this implication one probably will introduce a predicate P . To apply the (assumed) left-hand side of the implication one may select any predicate desired, not necessarily P : use predicate Q defined by, for all $v \in V$: $Q(v) = (\forall u: u R^* v: P(u))$.

26. On the natural numbers a distinction is often made between (so-called) “weak” (or “step-by-step”) induction and “strong” (or “course-of-values”) induction. Weak induction is the property that, for every predicate P on \mathbb{N} :

$$P(0) \wedge (\forall n:: P(n) \Rightarrow P(n+1)) \Rightarrow (\forall n:: P(n)) \text{ ,}$$

whereas strong induction is the property that, for every predicate P on \mathbb{N} :

$$(\forall n:: (\forall m: m < n: P(m)) \Rightarrow P(n)) \Rightarrow (\forall n:: P(n)) \text{ ,}$$

Show that the proposition:

$$\text{“weak induction”} \Rightarrow \text{“strong induction”} \text{ ,}$$

is a special case of the proposition in the previous exercise.

27. Suppose that endorelation R satisfies $I \cap R^+ = \perp$. What does this mean?
28. Which of the following relations on set U , with $U = \{1, 2, 3, 4\}$, is reflexive, irreflexive, symmetric, antisymmetric, or transitive?

- (a) $\{(1, 3), (2, 4), (3, 1), (4, 2)\}$;
 - (b) $\{(1, 3), (2, 4)\}$;
 - (c) $\{(1, 1), (2, 2), (3, 3), (4, 4), (1, 3), (2, 4), (3, 1), (4, 2)\}$;
 - (d) $\{(1, 1), (2, 2), (3, 3), (4, 4)\}$;
 - (e) $\{(1, 1), (2, 2), (3, 3), (4, 4), (1, 2), (2, 3), (3, 4), (4, 3), (3, 2), (2, 1)\}$.
29. Construct an example, as simple as possible, illustrating that relation composition is not *commutative*, which means that it is *not* true that: $R;S = S;R$, for all relations R, S .
30. Construct for each of the relations in Exercise 20 an adjacency matrix.
31. Construct for each relation, named R here, in Exercise 20 an adjacency matrix for R^2 .
32. Compute for each of the relations in Exercise 20 their reflexive, symmetric, and transitive closures.
33. Prove that every reflexive and transitive endorelation R satisfies: $R^2 = R$.
34. Suppose R and S are finite relations with adjacency matrices A and B , respectively. Define adjacency matrices, in terms of A and B , for the relations $R \cup S$, $R \cap S$, $R \setminus S$, and R^C .
35. Suppose R and S are endorelations. Prove or disprove:
- (a) If R and S are reflexive, then so is $R;S$.
 - (b) If R and S are irreflexive, then so is $R;S$.
 - (c) If R and S are symmetric, then so is $R;S$.
 - (d) If R and S are antisymmetric, then so is $R;S$.
 - (e) If R and S are transitive, then so is $R;S$.
 - (f) If R and S are equivalence relations, then so is $R;S$.
36. Prove that every endorelation R satisfying $R \subseteq I$ satisfies:
- (a) R is symmetric.
 - (b) R is antisymmetric.
 - (c) R is transitive.
37. Implement Warshall's algorithm in your favorite programming language.

2 Graphs

2.1 Directed Graphs

Both in Computer Science and in the rest of Mathematics *graphs* are studied and used frequently. Graphs come in two flavours, *directed* graphs and *undirected* graphs.

There is no fundamental difference between directed graphs and relations: a directed graph just *is* an endorelation on a given set. If V –for Vertices– is this set and if E –for Edges– is the relation we call the pair (V, E) a directed graph. Usually, set V will be finite, but this is not really necessary: infinite graphs are conceivable too. A directed graph (V, E) is finite if and only if V is finite. Unless stated otherwise, we confine our attention to finite graphs. Always set V will be *nonempty*.

Traditionally, the elements of set V are called “vertices” or “nodes”, whereas the elements of E , that is, the pairs (u, v) satisfying uEv , are called “directed edges” or “arrows”. In this terminology we say that the graph contains “an arrow from u to v ” if and only if uEv . Also, in this case, we say that u is a “predecessor” of v and that v is a “successor” of u .

Graphs can be represented by pictures, in the following way. Every vertex is drawn as a small circle with its name inside the circle, and every arrow from u to v is drawn as an arrow from u ’s circle to v ’s circle. Such a picture may be attractive because it enables us to comprehend, in a single glance, the whole structure of a graph, but, of course, drawing such pictures is only feasible if the set of vertices is not too large. Figures 1 and 2 give simple examples.



Figure 1: The smallest directed graphs: $V = \{a\}$ with $E = \emptyset$ and $E = \{(a, a)\}$

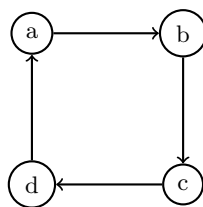


Figure 2: The graph of relation $\{(a, b), (b, c), (c, d), (d, a)\}$

If we are only interested in the pattern of the arrows we may omit the names of

the vertices and simply draw the vertices as dots. The resulting picture is called an “unlabelled” (picture of the) graph.

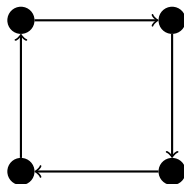


Figure 3: The same graph, unlabelled

Relation E may be such that uEu , for some u . In terms of graphs this means that a vertex may have an arrow from itself to itself. This is perfectly admissible, although in some applications such “auto-arrows” may be undesirable. Notice that the property “having no auto-arrows” is the directed-graph equivalent of the relational property “being irreflexive”.

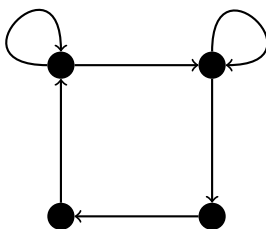


Figure 4: The unlabelled graph of relation $\{ (a, a), (a, b), (b, b), (b, c), (c, d), (d, a) \}$

2.2 Undirected Graphs

Sometimes we are only interested in the (symmetric) concept of nodes being connected, independent of any notion of direction. An “undirected graph” is a symmetric (endo)relation E on a set V . As before, we call the elements of V “nodes” or “vertices”. The pairs (u, v) satisfying uEv are now called “edges”; we also say that such u and v are “directly connected” or “neighbours”.

Relation E being symmetric means that uEv is equivalent to vEu ; hence, being neighbours is a symmetric notion: edge (u, v) is the same as edge (v, u) . In this view an undirected graph just is a special case of a directed graph, with this characteristic property: the graph contains an arrow from u to v if and only if the graph contains an arrow from v to u . So, arrows occur in pairs. See Figure 5, for a simple example. A more concise rendering of an undirected graph is obtained by combining every such pair of arrows into a single, undirected edge, as in Figure 6.

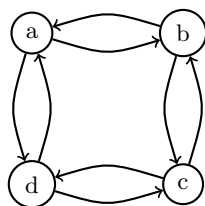


Figure 5: The graph of $\{ (a, b), (b, a), (b, c), (c, b), (c, d), (d, c), (d, a), (a, d) \}$

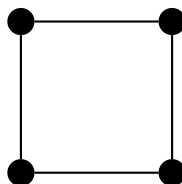


Figure 6: The same graph, with undirected edges and unlabelled

There is no fundamental reason why undirected graphs might not also contain edges connecting a node to itself. Such edges are called “auto loops”. That is, if uEu then u is directly connected to itself, so u is a neighbour to itself.¹ It so happens, however, that in undirected graphs auto-loops are more a nuisance than useful: many properties and theorems obtain a more pleasant form in the absence of auto-loops.

Therefore, we adopt the convention that undirected graphs contain no auto-loops. Formally, this means that an undirected graph is an *irreflexive* and symmetric relation.

In the case of finite graphs we sometimes wish to *count* the number of arrows or edges. We adopt the convention that, in an undirected graph, every pair of directly connected nodes counts as a single edge, even though this single edge corresponds to *two* arrows in the corresponding undirected graph. This reflects the fact that, in a symmetric relation, the pairs (u, v) and (v, u) are indistinguishable. For example, according to this convention, the undirected graph in Figure 6 has four edges.

* * *

We have defined an undirected graph as an irreflexive and symmetric directed graph. Every directed graph can be transformed into an undirected one, just by “ignoring the directions of the arrows”. In terms of relations this amounts to taking the symmetric closure of the relation and removal of the auto-arrows: in the undirected graph nodes u and v are neighbours if and only if, in the directed graph, there is an arrow from

¹This shows that we should not let ourselves be confused by the connotations of the everyday-life word “neighbour”: here the word is used in a strictly technical meaning.

u to v or from v to u (or both), provided $u \neq v$. For example, the directed graph in Figure 4 can thus be transformed into the undirected graph in Figure 6.

2.3 A more compact notation for undirected graphs

We have defined an undirected graph as an irreflexive –no edge between a node and itself– and symmetric relation. Although this is correct mathematically, it is not very practical. For example, the set of edges of the graph in Figure 5 now is $\{ (a, b), (b, a), (b, c), (c, b), (c, d), (d, c), (d, a), (a, d) \}$, in which every edge occurs *twice*: that nodes a and b , for instance, are connected is represented by the presence of both (a, b) and (b, a) in the set of edges. Yet, we do wish to consider this connection as a *single* undirected edge. It is awkward, then, to have to write down both (a, b) and (b, a) to represent this single edge. We would rather not be forced to distinguish these pairs.

We obtain a more convenient representation by using by two-element² sets $\{u, v\}$: as set $\{u, v\}$ equals set $\{v, u\}$ we only need to write this down once. So, in the sequel, an undirected graph will be a pair (V, E) , where V is the set of nodes, as usual, and where E is a set of pairs $\{u, v\}$, with $u, v \in V$ and $u \neq v$, and such that:

$$\{u, v\} \in E \Leftrightarrow \text{“}u \text{ and } v \text{ are connected”} \quad .$$

For example, the set of edges of the graph in Figure 5 can now be written as:

$$\{ \{a, b\}, \{b, c\}, \{c, d\}, \{d, a\} \} .$$

2.4 Additional notions and some properties

Occasionally, we use infix operators for the relations in directed and undirected graphs. That is, sometimes we write uEv as $u \rightarrow v$ and we speak of directed graph (V, \rightarrow) instead of (V, E) . Similarly, for symmetric relations we sometimes use $u \sim v$ instead of uEv and we speak of undirected graph (V, \sim) . So, in this nomenclature, $u \rightarrow v$ means “the graph has an arrow from u to v ” and $u \sim v$ means “in the graph u and v are neighbours”.

In a directed graph (V, \rightarrow) , for every node u the *number of* nodes v satisfying $u \rightarrow v$ is called the “out-degree” of u , whereas the number of nodes u satisfying $u \rightarrow v$ is called the “in-degree” of v , provided these numbers are *finite*. Notice that if V is finite the in-degree and out-degree of every node are finite too. An auto-arrow adds 1, both to the in-degree and the out-degree of its node.

If relation \rightarrow is symmetric, so $u \rightarrow v \Leftrightarrow v \rightarrow u$ for all u, v , then the in-degree of every node equals its out-degree.

In an undirected graph (V, \sim) , the “degree” of a node u is its number of neighbours, that is, the number of nodes v with $u \sim v$. Thus, the degree of a node in an undirected graph equals the in-degree and the out-degree of that node in the underlying directed graph.

²Undirected graphs contain no auto-edges, so the pair (u, u) is not an edge.

With in , out , and deg for “in-degree”, “out-degree”, and “degree” respectively, the following properties hold.

Properties: For all $u, v \in V$:

$$in(v) = \#\{u \mid u \rightarrow v\}$$

$$out(u) = \#\{v \mid u \rightarrow v\}$$

$$deg(u) = \#\{v \mid u \sim v\}$$

By straightforward addition we obtain:

$$(\Sigma v :: in(v)) = \#\{(u, v) \mid u \rightarrow v\} \text{ , and also:}$$

$$(\Sigma u :: out(u)) = \#\{(u, v) \mid u \rightarrow v\} \text{ .}$$

Because the two right-hand side expressions are equal, we conclude:

$$(\Sigma v :: in(v)) = (\Sigma u :: out(u)) \text{ .}$$

In an undirected graph every edge both is an in-arrow and an out-arrow, so in an undirected graph we have, for every node u :

$$deg(u) = in(u) \wedge deg(u) = out(u) \text{ .}$$

In an undirected graph (V, \sim) we have this relation:

$$(\Sigma u :: deg(u)) = 2 * \#\{\{u, v\} \mid u \sim v\} \text{ .}$$

□

* * *

With N for the size of V , so N equals the number of vertices in the graph, we have that the degree of every node is at most $N-1$. If the degree of a node equals $N-1$ then this node is a neighbour of *all other* nodes. If every node in an undirected graph has this property, the graph is called “complete”. Similarly, a directed graph is complete if it contains an arrow from every node to every node. Thus, the complete directed graph corresponds to the complete relation \top , whereas the complete undirected graph corresponds to the relation $\top \setminus I$ (because of the omission of auto-loops). The complete undirected graph with N nodes is called the complete N -graph. Figure 7, for example, gives a picture of the complete 5-graph.

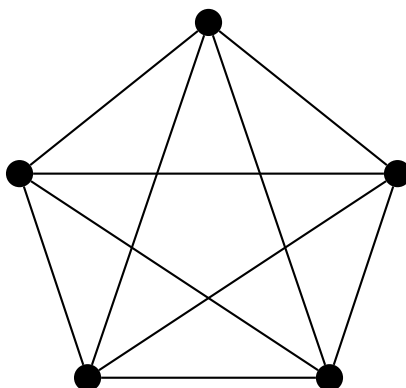


Figure 7: The complete 5-graph, unlabelled

2.5 Connectivity

2.5.1 Paths

We simultaneously consider a directed graph (V, \rightarrow) and an undirected graph (with the same set of nodes) (V, \sim) . A *directed path from node u to node v* is a finite sequence $[s_0, \dots, s_n]$ consisting of $n+1$, $0 \leq n$, nodes satisfying:

$$u = s_0 \wedge (\forall i: 0 \leq i < n: s_i \rightarrow s_{i+1}) \wedge s_n = v .$$

Although this path contains $n+1$ nodes, it pertains to only n arrows, namely the n pairs (s_i, s_{i+1}) , for all $i: 0 \leq i < n$. Therefore, we say that the *length* of this path equals n : the length of a path is the number of arrows in it. If $n=0$ the path contains no arrows and we have $u=v$: the only paths of length 0 are the one-element sequences $[u]$ which are paths from u to u , for every node u . Paths of length 0 are called “empty” whereas paths of positive length are called “non-empty”.

Similarly, in an undirected graph an *undirected path from node u to node v* is a finite sequence $[s_0, \dots, s_n]$ consisting of $n+1$, $0 \leq n$, nodes satisfying:

$$u = s_0 \wedge (\forall i: 0 \leq i < n: s_i \sim s_{i+1}) \wedge s_n = v .$$

Again, the length of this path is n , being the number of edges in it.

Whenever no confusion is possible, we simply use “path” instead of “directed path” or “undirected path”. In any, directed or undirected graph, we call nodes u and v “connected” if the graph contains a path from u to v . Every node is connected to itself, because we have seen that for every node u a path, of length 0, exists from u to u .

In relational terms being connected means being related by the reflexive-transitive closure of the relation. In what follows, we denote the reflexive-transitive closures of relations \rightarrow and \sim by $\overset{*}{\rightarrow}$ and $\overset{*}{\sim}$, respectively, and we denote their transitive closures by $\overset{+}{\rightarrow}$ and $\overset{+}{\sim}$, respectively.

2.1 Lemma. In a directed graph the relation “is connected to” equals $\overset{*}{\rightarrow}$.

Proof. From the chapter on relations we recall the property $R^* = (\bigcup_{n:0 \leq n} R^n)$; in terms of \rightarrow this can be written as: $\overset{*}{\rightarrow} = (\bigcup_{n:0 \leq n} \overset{n}{\rightarrow})$, where $\overset{n}{\rightarrow}$ denotes the equivalent of R^n . This means that $u \overset{*}{\rightarrow} v$ is equivalent to $(\exists n: 0 \leq n: u \overset{n}{\rightarrow} v)$, whereas “ u is connected to v ” is equivalent to

$(\exists n: 0 \leq n: \text{“}u \text{ is connected to } v \text{ by a path of length } n\text{”})$. We now prove the equivalence of these two characterisations term-wise; that is, for all natural n we prove that $u \overset{n}{\rightarrow} v$ is equivalent to “ u is connected to v by a path of length n ”. We do so by Mathematical Induction on n :

$$\begin{aligned}
& u \overset{0}{\rightarrow} v \\
\Leftrightarrow & \quad \{ \text{definition of } \overset{0}{\rightarrow} \} \\
& u I v \\
\Leftrightarrow & \quad \{ \text{definition of } I \} \\
& u = v \\
\Leftrightarrow & \quad \{ \text{definition of path } \} \\
& \text{“the path } [u], \text{ of length } 0, \text{ connects } u \text{ to } v\text{”} \\
\Leftrightarrow & \quad \{ \text{definition of “connected”, see below } \} \\
& \text{“}u \text{ is connected to } v \text{ by a path of length } 0\text{”} .
\end{aligned}$$

As to the logical equivalence in the last step of this derivation: in the direction “ \Rightarrow ” this is just \exists -introduction; in the direction “ \Leftarrow ” we observe: for *every* path $[x]$, of length 0, we have that if $[x]$ connects u to v then $x = u$, hence $[x] = [u]$. (That is, the path of length 0 connecting u to v is *unique*.)

Furthermore, we derive, for $0 \leq n$ and for nodes u, w :

$$\begin{aligned}
& u \overset{n+1}{\rightarrow} w \\
\Leftrightarrow & \quad \{ \text{definition of } \overset{n+1}{\rightarrow} \} \\
& (\exists v :: u \overset{n}{\rightarrow} v \wedge v \rightarrow w) \\
\Leftrightarrow & \quad \{ \text{Induction Hypothesis } \} \\
& (\exists v :: \text{“}u \text{ is connected to } v \text{ by a path of length } n\text{”} \wedge v \rightarrow w) \\
\Leftrightarrow & \quad \{ \text{definition of connected } \} \\
& (\exists v :: (\exists s : \text{“}s \text{ is a path of length } n\text{”} : u = s_0 \wedge s_n = v) \wedge v \rightarrow w) \\
\Leftrightarrow & \quad \{ \wedge \text{ over } \exists \} \\
& (\exists v :: (\exists s : \text{“}s \text{ is a path of length } n\text{”} : u = s_0 \wedge s_n = v \wedge v \rightarrow w)) \\
\Leftrightarrow & \quad \{ \text{dummy unnesting } \} \\
& (\exists s, v : \text{“}s \text{ is a path of length } n\text{”} : u = s_0 \wedge s_n = v \wedge v \rightarrow w) \\
\Leftrightarrow & \quad \{ \text{if } s \text{ is a path of length } n \text{ then } s ++ [v, w] \text{ is a path of length } n+1 :
\end{aligned}$$

$$\begin{aligned}
& \text{dummytransformation } \} \\
& (\exists t: \text{“}t \text{ is a path of length } n+1\text{”} : u = t_0 \wedge t_{n+1} = w) \\
\Leftrightarrow & \quad \{ \text{definition of connected} \} \\
& \text{“}u \text{ is connected to } w \text{ by a path of length } n+1\text{”} .
\end{aligned}$$

□

In a very similar way we can prove that the relation “is connected by a non-empty path length” is equivalent to $\overset{\pm}{\rightarrow}$. Moreover, the proof of the above lemma does not depend on particular properties of the directed relation \rightarrow : the lemma and its proof also are valid for undirected graphs, provided, of course, we replace $\overset{*}{\rightarrow}$ and $\overset{\pm}{\rightarrow}$ by $\overset{\sim}{\rightarrow}$ and $\overset{\simeq}{\rightarrow}$ respectively.

* * *

Note that being connected in an undirected graph is a symmetric relation: u is connected to v if and only if v is connected to u , because $[s_0, \dots, s_n]$ is a path from u to v if and only if the *reverse* of s , that is, the sequence $[s_n, \dots, s_0]$, is a path from v to u .

In directed graphs, being connected is not necessarily symmetric, of course: the existence of a *directed path* (usually) does not imply the existence of directed path in the reverse direction.

2.5.2 Path concatenation

Let s be a directed path of length m from node u to node v , and let t be a directed path of length n from node v to node w . So, the end point of s , which is v , equals the starting point of t , that is, we have $s_m = t_0$.

From s and t we can now construct a directed path, of length $m+n$, from node u to node w ; this is called the “concatenation” of s and t , and we denote it by $s++t$. For s and t paths of length m and n , respectively, their concatenation $s++t$ is a path of length $m+n$, defined as follows:

$$\begin{aligned}
(s++t)_i &= s_i, \text{ for } 0 \leq i \leq m \\
(s++t)_{m+i} &= t_i, \text{ for } 0 \leq i \leq n
\end{aligned}$$

Keep in mind that $s++t$ is defined *only if* $s_m = t_0$, and this is implied by this definition: on the one hand $(s++t)_m = s_m$, on the other hand $(s++t)_m = t_0$. In this case, $s++t$ is a path from u to w indeed. This we prove as follows:

$$\begin{aligned}
& (s++t)_0 \\
= & \quad \{ \text{definition of } ++ \} \\
& s_0 \\
= & \quad \{ s \text{ is a path from } u \text{ to } v \}
\end{aligned}$$

u ,

as required; and, for $0 \leq i < m$:

$$\begin{aligned} & (s \text{++} t)_i \rightarrow (s \text{++} t)_{i+1} \\ = & \quad \{ \text{definition of ++} \} \\ & s_i \rightarrow s_{i+1} \\ = & \quad \{ s \text{ is a path of length } m \} \\ & \text{true} \end{aligned}$$

as required; and, for $0 \leq i < n$:

$$\begin{aligned} & (s \text{++} t)_{m+i} \rightarrow (s \text{++} t)_{m+i+1} \\ = & \quad \{ \text{definition of ++} \} \\ & t_i \rightarrow t_{i+1} \\ = & \quad \{ t \text{ is a path of length } n \} \\ & \text{true} \end{aligned}$$

as required; and, finally:

$$\begin{aligned} & (s \text{++} t)_{m+n} \\ = & \quad \{ \text{definition of ++} \} \\ & t_n \\ = & \quad \{ t \text{ is a path from } v \text{ to } w \} \\ & w \text{ ,} \end{aligned}$$

as required.

Concatenation of undirected paths is defined in exactly the same way: here concatenation is actually an operation on sequences of nodes, and the difference between \rightarrow and \sim , that is, the difference between directed and undirected, only plays a role in the interpretation of such sequences as paths.

We now conclude that, both in directed and in undirected graphs, if a path s , say, exists from node u to node v and if a path t , say, exists from node v to node w , then also a path exists from node u to node w , namely $s \text{++} t$. Thus we have proved the following lemma.

2.2 Lemma. Both in directed and in undirected graphs, the relation “is connected to” is transitive.

□

2.5.3 The triangular inequality

Every path in a graph has a *length*, which is a natural number. Every non-empty set of natural numbers has a smallest element. Therefore, if node u is connected to node v we can speak of the *minimum* of the *lengths* of all paths from u to v . This we call the “distance” from u to v . Because, in undirected graphs, connectedness is

symmetric, we have, in undirected graphs, that the distance from u to v is equal to the distance from v to u .

If u is *not* connected to v we define, for the sake of convenience, the distance from u to v to be ∞ (“infinity”), because ∞ can be considered, more or less, as the identity element of the minimum-operator. Note, however, that ∞ is not a natural number and that we must be very careful when attributing algebraic properties to it. For example, it is viable to define $\infty + n = \infty$, for every natural n , and even $\infty + \infty = \infty$, but $\infty - \infty$ cannot be defined in a meaningful way. An important property is:

$$(8) \quad n < \infty \quad , \quad \text{for all } n \in \text{Nat};$$

$$(9) \quad n \leq \infty \quad , \quad \text{for all } n \in \text{Nat} \cup \{\infty\}.$$

We denote the distance from u to v by $\text{dist}(u, v)$. Then, function dist is defined as follows, for all nodes u, v :

$$\begin{aligned} \text{dist}(u, v) &= \infty \quad , \quad \text{if } u \text{ is not connected to } v; \\ \text{dist}(u, v) &= (\min n : 0 \leq n \wedge \text{“a path of length } n \text{ exists from } u \text{ to } v\text{”} : n) \quad , \\ &\quad \text{if } u \text{ is connected to } v. \end{aligned}$$

Function dist now satisfies what is known in Mathematics as the “triangular inequality”. This lemma holds for both directed and undirected graphs.

2.3 Lemma. All nodes u, v, w satisfy: $\text{dist}(u, w) \leq \text{dist}(u, v) + \text{dist}(v, w)$.

Proof. By (unavoidable) case analysis. If $\text{dist}(u, v) = \infty$ or $\text{dist}(v, w) = \infty$ then also $\text{dist}(u, v) + \text{dist}(v, w) = \infty$; now, by property (9), we have $\text{dist}(u, w) \leq \infty$, so we conclude, for this case: $\text{dist}(u, w) \leq \text{dist}(u, v) + \text{dist}(v, w)$, as required.

Remains the case $\text{dist}(u, v) < \infty$ and $\text{dist}(v, w) < \infty$. In this case, paths exist from u to v and from v to w . Let s be a path, of length m , from u to v and let t be a path, of length n , from v to w . Then, as we have seen in the previous subsection, $s + t$ is a path, of length $m + n$, from u to w . By the definition of dist , we conclude: $\text{dist}(u, w) \leq m + n$. As this inequality is true for *all* such paths s and t , it is true for paths of minimal length as well. Hence, also for this case we have: $\text{dist}(u, w) \leq \text{dist}(u, v) + \text{dist}(v, w)$, as required.

□

2.5.4 A lemma and its proof

In Section 1.4, on Warshall’s algorithm, we have defined a sequence of auxiliary relations W_n , for $0 \leq n \leq N$. Relation W_n is the transitive closure of R – the relation under consideration –, not based on ordinary composition, however, but on the restricted composition “; _{n} ” instead. That is, W_n is the smallest of all relations X satisfying:

$$R \cup X ;_n X \subseteq X .$$

Just as $\xrightarrow{*}$ and $\xrightarrow{+}$, the relations W_n can be interpreted in terms of paths; for this we need the auxiliary notion of a, so-called, “ n -path”, for all $0 \leq n \leq N$. A *directed n -path* is a directed path whose *internal* nodes are at most n . That is, if $[s_0, \dots, s_m]$ is a directed path, of length m , then s is a directed n -path if:

$$(\forall i: 1 \leq i < m: s_i \leq n) .$$

In very much the same way as we did in Lemma 2.1, we can now prove that the relations W_n have the following property.

2.4 Properties. For all $n, 0 \leq n \leq N$, relation W_n is equal to the relation “is connected by a non-empty n -path”.

□

Using the notion of n -paths we are now ready to prove the following lemma.

2.5 Lemma. For all $n, 1 \leq n \leq N$: $W_n = W_{n-1} \cup W_{n-1} \langle n \rangle W_{n-1}$.

Proof. By mutual set inclusion.

“ \supseteq ”: We must prove $W_n \supseteq W_{n-1}$ and $W_n \supseteq W_{n-1} \langle n \rangle W_{n-1}$; we do so separately. Obviously, $W_n \supseteq W_{n-1}$ because $X ;_{n-1} X \Rightarrow X ;_n X$, for all X . If nodes u, v satisfy $u (W_{n-1} \langle n \rangle W_{n-1}) v$ then, by definition of $\langle n \rangle$, we have: $u W_{n-1} n$ and $n W_{n-1} v$. So, an $(n-1)$ -path exists from u to n and an $(n-1)$ -path exists from n to v . The concatenation of these two paths then is an n -path from u to v . Hence, $u W_n v$.

“ \subseteq ”: Assume $u W_n v$, for nodes u, v . We must show that either $u W_{n-1} v$ or $u (W_{n-1} \langle n \rangle W_{n-1}) v$. Because $u W_n v$, an n -path exists from u to v . If no internal node of this paths equals n then this path is an $(n-1)$ -path as well, so then $u W_{n-1} v$. Remains, on the other hand, the case that n occurs at least once as an internal node in this path. Let $[s_0, \dots, s_m]$ be this path, of length m . Let p be the *smallest* and let q be the *largest* of all indices $i, 1 \leq i < m$, with $s_i = n$. Then we have: $s_p = n$ and $s_q = n$, but also $(\forall i: 1 \leq i < p: s_i < n)$ and $(\forall i: q+1 \leq i < m: s_i < n)$. Therefore, $[s_0, \dots, s_p]$ and $[s_q, \dots, s_m]$ are $(n-1)$ -paths from u to n and from n to v respectively, so we conclude $u W_{n-1} n$ and $n W_{n-1} v$, and, hence, $u (W_{n-1} \langle n \rangle W_{n-1}) v$ as well.

□

2.5.5 Connected components

A directed graph (V, \rightarrow) is *strongly connected* if every node is connected to every node, that is, if there is a directed path from every node u to every node v . In relational terms, this means that $\xrightarrow{*} = \top$. The adverb “strongly” stresses the fact that, in directed graphs, strong connectedness is a symmetric notion: for every two nodes u, v there is a path from u to v and there is path from v to u .

An undirected graph is *connected* if every pair of nodes is connected by a path. Relationally, a graph is connected if and only if $\sim^* = \top$. As we have seen, in undirected graphs connectedness is symmetric. It even is an equivalence relation. A *connected component* is a *maximal subset* of the nodes of the graph that is connected: the connected components of an undirected graph are the equivalence classes of \sim^* .

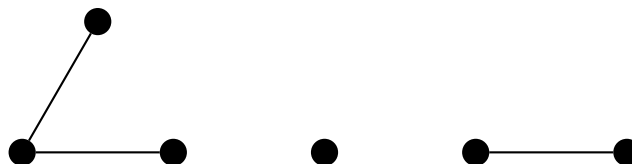


Figure 8: An undirected graph with 3 connected components

2.6 Cycles

A *cycle* in a graph is a non-empty path from a node to itself. Particularly in undirected graphs a proper definition of what constitutes a cycle is not as simple as it may seem, though. Generally, a graph may contain few cycles, many cycles, or no cycles at all. In the latter case the graph is called *acyclic*.

2.6.1 Directed cycles

In a directed graph a *cycle* is a (directed) path from a node to itself. For example, if $a \rightarrow b$ and $b \rightarrow a$ then the path $[a, b, a]$ is a cycle, and so is the path $[b, a, b]$. Although these are different paths they constitute, in a way, the same cycle. The simplest possible case of a directed cycle is $[a, a]$, namely if $a \rightarrow a$.



Figure 9: A simple directed cycle



Figure 10: An even simpler cycle

2.6.2 Undirected cycles

In undirected graphs the notion of cycles is somewhat more complicated. For example, if, in undirected graph (V, \sim) , we have $a \sim b$ and, hence, also $b \sim a$, then $[a, b, a]$ is a path from node a to itself. Yet, we do not wish to consider this a cycle. More generally, we do not wish the pattern $[\dots, a, b, a, \dots]$ to occur anywhere in a cycle: in a cycle, every next edge should be different from its predecessor. As a consequence, in an undirected graph the smallest possible cycle involves at least *three* nodes and *three* edges.

These considerations give rise to the following definition. An undirected cycle is a path $[s_0, \dots, s_n]$, of length n , with the following additional properties:

$$\begin{aligned} 3 &\leq n \\ s_0 &= s_n \\ (\forall i: 0 \leq i \leq n-2: s_i &\neq s_{i+2}) \wedge s_{n-1} \neq s_1 \end{aligned}$$

The first of these conditions expresses that a cycle comprises at least 3 nodes, the second condition expresses that the path's last node equals its first node – thus “closing the cycle” –, and the last condition precludes that every two successive edges in the cycle are different. The conjunct $s_{n-1} \neq s_1$ really is needed here: the “last” edge, $\{s_{n-1}, s_n\}$, which is the same as $\{s_{n-1}, s_0\}$, and the “first” edge, $\{s_0, s_1\}$, are successive too, which must be different as well.

In Figure 14, for example, we have that $[a, b, d, b, c, a]$ is not a cycle, because it contains edge $\{b, d\}$ twice in succession. Without the conjunct $s_{n-1} \neq s_1$, however, the path $[d, b, c, a, b, d]$ would be a cycle, which is undesirable: whether or not a certain collection of nodes constitutes a cycle should not depend on which node is the first node of the path representing that cycle.

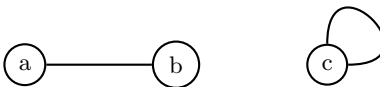


Figure 11: No cycles at all



Figure 12: Not even an (undirected) graph

Thus we obtain the following lemma, which expresses that cycles are “invariant under rotation”. This lemma is useful because it allows us to let any node in a cycle be the starting node of the path representing that cycle.

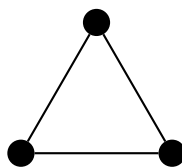
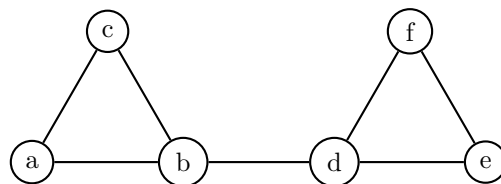


Figure 13: The smallest undirected cycle

Figure 14: $[a, b, d, e, f, d, b, c, a]$ is a cycle, $[a, b, d, b, c, a]$ is not

2.6 Lemma. [Rotation Lemma] For every natural n , $3 \leq n$, a path $[s_0, s_1, \dots, s_{n-1}, s_0]$ is a cycle if and only if the path $[s_1, \dots, s_{n-1}, s_0, s_1]$ is a cycle.

□

2.7 Euler and Hamiltonian cycles

2.7.1 Euler cycles

In an undirected graph a cycle with the property that it contains *every edge* of the graph exactly once is called an *Euler cycle*.

2.7 Theorem. For every connected graph (V, \sim) :

$$\text{“}(V, \sim) \text{ contains an Euler cycle”} \Leftrightarrow (\forall v : v \in V : \text{“}deg(v) \text{ is even”}) .$$

Proof. By mutual implication.

“ \Rightarrow ”: We consider an Euler cycle in graph (V, \sim) . Let v be a node. Wherever v occurs in the Euler cycle v has a predecessor u , say, in the cycle and a successor w , say, in the cycle. This means that u, v, w are all different and $u \sim v$ and $v \sim w$. Thus, all edges associated with v occurring in the Euler cycle occur in pairs; hence, the total number of edges associated with v occurring in the Euler cycle is even. Because the cycle is an Euler cycle *all* of v 's edges occur in the Euler cycle; hence, $deg(v)$ is even.

“ \Leftarrow ”: Assuming $(\forall v : v \in V : \text{“}deg(v) \text{ is even”})$ we prove the existence of an Euler cycle by sketching an algorithm for the construction of an Euler cycle. This algorithm consists of two phases. In the first phase a collection of (one or more) cycles is formed such that *every* edge of the graph occurs exactly once in exactly one of these cycles.

In the second phase, the cycles in this collection are combined into larger cycles, thus reducing the number of cycles in the collection while retaining the property that *every* edge of the graph occurs exactly once in exactly one of the cycles in the collection. As soon as this collection contains only one cycle, this one cycle is a Euler cycle.

first phase: Initially all edges are white. The property $(\forall v : v \in V : \text{“deg}(v) \text{ is even”})$ will remain valid for the subgraph formed by V and the white edges only: it is an invariant of this phase. Another invariant is that all red edges form a collection of cycles with the property that *every red edge* of the graph occurs exactly once in exactly one of these cycles. Initially this is true because there are no red edges: initially the collection of red cycles is empty. If, on the other hand, all edges are red the collection of red cycles comprises all edges of the graph, and the first phase terminates. As long as the graph contains at least one white edge, the following step is executed.

Select a white edge, $\{s_0, s_1\}$, say. Because $\text{deg}(s_1)$ is even, node s_1 has a neighbour s_2 , say, that differs from s_0 and such that edge $\{s_1, s_2\}$ is white as well. Repeating this indefinitely yields an infinite sequence $s_i : 0 \leq i$ of nodes, pairwise connected by white edges. As the graph is finite, this sequence contains a sub-path $[s_p, \dots, s_q]$, for some p, q with $0 \leq p < q$, that is a cycle, comprising white edges only. Now all white edges in this cycle are turned red. Because, for every node in this cycle, its associated edges occur in pairs, the number of white edges associated with any node in this cycle is even and, as a result, the degree of all nodes remains even under reddening of the white edges in this cycle. Because in a undirected graph every cycle contains at least 3 edges the number of white edges thus decreases (by at least 3), guaranteeing termination of the first phase.

second phase: The second phase terminates if the collection of red cycles contains only one cycle. As long as this collection contains at least two cycles it also contains two cycles that have a node in common: for any node u on one cycle and any node x *not* on this cycle, we have that the graph has a path connecting u and x . Let w be the node on this path that is *closest* to u that is not on the cycle. Then w 's predecessor v , say, on this path is on the cycle; so, now we have the following situation: node v is on one cycle, node w is not on this cycle and $\{v, w\}$ is an edge of the graph, which, therefore, also is not part of the cycle we started with. So, this edge is part of one of the other cycles in the collection. Thus, we have identified two cycles having node v in common. Let s be a path, connecting v to v , representing the one cycle and let t be a path, also connecting v to v , representing the other cycle. Then their concatenation $s \text{++} t$ is a cycle connecting v to v too, and $s \text{++} t$ contains all edges from s and t together exactly once. Thus, replacing, in the collection of red cycles, s and t by $s \text{++} t$ respects the invariant that every edge of the graph occurs exactly once in exactly one cycle; moreover, this replacement decreases the number of cycles in the collection.

□

2.7.2 Hamiltonian cycles

In a (directed or undirected) graph a cycle with the property that it contains *every node* of the graph exactly once is called a *Hamiltonian cycle*.

A naive algorithm to compute whether a given graph contains a Hamiltonian cycle is conceptually simple: enumerate all cycles and check whether any of them is a Hamiltonian cycle. This naive algorithm is quite inefficient, of course, but really efficient algorithms are not (yet) known: the problem to decide whether a graph contains a Hamiltonian cycle is *NP-hard*, which in practice means that all algorithms will require an amount of computation time that grows exponentially with the size of the graph.

Notice the contrast in complexity between the notion of Euler and Hamiltonian cycles. On the one hand, Theorem 2.7 provides a simple algorithm to evaluate the existence of an Euler cycle – just calculate the degrees of the nodes –, and its proof contains a relatively straightforward algorithm for the construction of an Euler cycle. On the other hand, calculating the existence of an Hamiltonian cycle, let alone construction of one, is NP-hard.

Thus, two seemingly similar notions – Euler cycles and Hamilton cycles – happen to have essentially different properties.

2.7.3 A theorem on Hamiltonian cycles

We consider finite, undirected graphs, with at least 4 nodes. We present a theorem giving a sufficient condition for the existence of Hamiltonian cycles, namely if the graph contains “sufficiently many” edges. In our case the notion of “sufficiently many” and the theorem take the following shape.

2.8 Theorem. We consider an undirected graph with n nodes, $4 \leq n$. If, for every two unconnected nodes, the sum of their degrees is at least n , then the graph contains a Hamiltonian cycle.

□

To formalize this, let V be a (fixed) set of nodes, with $n = \#V$, $4 \leq n$. In what follows variables u, v, p, q range over V , with $p \neq q$. The set E of edges is variable; that is, as a function of E we define predicates P and H , as follows:

$$P(E) = (\forall u, v: u \neq v: \{u, v\} \notin E \Rightarrow \text{deg}(u) + \text{deg}(v) \geq n) \text{ , and:}$$

$$H(E) = \text{“graph } (V, E) \text{ contains a Hamiltonian cycle” .}$$

Predicate P formalizes our particular version of “sufficiently many”: P expresses that, for every two unconnected nodes, the sum of their degrees is at least n .

Both P and H are *monotonic*, as follows:

monotonicity: For all E and for any two nodes p, q :

$P(E) \Rightarrow P(E \cup \{\{p, q\}\})$, and:

$H(E) \Rightarrow H(E \cup \{\{p, q\}\})$.

□

In addition, for the extreme cases, the empty graph \perp and the complete graph \top , we have:

$$\neg(P(\perp)) \wedge \neg(H(\perp)) \wedge P(\top) \wedge H(\top) .$$

The theorem now states that every graph satisfying predicate P contains at least one Hamiltonian cycle.

2.9 Theorem. $(\forall E :: P(E) \Rightarrow H(E))$.

□

We present two proofs for this theorem. These proofs are essentially the same, but they differ in their formulation. The crucial part in both proofs is the following:

Core Property: For set E of edges and for any two nodes p, q with $\neg(\{p, q\} \in E)$:

$$P(E) \wedge H(E \cup \{\{p, q\}\}) \Rightarrow H(E) .$$

□

Notice that the Core Property also holds if $\{p, q\} \in E$, but in a trivial way only: then $H(E \cup \{\{p, q\}\}) = H(E)$, so in this case the property is void.

We will present a proof for the Core Property later, but first we will show how it is used in the proofs of the Theorem.

2.7.4 A proof by contradiction

The first proof runs as follows, by contradiction. That is, we suppose that the Theorem is false. Then, there exists a set F of edges such that $P(F)$ and $\neg(H(F))$. Because $\neg(H(F))$ and $H(\top)$, and because $F \subseteq \top$, there also exists a “turning point”, that is, a set E of edges and a pair p, q of nodes such that:

$$F \subseteq E \wedge \neg(H(E)) \wedge H(E \cup \{\{p, q\}\}) .$$

Notice that, because of $\neg(H(E)) \wedge H(E \cup \{\{p, q\}\})$, we have –Leibniz!– that $E \neq E \cup \{\{p, q\}\}$, hence $\{p, q\} \notin E$.

Because of the monotonicity of P , and because $P(F)$ and $F \subseteq E$, set E satisfies $P(E)$ too. Now, from $P(E)$ and $H(E \cup \{\{p, q\}\})$ we conclude, using the Core Property, $H(E)$. In conjunction with the assumed $\neg(H(E))$ we obtain the desired contradiction.

2.7.5 A more explicit proof

The reasoning in the previous proof is somewhat strange: the assumption $\neg(H(E))$ is not really used in the proof proper: it is only used to conclude a contradiction. Therefore, we should be able to construct a more direct proof. In addition what does “there exists a ‘turning point’” really mean, mathematically speaking?

Because of the monotonicity of P we have $P(E) \Rightarrow P(E \cup \{\{p, q\}\})$; therefore, by means of elementary propositional calculus, the Core Property can be rewritten thus:

$$(10) \quad (P(E \cup \{\{p, q\}\}) \Rightarrow H(E \cup \{\{p, q\}\})) \Rightarrow (P(E) \Rightarrow H(E)) \quad ,$$

and this smells very strongly of a proof by Mathematical Induction. As a matter of fact, this *is* Mathematical Induction, albeit in a somewhat unusual direction, namely from larger towards smaller.

Firstly, we have $H(\top)$ –the complete graph contains a Hamiltonian cycle, very many even–, so we also have $P(\top) \Rightarrow H(\top)$. This is the basis of the induction.

Secondly, property (10) now represents the induction step. Because every set E of edges can be obtained from a larger set $E \cup \{\{p, q\}\}$, with $\{p, q\} \notin E$, we are done.

Notice that the fact that the collection of all possible sets of edges is *finite*³ is of no consequence: although usually applied to infinite sets the principle of Mathematical Induction is perfectly valid in a finite setting.

2.7.6 Proof of the Core Property

We repeat the Core Property, which is the essential part of both proofs of the Theorem.

Core Property: For set E of edges and for any two nodes p, q with $\{p, q\} \notin E$:

$$P(E) \wedge H(E \cup \{\{p, q\}\}) \Rightarrow H(E) \quad .$$

□

To prove this we assume that E is a set of edges and p, q are different nodes, such that $\{p, q\} \notin E$, satisfying $P(E)$ and $H(E \cup \{\{p, q\}\})$. The latter means that the graph $(V, E \cup \{\{p, q\}\})$ contains a Hamiltonian cycle. If such a Hamiltonian cycle does *not* contain edge $\{p, q\}$, then it also is a Hamiltonian cycle in the graph (V, E) ; hence, $P(E)$ and in this case we are done.

So, remains the case that $(V, E \cup \{\{p, q\}\})$ contains a Hamiltonian cycle that does contain edge $\{p, q\}$. Now we have to prove $P(E)$, that is, we must prove that (V, E) contains a Hamiltonian cycle as well, that is, *without* edge $\{p, q\}$.

For this purpose, let $[s_0, s_1, \dots, s_n]$ be a Hamiltonian cycle in $(V, E \cup \{\{p, q\}\})$. This means that $\{s_i \mid 0 \leq i < n\} = V$ –recall that $n = \#V$ –, that $s_n = s_0$, and that $(\forall i: 0 \leq i < n: \{s_i, s_{i+1}\} \in E \cup \{\{p, q\}\})$. We assume that this cycle contains edge $\{p, q\}$ and, without loss of generality, we assume that $s_0 = p$ and $s_1 = q$.

³for our *fixed*, finite set of nodes

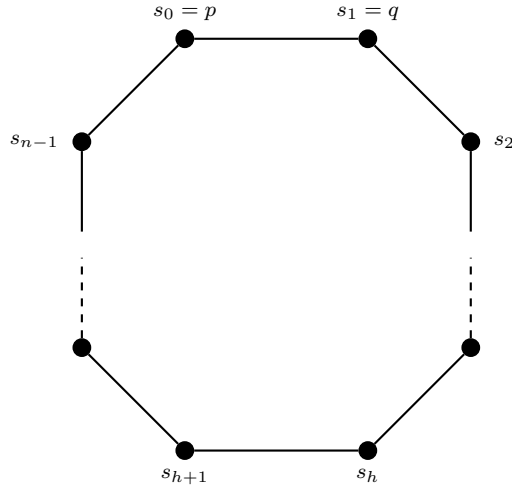


Figure 15: a Hamiltonian cycle, with edge $\{p, q\}$

In this setting we prove that (V, E) contains a Hamiltonian cycle. To construct a Hamiltonian cycle *not* containing edge $\{p, q\}$ we take the Hamiltonian cycle introduced above, containing edge $\{p, q\}$, as a starting point. Removal of edge $\{p, q\}$ destroys the cycle, and what remains is a path connecting s_1 , that is q , to s_0 , that is p , that still contains all nodes of the graph and all edges of which are in E .

Now we must restore the cycle by somehow reconnecting p and q , using edges in E only. We do so by selecting an index h in the interval $[2..n-1]$ such that both $\{p, s_h\} \in E$ and $\{q, s_{h+1}\} \in E$. To show that this is possible we need the theorem's assumption $P(E)$, which was defined as:

$$(\forall u, v : u \neq v : \{u, v\} \notin E \Rightarrow \deg(u) + \deg(v) \geq n) .$$

Applying this to p, q and using $\{p, q\} \notin E$ we obtain:

$$(11) \quad \deg(p) + \deg(q) \geq n .$$

Let $x = \#\{i \in [2..n-1] \mid \{p, s_i\} \in E\}$ and let $y = \#\{j \in [2..n-1] \mid \{q, s_{j+1}\} \in E\}$; now we calculate:

$$\begin{aligned} & \deg(p) + \deg(q) \geq n \\ \Leftrightarrow & \quad \{ \text{definition of } \deg \text{ (twice), using } s_0 = p \text{ and } s_1 = q \} \\ & \#\{i \in [1..n] \mid \{s_0, s_i\} \in E\} + \#\{j \in [2..n] \mid \{s_1, s_j\} \in E\} \geq n \\ \Leftrightarrow & \quad \{ \text{split off } i=1 \text{ and } i=n-1, \text{ using } \{s_0, s_1\} \notin E \text{ and } \{s_{n-1}, s_n\} \in E \} \\ & 1 + \#\{i \in [2..n-1] \mid \{s_0, s_i\} \in E\} + \#\{j \in [2..n] \mid \{s_1, s_j\} \in E\} \geq n \\ \Leftrightarrow & \quad \{ \text{split off } j=2 \text{ and } j=n, \text{ using } \{s_1, s_2\} \in E \text{ and } \{s_1, s_n\} \notin E \} \end{aligned}$$

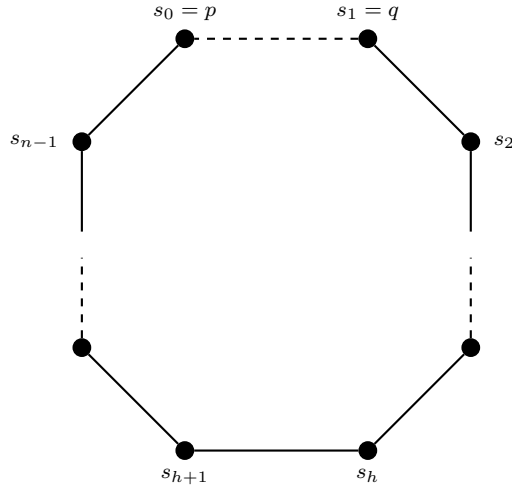


Figure 16: The remains of the cycle, after removal of edge $\{p, q\}$

$$\begin{aligned}
& 1 + \#\{i \in [2..n-1) \mid \{s_0, s_i\} \in E\} + 1 + \#\{j \in [3..n) \mid \{s_1, s_j\} \in E\} \geq n \\
\Leftrightarrow & \quad \{ \text{dummy transformation } j := j+1 \} \\
& 1 + \#\{i \in [2..n-1) \mid \{s_0, s_i\} \in E\} + 1 + \#\{j \in [2..n-1) \mid \{s_1, s_{j+1}\} \in E\} \geq n \\
\Leftrightarrow & \quad \{ \text{definitions of } x \text{ and } y, \text{ using } s_0 = p \text{ and } s_1 = q \} \\
& 1 + x + 1 + y \geq n \\
\Leftrightarrow & \quad \{ \text{calculus} \} \\
& x + y \geq n - 2 .
\end{aligned}$$

So, the number of indices i in the interval $[2..n-1)$ for which $\{p, s_i\} \in E$ plus the number of indices i in the range $[2..n-1)$ for which $\{q, s_{i+1}\} \in E$ is at least $n-2$. The size of the interval $[2..n-1)$, however, only is $n-3$; hence the two sets of indices have a non-empty intersection: there exists an index h , $h \in [2..n-1)$, such that both $\{p, s_h\} \in E$ and $\{q, s_{h+1}\} \in E$.

For every such an index h , $[s_0, s_h, \dots, s_2, s_1, s_{h+1}, \dots, s_{n-1}, s_n]$ is a Hamiltonian cycle in the graph (V, E) . Because we have shown the existence of such an h , we conclude the existence of a Hamiltonian cycle in (V, E) , which was our goal.

remark: The existence of an index h in the interval $[2..n-1)$ implies that this interval is nonempty, that is, $2 < n-1$, which boils down to $4 \leq n$. Hence, the proofs of the Theorem presented here are only valid for graphs with at least 4 nodes. It can be easily verified that the Theorem also holds for $n=3$: the complete 3-graph – a “triangle” – is the only one satisfying predicate P , and a “triangle” is a Hamiltonian cycle. For $n < 3$ the Theorem does not

hold.

□

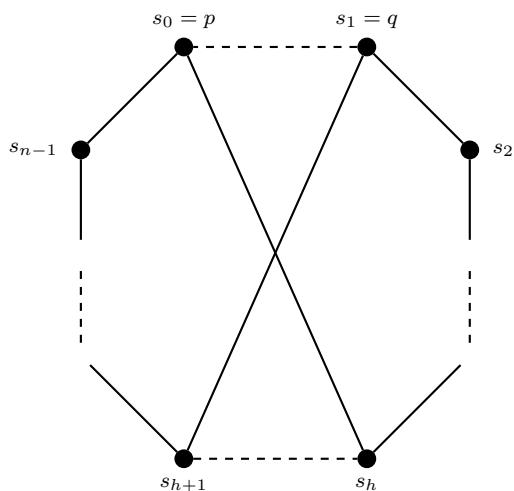


Figure 17: a Hamiltonian cycle, without edge $\{p, q\}$

2.8 Ramsey's theorem

2.8.1 Introduction

We are having a party at which every two guests either do know each other or do not know each other. If the number of guests at the party is “large enough” then the party has at least 5 guests all of which either do know one another or do not know one another. How large must the party be for this to be true?

F.P. Ramsey has developed some theory for the treatment of problems like this. This theory makes it possible to draw rather global conclusions about undirected graphs, independently of their actual structure.

To illustrate this we present a simple theorem that represents his work. In the above example the guests at the party can be considered the nodes of an undirected graph. Any two nodes are connected by an edge if and only if the two corresponding guests do know each other. A set of 5 guests all of which do know one another then amounts to a subgraph of size 5 that is *complete*, that is, we say that the whole graph *contains a complete 5-graph*. How do we formulate, on the other hand, that from a set of 5 guests every two guests do *not* know each other? Well, this means that the whole graph contains 5 nodes every two of which are *not* connected.

We might as well, however, consider the *complement* graph, in which two nodes are connected by an edge if and only if the two corresponding guests do *not* know

each other. As a matter of fact, the problem as stated is *symmetric* in the notions of “knowing each other” and “not knowing each other”. It is, therefore, awkward to destroy this symmetry by representing the one concept by the presence of edges and the other concept by their absence. Moreover, we have two possibilities here, the choice between which is irrelevant.

To restore the symmetry we, therefore, consider a complete undirected graph, of which the set of edges has been partitioned into two subsets – or more than two, in the more general case of Ramsey’s theory –. The one subset then represents the pairs of guests who know each other and the other subset represents the pairs of guests who do not know each other.

Partitioning a set into (disjoint) subsets can be represented conveniently by *colouring*. In our case, partitioning the edges of a complete graph into two subsets can be represented by colouring each edge with one out of two colours. (And, of course, with more than two colours we can represent partitionings into more than two subsets.) Now, an edge of the one colour may represent a pair of guests who do know each other, whereas an edge of the other colour may represent a pair of guests who do not know each other. Thus, the symmetry between “knowing” and “not knowing” is restored.

2.8.2 Ramsey’s theorem

We consider finite, complete, undirected graphs only. For the sake of brevity, we will use “ k -graph” for the “complete k -graph”, for any natural k , $2 \leq k$. Formally, a colouring of a graph’s edges is a function from the set of edges to the set of colours used, $\{\text{red}, \text{blue}\}$, say, if two colours are sufficient. So, a colouring is a function of type $E \rightarrow \{\text{red}, \text{blue}\}$, and if c is such a colouring, then for any edge $\{u, v\}$ we have either $c(\{u, v\}) = \text{red}$ or $c(\{u, v\}) = \text{blue}$, but not both simultaneously, as we presume that $\text{red} \neq \text{blue}$: every edge has only one colour. In what follows we use variables c and d to denote colourings.

Again for brevity’s sake, we say that the k -graph “contains a red m -graph” if the nodes of the k -graph contain a subset of m nodes such that all edges connecting these nodes are red; that is, these m nodes together with their edges constitute a completely red m -graph as a subgraph of the k -graph, for any k, m with $2 \leq m \leq k$.

As an example, notice that the 2-graph has two nodes only, connected by one single edge; hence, the proposition “the k -graph contains a red 2-graph” is equivalent to the proposition “the k -graph contains at least one red edge”.

The proposition “the k -graph contains a red m -graph” depends on the parameters k and m , of course, but also on the actual colouring. So, it is a predicate with three parameters. Calling this predicate Rd , we define it as follows, together with a similar predicate Bl , for the colour blue, for all k, c, m with $2 \leq m \leq k$:

$Rd(k, c, m) \Leftrightarrow$ “the k -graph with colouring c contains a red m -graph”, and:

$Bl(k, c, m) \Leftrightarrow$ “the k -graph with colouring c contains a blue m -graph” .

These predicates are monotonic, in the following way. Suppose $Rd(k, c, m)$, for some k, c, m . Then we also have $Rd(k+1, c, m)$, provided we consider the colouring of

the $(k+1)$ -graph as an extension of the colouring of the k -graph – both colourings being denoted here by the very same c –, just as the $(k+1)$ -graph can be viewed as an extension of the k -graph. For this purpose we consider the $k+1$ nodes of the $(k+1)$ -graph as a set of k nodes, forming a k -graph, plus one additional node, which may remain anonymous. Every colouring of the $(k+1)$ -graph thus induces a colouring of the k -graph; as stated, function c denotes either colouring.

Ramsey's theorem now is about a function R , say, defined as follows, for all m, n with $2 \leq m$ and $2 \leq n$:

$$R(m, n) = \text{“the smallest of all natural numbers } k \text{ satisfying:} \\ (\forall c :: Rd(k, c, m) \vee Bl(k, c, n) \text{)”} .$$

The function value $R(m, n)$ is only well-defined, of course, if at least one natural number k exists satisfying $(\forall c :: Rd(k, c, m) \vee Bl(k, c, n))$: only then we can speak of the smallest such number. Notice that, by definition, if $R(m, n) = p$ then for every k , $p \leq k$, the k -graph contains at least one red m -graph or contains at least one blue n -graph (or both).

The following theorem states that such natural numbers exist and provides an upperbound for $R(m, n)$.

2.10 Theorem. (Ramsey)

$$R(m, n) \leq \binom{m+n-2}{m-1} , \text{ for all } m, n : 2 \leq m \wedge 2 \leq n .$$

□

Notice that, by definition, $R(m, n)$ is symmetric in m and n , that is, we have: $R(m, n) = R(n, m)$, because for every colouring c satisfying $Rd(k, c, m) \vee Bl(k, c, n)$ a colouring d exists – which one? – satisfying $Rd(k, d, n) \vee Bl(k, d, m)$. The expression $\binom{m+n-2}{m-1}$ does not look symmetric, at least, not at first sight. Yet, it is, because binomial coefficients satisfy the following, general property:

$$\binom{m+n}{m} = \binom{m+n}{n} , \text{ for all } m, n : 1 \leq m \wedge 1 \leq n ,$$

as a result of which we also have: $\binom{m+n-2}{m-1} = \binom{m+n-2}{n-1}$.

Proof of the Theorem: By Mathematical Induction on the value of $m+n$; that is, the Induction Hypothesis is:

$$R(p, q) \leq \binom{p+q-2}{p-1} , \text{ for all } p, q : 2 \leq p \wedge 2 \leq q \wedge p+q < m+n .$$

We distinguish 3 cases.

Firstly, $2 \leq m \wedge n = 2$: We consider the m -graph. Let c be the colouring in which all edges of the m -graph are red, so this particular c yields $Rd(m, c, m)$. For every *other* colouring c we have that not all edges of the m -graph are red, so the m -graph contains at least one blue edge, which means $Bl(m, c, 2)$, for all other c . Combining these cases we obtain $(\forall c :: Rd(m, c, m) \vee Bl(m, c, 2))$, from which we conclude that $R(m, n) \leq m$. (Actually, we have $R(m, n) = m$ because no smaller graph contains a red m -graph, but the upper bound is all we need.) Now $m = \binom{m+2-2}{m-1}$, so we conclude $R(m, n) \leq \binom{m+2-2}{m-1}$, as required.

Secondly, $m = 2 \wedge 2 \leq n$: By symmetry with the previous case.

Thirdly, $3 \leq m \wedge 3 \leq n$: We need an additional property, to be proved later; the need of this property is inspired by a well-known property of binomial coefficients:

$$\begin{aligned}
& R(m, n) \\
\leq & \quad \{ \bullet \text{ property of } R, \text{ see below, using } 3 \leq m \wedge 3 \leq n \} \\
& R(m-1, n) + R(m, n-1) \\
\leq & \quad \{ \text{Induction Hypothesis (twice)} \} \\
& \binom{m+n-3}{m-2} + \binom{m+n-3}{m-1} \\
= & \quad \{ \text{property of binomial coefficients} \} \\
& \binom{m+n-2}{m-1} .
\end{aligned}$$

□

In the above proof of the Theorem we have used the following property of R , which constitutes the core of the proof.

property: $R(m, n) \leq R(m-1, n) + R(m, n-1)$, for all $m, n : 3 \leq m \wedge 3 \leq n$.

proof: Let $k = R(m-1, n) + R(m, n-1)$. To prove that $R(m, n) \leq k$ it suffices to prove $Rd(k, c, m) \vee Bl(k, c, n)$, for all colourings c . Therefore, let c be a colouring of the k -graph. Let v be a node of the k -graph and in what follows dummy u also ranges over the nodes of the k -graph. We define subsets X and Y of the nodes, as follows:

$$X = \{ u \in V \mid u \neq v \wedge c(\{u, v\}) = \text{red} \} , \text{ and:}$$

$$Y = \{u \in V \mid u \neq v \wedge c(\{u, v\}) = \text{blue}\} .$$

Then X and Y and $\{v\}$ partition the nodes of the k -graph, so we have:

$$\#X + \#Y + 1 = R(m-1, n) + R(m, n-1) .$$

From this it can be derived that $R(m-1, n) \leq \#X \vee R(m, n-1) \leq \#Y$, by contraposition:

$$\begin{aligned} & \#X < R(m-1, n) \wedge \#Y < R(m, n-1) \\ \Leftrightarrow & \quad \{ \text{all values here are integers} \} \\ & \#X \leq R(m-1, n) - 1 \wedge \#Y \leq R(m, n-1) - 1 \\ \Rightarrow & \quad \{ \text{monotonicity of addition} \} \\ & \#X + \#Y \leq R(m-1, n) + R(m, n-1) - 2 \\ \Leftrightarrow & \quad \{ \text{all values here are integers} \} \\ & \#X + \#Y + 1 < R(m-1, n) + R(m, n-1) \\ \Rightarrow & \quad \{ < \text{ is irreflexive} \} \\ & \#X + \#Y + 1 \neq R(m-1, n) + R(m, n-1) , \end{aligned}$$

which settles the issue.

We now prove the required property, $Rd(k, c, m) \vee Bl(k, c, n)$, by distinguishing the two cases of this disjunction.

Case $R(m-1, n) \leq \#X$: From the definition of R we conclude, for our colouring c , that either $Rd(\#X, c, m-1)$ or $Bl(\#X, c, n)$. If $Rd(\#X, c, m-1)$ then X contains a red $(m-1)$ -graph. By definition of X , we also have $c(\{u, v\}) = \text{red}$, for all $u \in X$; hence, $X \cup \{v\}$ contains a red m -graph, which implies $Rd(k, c, m)$ as well. If, on the other hand, $Bl(\#X, c, n)$ then we also have $Bl(k, c, n)$. In both cases we have $Rd(k, c, m) \vee Bl(k, c, n)$, which concludes this case.

Case $R(m, n-1) \leq \#Y$: By symmetry.

□

2.8.3 A few applications

A party containing 5 guests all knowing one another or all not knowing one another can now be represented as a complete graph containing a red 5-graph or a blue 5-graph. So, asking for the smallest such party is asking for the value of $R(5, 5)$.

As upper bound for $R(5, 5)$, Ramsey's theorem now gives $\binom{8}{4}$, which equals 70.

Further investigation of this problem has revealed that $R(5, 5) \in [43..49]$; what is the actual value of $R(5, 5)$ still is an open problem!

* * *

The smallest complete graph containing, independently of the colouring, at least one *monochrome* triangle is the complete 6-graph. Ramsey's theorem yields $R(3, 3) \leq 6$. For the complete 5-graph a colouring exists such that the graph does *not* contain a monochrome triangle; hence, $R(3, 3) \geq 6$. So, $R(3, 3) = 6$.

2.9 Trees

2.9.1 Undirected trees

An (*undirected*) *tree* is an undirected graph that is connected and acyclic. As we will see, on the one hand trees are the *smallest* connected graphs: removal of an edge from a tree always results in a graph that is not connected anymore. On the other hand, trees are the *largest* acyclic graphs: adding an additional edge to a tree results in a graph containing at least one cycle.

Although trees may be infinite, usually only finite trees are considered. Without further notice we confine our attention to finite trees.

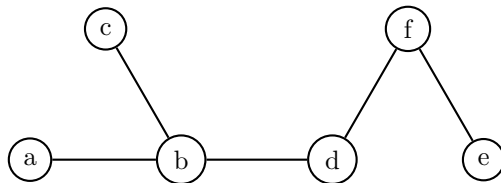


Figure 18: A (labelled) tree

In a connected graph, a *leaf* is a node with exactly one neighbour, that is, a node the degree of which equals 1. In Figure 18, for example, the leaves are a, c, e .

2.11 Lemma. Removal of a leaf and its (unique) associated edge from a connected graph yields a connected graph.

Proof. Let v be a leaf in a graph, and let u, w be nodes different from v . Then, no path connecting u and w contains v . Hence, every such path still exists in the graph resulting from removal of v and its associated edge.

□

2.12 Lemma. Every finite, acyclic, and connected graph with at least 2 nodes contains at least one leaf.

Proof. By contraposition: we prove that a finite, connected graph without leaves contains at least one cycle. Let (V, \sim) be such a graph, with $\#V \geq 2$. Because the graph is connected every node has at least one neighbour, so $\text{deg}(v) \geq 1$, for every node v ; because the graph contains no leaves we even have $\text{deg}(v) \geq 2$, for all v .

Now we construct an infinite sequence $s_{i:0 \leq i}$ of nodes, as follows. Choose $s_0 \in V$ arbitrarily, and choose $s_1 \in V$, such that $s_0 \sim s_1$. Note that this is possible because

V is assumed to have at least 2 nodes, and because $\deg(s_0) \geq 2$. Next, for all $i, 0 \leq i$, we choose $s_{i+2} \in V$ such that $s_{i+1} \sim s_{i+2}$ and $s_i \neq s_{i+2}$. This is possible because $\deg(s_{i+1}) \geq 2$, which follows from the assumption that the degree of every node is at least 2.

Thus, we have defined an infinite path s , starting at s_0 and with the property that $s_i \neq s_{i+2}$, for all $i, 0 \leq i$. The set V of nodes, however, is finite. Therefore⁴, we have $s_p = s_q$, for some p, q with $0 \leq p < q$. Hence, the sub-path $[s_p, \dots, s_q]$ is a cycle connecting s_p to itself, which concludes the proof.

□

A direct corollary of this lemma is that every tree with at least 2 nodes contains at least one leaf; after all, every tree is acyclic and connected.

2.13 Theorem. A tree with $n, 1 \leq n$, nodes contains $n-1$ edges.

Proof. By Mathematical Induction on n . A tree with 1 node has 0 edges – after all, every edge connects two *different* nodes –, and $1-1=0$. Now let (V, \sim) be a tree with $\#V = n+1$, where $1 \leq n$. By (the corollary to) the previous lemma this tree has a leaf u , say, so $\deg(u) = 1$. Hence, there is exactly one node v , say, with $u \sim v$, so the one-and-only edge involving u is $\{u, v\}$. Now let (W, \approx) be the graph obtained from (V, \sim) by removal of leaf u and its edge $\{u, v\}$. This means that $W = V \setminus \{u\}$ and that $w \approx x \Leftrightarrow w \sim x$, for all $w, x \in W$.

Because $v \in V$ we have $\#W = \#V - 1$, so $\#W = n$. The graph (W, \approx) is a tree because removal of node u and its edge $\{u, v\}$ maintains connectness of the remaining graph and, obviously, introduces no cycles. By Induction Hypothesis, tree (W, \approx) contains n edges, hence the original tree (V, \sim) contains $n+1$ edges.

□

Actually, (finite) trees can be characterized in many different way. This is illustrated by the following theorem, of which the above theorem is a special case.

2.14 Theorem. For a connected, undirected graph (V, E) the following propositions are mutually equivalent.

- (a) (V, E) is acyclic.
- (b) For every $e \in E$ the graph $(V, E \setminus \{e\})$ is not connected.
- (c) $\#E = \#V - 1$.
- (d) For all nodes u, v a *unique* path exists connecting u to v .

□

In addition the following properties deserve to be mentioned, as they are sometimes useful. Recall that, by definition, a graph is connected if every two nodes are connected by *at least* one path.

⁴See the discussion on finite and infinite, in the chapter on functions.

2.15 Lemma. An undirected graph is acyclic if and only if every two nodes are connected by *at most* one path.

2.16 Lemma. An undirected graph is a tree if and only if every two nodes are connected by *exactly* one path.

2.17 Lemma. For every undirected graph (V, E) :

- (a) $(\forall v: v \in V: \text{deg}(v) \geq 2) \Rightarrow \#E \geq \#V$
- (b) “ (V, E) is connected” $\Rightarrow \#E \geq \#V - 1$
- (c) “ (V, E) is acyclic” $\Rightarrow \#E \leq \#V - 1$

□

Notice that the latter two propositions provide another proof that if (V, E) is a tree then $\#E = \#V - 1$. Also, notice that the last proposition can also, by contraposition, be formulated as:

$$\#E \geq \#V \Rightarrow “(V, E) \text{ contains a cycle}” .$$

As a consequence, by combination with proposition (a) of this lemma, we obtain:

$$(\forall v: v \in V: \text{deg}(v) \geq 2) \Rightarrow “(V, E) \text{ contains a cycle}” .$$

2.9.2 Rooted trees

A *rooted tree* is a tree in which one node is identified separately. This designated node is called the *root* of the tree. For every node in a rooted tree we can define its *distance* as the length of the unique path connecting that node and the root. Thus, for example, the root itself has distance 0, and for every two neighbouring nodes, their distances differ by 1. In the latter case, the node with the smaller distance is *closer* to the root than the node with the larger distance. All edges in a rooted tree can now be turned into directed arrows, either by directing them *towards* the root or by directing them *away from* the root. The choice between these two options is rather irrelevant but must be made; therefore, in this text we adopt the convention that all arrows are directed away from the root. For example, Figure 16 gives the tree from Figure 15, as a rooted tree.

As an example of an infinite rooted tree, Figure 17 shows the tree obtained with the natural numbers as nodes, 0 as the root, and $\{(n, n+1) \mid 0 \leq n\}$ as the (directed) edges. Such a linear arrangement hardly deserves to be called a “tree”, of course, but formally it *is* a tree. Usually, however, such a linear arrangement is called a “list”.

A more interesting example is obtained as follows. The nodes are the positive naturals, the root is 1, and, for positive m, n , the pair (n, m) is an arrow if and only if $2 * n = m \vee 2 * n + 1 = m$. (This relation can also be formulated as $n = m \text{ div } 2$.) The graph thus obtained is a rooted tree: via the arrows every positive natural number can be obtained from 1 in a unique way. See Figure 18.

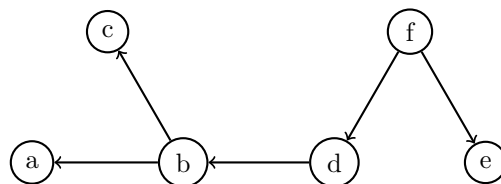
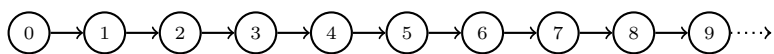
Figure 19: A rooted tree, with root f 

Figure 20: (Part of) the linear structure of the naturals

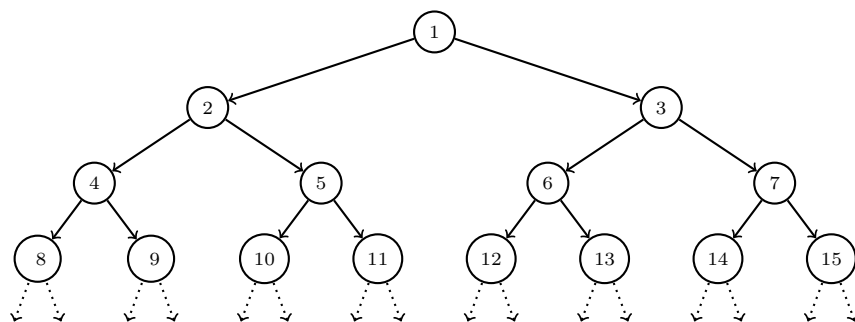


Figure 21: (Part of) the rooted, binary tree of the positive naturals

2.10 Exercises

1. How many edges does the complete n -graph have, for all $n \geq 1$? Prove the correctness of your answer.
2. An undirected graph (V, E) is called “regular of degree d ”, for some natural number d , if $\deg(v) = d$, for all $v \in V$. Prove that such a graph satisfies $d * \#V = 2 * \#E$.
3. Let (V, E) be an undirected graph satisfying $\#V = 9$ and $\#E \geq 14$. Prove that V contains at least one node the degree of which is at least 4.
4. Prove Lemma 2.6.
5. Prove Lemma 2.15.
6. Prove Lemma 2.16.
7. Prove Lemma 2.17.
8. Given are two connected (undirected) graphs (V, E) and (W, F) , such that $V \cap W = \emptyset$. Let $v \in V$ and $w \in W$. Prove that the graph $(V \cup W, E \cup F \cup \{\{v, w\}\})$ is connected.
9. Let (V, \sim) be a connected undirected graph such that $v, w \in V$ with the following properties: $\deg(v)$ and $\deg(w)$ are odd and $\deg(u)$ is even for all *other* nodes $u \in V$. Prove that the graph contains an Euler-path connecting v and w , that is, a path containing every edge of the graph exactly once.
10. A *chain* in a directed graph (V, \rightarrow) is an infinite sequence $s_i: 0 \leq i$ of nodes – that is, a function of type $\mathbb{N} \rightarrow V$ – with the property $(\forall i: 0 \leq i: s_i \rightarrow s_{i+1})$.
 - (a) Prove that finite and acyclic directed graphs do not contain chains.
 - (b) As a consequence, prove that every finite and acyclic directed graph contains at least one node the out-degree of which is zero.
11. We consider a (finite) undirected graph in which the degree of every node is at least 3. Prove that this graph contains a cycle containing at least 4 nodes.
12. We consider an (finite) undirected graph with n nodes, for $n \geq 3$. The degree of every node in this graph is at least n and the graph contains a node of degree $n-2$. Prove that this graph is connected.
13. We consider two undirected trees (V, E) and (W, F) , with $V \cap W = \emptyset$. Let $v_0, v_1 \in V$ and $w_0, w_1 \in W$. Prove that the graph $(V \cup W, E \cup F \cup \{\{v_0, w_0\}, \{v_1, w_1\}\})$ contains a cycle.
- * 14. We consider the complete 6-graph in which every edge has been coloured either red or blue. Prove that, independent of the chosen colouring, this graph contains at least 2 *monochrome* triangles.

15. Prove that an undirected graph with n nodes and in which the number of edges is greater than $n^2/4$ contains at least one triangle.
16. Give an example of an undirected graph, with at least 4 nodes, and containing a cycle that is both an Euler cycle and a Hamiltonian cycle.
17. Give an example of an undirected graph containing an Euler cycle and a Hamiltonian cycle that are different.
18. Give an example of a connected, undirected graph with 6 nodes, in which the degree of every node equals 3. Also give an example of such a graph in which the degree of every node equals 4.
19. Give an example of an undirected graph with 7 nodes, in which the degree of every node equals 2, and consisting of exactly 2 connected components.
20. Prove that every undirected graph, with 5 nodes and in which the degree of every node equals 2, is connected.
21. Prove that every acyclic, undirected graph, with n nodes and $n-1$ edges, is connected.
22. Prove that every undirected graph, with n nodes and at least $(n^2 - 3 * n + 4) / 2$ edges, is connected.
- * 23. Prove that every undirected graph, with n nodes and at least $(n^2 - 3 * n + 6) / 2$ edges, contains a Hamiltonian cycle.

3 Functions

3.1 Functions

Functions are about the most important building blocks in mathematical reasoning. Functions are almost everywhere. Examples are the well known functions $f \in \mathbb{R} \rightarrow \mathbb{R}$ given by $f(x) = x^2$, $f(x) = \sin(x)$, or $f(x) = \frac{1}{x^2+1}$. Actually, such functions are relations on \mathbb{R} : we say that $x \in \mathbb{R}$ is related to $y \in \mathbb{R}$ if and only if $f(x) = y$. Thus, every function $f \in \mathbb{R} \rightarrow \mathbb{R}$ corresponds to the relation $\{(x, y) \mid f(x) = y\}$. This is the traditional mathematical way to define functions.

3.1 Definition. A relation R from a set B to a set V is a *function* – also called a *mapping* – if (and only if) it has the following two properties.

- (a) For every $b \in B$ there is *at most one* $v \in V$ with bRv ;
- (b) For every $b \in B$ there is *at least one* $v \in V$ with bRv ;

Requirement (a) is called the requirement of *functionality*. This can be formalised as follows: $bRu \wedge bRv \Rightarrow u = v$, for all $b \in B$ and $u, v \in V$. In words, if $b \in B$ is in relation to some element of V this element is *unique*.

Requirement (b) is called the requirement of *totality*. It states that *every* $b \in B$ is related to some element of V .

If relation R is functional but not total then R is called a *partial function*, so a (truly) partial function only satisfies (a). To emphasize that a function is not partial, functions satisfying (a) and (b) also are called *total functions*. Notice that requirements (a) and (b) together state that for every $b \in B$ there is *exactly one* $v \in V$ satisfying bRv .

□

By default, functions are total, unless stated otherwise. For functions we usually (but not always) use names f, g, h, \dots or F, G, H, \dots . If f is a (partial or total) function from B to V we write $f(b)$ for the unique element $v \in V$ for which bfv , if it exists: so, $f(b) = v \Leftrightarrow bfv$.

If f is a total function we call B the *domain* of f . Now we have $f(b) \in V$ for *every* $b \in B$. If f is a partial function then f 's domain is the *largest* subset of B on which f is defined. That is, subset $A: A \subseteq B$ is f 's domain if and only if it satisfies both:

$$(\forall b: b \in A: (\exists v: v \in V: bfv)) \quad , \text{ and:}$$

$$(\forall b: b \in B \setminus A: \neg(\exists v: v \in V: bfv)) \quad .$$

So, for partial function f from B to V with domain A we have $f(b) \in V$ for *every* $b \in A$ and $f(b)$ is *undefined* for all b not in A . Notice that every partial function from B to V with domain A is a total function from A to V .

Combining these observations we conclude that every (partial or total) function f from B to V with domain A satisfies $(\forall b: b \in A: f(b) \in V)$, and $A = B$ if and only if f is total.

The set of all functions from B to V is denoted by $B \rightarrow V$, also by V^B . If f is a function in this set we also say that “ f has type $B \rightarrow V$ ”; also, $f \in B \rightarrow V$ and $f: B \rightarrow V$ are common ways to denote this.

3.2 Example. We already have encountered numerous examples of functions. Here are some familiar ones.

- (a) polynomial functions like $f \in \mathbb{R} \rightarrow \mathbb{R}$, with $f(x) = x^3$, for all $x \in \mathbb{R}$.
- (b) goniometric functions like \cos , \sin and \tan .
- (c) $\sqrt{\cdot} \in \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$, mapping the non-negative reals to their square roots.
- (d) $\ln \in \mathbb{R}^+ \rightarrow \mathbb{R}$, the natural logarithm.

3.2 Equality of functions

If two (total) functions with the same domain have equal values in all points of their common domain then these functions are equal. This is rendered formally as follows:

Principle of Extensionality: Functions f and g in $B \rightarrow V$, for some sets B and V , satisfy:

$$(\forall b: b \in B: f(b) = g(b)) \Rightarrow f = g$$

□

3.3 Monotonicity of function types

If $V \subseteq W$ then every (partial or total) function in $B \rightarrow V$ also is a (partial or total) function in $B \rightarrow W$. Hence, the set $B \rightarrow V$ of functions to V is a subset of the set $B \rightarrow W$ of functions to W .

Similarly, if $A \subseteq B$ then every function f in $B \rightarrow V$ satisfies: $f(b) \in V$ for every $b \in A$. This way, f can be considered as a function of type $A \rightarrow V$. Notice, however, that this is not true from a strictly formal point of view. Recalling that a function is a relation, function f , of type $B \rightarrow V$, contains a pair (b, v) , for every $b \in B$, so, also if $\neg(b \in A)$, whereas the functions in $A \rightarrow V$ only contain pairs (b, v) with $b \in A$. By omitting all pairs (b, v) with $\neg(b \in A)$, however, function f can be *restricted* to domain A . Calling the function thus obtained g , it is defined relationally by:

$$g = \{ (b, v) \mid b \in A \wedge f(b) = v \} .$$

Now function g has type $A \rightarrow V$, and it satisfies:

$$(\forall b: b \in A: g(b) = f(b)) .$$

In practical situations, however, we usually will not introduce a separate name for this restricted function g , and we simply state that if $A \subseteq B$ then every function in $B \rightarrow V$ also has type $A \rightarrow V$.

These properties can be summarized as follows.

3.3 Properties. For sets A, B, V, W :

$$V \subseteq W \Rightarrow B \rightarrow V \subseteq B \rightarrow W$$

$$A \subseteq B \Rightarrow B \rightarrow V \subseteq A \rightarrow V$$

□

3.4 Example. Every function in $\mathbb{R} \rightarrow \mathbb{R}$ also is a function in $\mathbb{N} \rightarrow \mathbb{R}$, and every function in $\mathbb{N} \rightarrow \mathbb{N}$ also is a function in $\mathbb{N} \rightarrow \mathbb{Z}$.

3.4 Function composition

We have seen earlier that if S is a relation from set B to set V and if T is a relation from V to a set Z then their composition, as denoted by $S;T$, is a relation from B to Z . The following lemma states that relational composition of two functions itself is a function again.

3.5 Lemma. If f is a function in $B \rightarrow V$ and if g is a function in $V \rightarrow Z$ then the relation $f;g$ is a function in $B \rightarrow Z$.

□

It is common mathematical practice to write the composition of two functions f and g as $g \circ f$, instead of as $f;g$. Notice, however, that the difference is purely notational, as we now have: $g \circ f = f;g$. So, according to the above lemma, if f has type $B \rightarrow V$ and if g has type $V \rightarrow Z$ then $g \circ f$ is a function of type $B \rightarrow Z$; it is defined by:

$$(g \circ f)(b) = g(f(b)) \text{ , for all } b \in B \text{ .}$$

3.6 Properties. The following properties (b) and (c) are directly inherited from the corresponding properties of relational composition.

- (a) On set B the identity relation I_B is a function from B to B , and $I_B(b) = b$, for every $b \in B$. Not surprisingly, I_B is also called the identity function on B .
- (b) I is the identity of function composition: if $f \in B \rightarrow V$ then $f \circ I_B = f$ and $I_V \circ f = f$.
- (c) Function composition is associative: $h \circ (g \circ f) = (h \circ g) \circ f$, for functions f, g, h of appropriate types.

3.5 Lifting a function

Let B and V be sets and let $f \in B \rightarrow V$ be a function from B to V . As stated earlier, set B is called the *domain* of f . Also, set V is called f 's *codomain*.

For $b \in B$ element $f(b)$ (in V) is called the *image* of b under f , or the *value of function f in point b* . The subset of V containing all values $f(b)$, for all $b \in B$, is called the *image* of set B under f .

The notion of image can be generalized to arbitrary subsets of B . For any subset $A: A \subseteq B$ the image of A under f is a subset of V , namely:

$$\{ f(b) \mid b \in A \} .$$

This subset depends, in a unique way, on subset A . So, we can define a function F , say, from the set of all subsets of B to the set of all subsets of V , as follows:

$$F(A) = \{ f(b) \mid b \in A \} , \text{ for all } A: A \subseteq B .$$

In common mathematical language the set of all subsets of set B , also called B 's *power set*, is denoted by $\mathcal{P}(B)$. Similarly, the power set of V is $\mathcal{P}(V)$. Thus, function F has type $\mathcal{P}(B) \rightarrow \mathcal{P}(V)$.

This function F also depends on f , of course. For every function f in $B \rightarrow V$ there is corresponding function F in $\mathcal{P}(B) \rightarrow \mathcal{P}(V)$, as defined above. We call F the *lifted* version of f . Function F has interesting⁵ algebraic properties.

3.7 Properties.

- (a) $F(\emptyset) = \emptyset$
- (b) $F(\{b\}) = \{f(b)\}$, for all $b \in B$
- (c) F *distributes over* arbitrary unions; that is, for any collection Ω of subsets of B – so, $\Omega \subseteq \mathcal{P}(B)$ –, we have:

$$F\left(\bigcup_{A: A \in \Omega} A\right) = \left(\bigcup_{A: A \in \Omega} F(A)\right)$$

- (d) F is *monotonic*; that is, for all subsets $A0, A1$ of B :

$$A0 \subseteq A1 \Rightarrow F(A0) \subseteq F(A1)$$

- (e) If F is lifted f and if G is lifted g , then $G \circ F$ is lifted $g \circ f$.
- (f) For every subset $A \subseteq B$ and subset $U \subseteq V$:

$$F(A) \subseteq U \Leftrightarrow (\forall b: b \in A: f(b) \in U) .$$

□

⁵In the chapter on partial orders we will see why.

Notice that Property (b) shows that, in turn, function F uniquely determines function f from which F was obtained. In a (not too strict) way, F is a *generalisation* of f , of which f can be considered an instance.

A function f , of type $B \rightarrow V$, and its lifted version, of type $\mathcal{P}(B) \rightarrow \mathcal{P}(V)$ and called F above, are entirely different functions. Nevertheless, it is common mathematical practice to *denote both* by the same name f ; so, instead of $f(b)$ and $F(A)$ we write $f(b)$ and $f(A)$. This, so-called, *overloading* of the name f is rather harmless, but the interpretation of an expression like $f(x)$ now depends on the type of x : if $x \in B$ the expression just means $f(x)$, but if $x \subseteq B$, that is: $x \in \mathcal{P}(B)$, then the expression means $F(x)$.

We have called the set of function values $f(b)$, for all $b \in B$, the *image* or *range* of f . In terms of lifted f the image of f just is $f(B)$.

3.8 Properties. Functions $f \in B \rightarrow V$ and $g \in V \rightarrow Z$ satisfy:

$$(a) \quad (g \circ f)(B) = g(f(B)) \quad .$$

$$(b) \quad (g \circ f)(B) \subseteq g(V) \quad .$$

□

* * *

An element $b \in B$ satisfying $f(b) = v$, for some $v \in V$, is called an *original* of v under function f . As we know, because f is a function every $b \in B$ is related to a unique value, written as $f(b)$, in V . Conversely, it is not necessarily true that every $v \in V$ has a unique original in B : for any $v \in V$ set B may contain 0, 1, or many elements b satisfying $f(b) = v$. The *whole set* of such elements b , however, is unique. That is, by lifting again, we can define a function G , say, of type $\mathcal{P}(V) \rightarrow \mathcal{P}(B)$, by:

$$G(U) = \{ b \mid b \in B \wedge f(b) \in U \} \quad , \text{ for all } U : U \subseteq V \quad .$$

Set $G(U)$ is called the *pre-image* under f of set U . In particular, for any value $v \in V$, the set of all elements $b \in B$ satisfying $f(b) = v$ now is $G(\{v\})$; so, the pre-image of $\{v\}$ is the set of all originals of v .

In common mathematical notation $G(U)$ is written as $f^{-1}(U)$. As we will see later, some functions f have the property that, for every $v \in V$ there is a unique $b \in B$ with $f(b) = v$. Then, a function exists, of type $V \rightarrow B$, mapping every v to this unique b . This function is called f 's *inverse* and it is denoted by f^{-1} . Function G , as defined here on sets, then is the lifted version of f^{-1} , which is why it is also written as f^{-1} . So, generally we have:

$$f^{-1}(U) = \{ b \mid b \in B \wedge f(b) \in U \} \quad , \text{ for all } U : U \subseteq V \quad ,$$

and if function f has an inverse we have, for all $b \in B$ and $v \in V$:

$$b = f^{-1}(v) \Leftrightarrow f(b) = v \quad , \text{ and:}$$

$$\{f^{-1}(v)\} = f^{-1}(\{v\}) ,$$

3.9 Example.

- (a) Let $f \in \mathbb{R} \rightarrow \mathbb{R}$ with $f(x) = x^2$ for all $x \in \mathbb{R}$. Then $f^{-1}([0..4]) = [-2..2]$.
- (b) Consider the function $\text{mod } 8$ in $\mathbb{Z} \rightarrow \mathbb{Z}$. The originals of 3 then are the elements of the set $\{\dots, -13, -5, 3, 11, 19, \dots\}$.

3.10 Theorem.

Every function $f \in B \rightarrow V$ satisfies:

- (a) $A \subseteq f^{-1}(f(A))$, for all $A : A \subseteq B$;
- (b) $f(f^{-1}(U)) \subseteq U$, for all $U : U \subseteq V$.

□

3.11 Example.

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by $f(x) = x^2$, for all $x \in \mathbb{R}$. Then the range $f^{-1}(f([0..1]))$ equals $[-1..1]$, which properly contains $[0..1]$. Moreover, we have $f^{-1}(f([-4..4])) = [0..4]$, which is properly contained in $[-4..4]$. This shows that we can have strict inclusions in the above theorem.

3.6 Surjective, injective, and bijective functions

Some functions have additional and useful properties. Here we define some of them.

3.12 Definition.

A function f in $B \rightarrow V$ is *surjective* if every element in V is the value of f for *at least* one value in B , that is, if:

$$(\exists b : b \in B : f(b) = v) , \text{ for all } v \in V .$$

A function f in $B \rightarrow V$ is *injective* if every element in V is the value of f for *at most* one value in B , that is, if:

$$f(a) = f(b) \Rightarrow a = b , \text{ for all } a, b \in B .$$

A function f in $B \rightarrow V$ is *bijective* if it is both surjective and injective. Hence, a function is bijective if every element in V is the value of f for *exactly one* value in B .

□

3.13 Lemma.

For function f in $B \rightarrow V$: “ f is surjective” $\Leftrightarrow f(B) = V$.

□

3.14 Example.

This example illustrates that the “same” function is or is not surjective or injective, depending on which domain one considers.

- (a) The function $\sin : \mathbb{R} \rightarrow \mathbb{R}$ is neither surjective nor injective.
- (b) The function $\sin : [-\pi/2.. \pi/2] \rightarrow \mathbb{R}$ is injective and not surjective.

(c) The function $\sin : \mathbb{R} \rightarrow [-1..1]$ is surjective and not injective.

(d) The function $\sin : [-\pi/2.. \pi/2] \rightarrow [-1..1]$ is bijective.

□

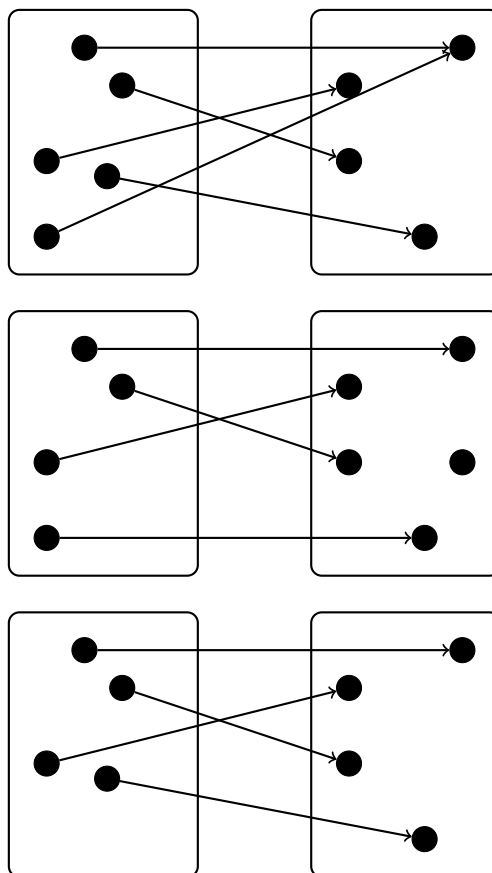


Figure 22: Surjective, injective and bijective functions

* * *

In the definition of “function”, in the beginning of this chapter, we have stated the requirements of “totality” and “functionality”. The notion of “surjectivity” is equivalent to the notion of “totality”, but with sets B and V interchanged. Similarly, the notion of “injectivity” is equivalent to the notion of “functionality”, but with B and

V interchanged⁶. Recall that for every relation R from B to V its transposition R^T is the relation from V to B defined by:

$$v R^T b \Leftrightarrow b R v, \text{ for all } b, v .$$

Function f being surjective now means that relation f^T is total, and f being injective means that relation f^T is functional. From this we conclude immediately that if f is bijective relation then f^T is a function from V to B . This function happens to be f 's *inverse*, and it is denoted by f^{-1} .

3.15 Definition. Function g in $V \rightarrow B$ is an *inverse* of function f in $B \rightarrow V$ if:

$$g \circ f = I_B \wedge f \circ g = I_V .$$

3.16 Lemma. A function has *at most one* inverse, that is, if g and h both are inverses of f , then $g = h$.

3.17 Lemma. Function g is the *inverse* of function f if and only if f is the inverse of g .

3.18 Lemma. A function f has an inverse if and only if f is bijective; for bijective f its inverse is f^{-1} .
□

3.19 Theorem. All functions f in $B \rightarrow V$ and g in $V \rightarrow Z$ satisfy:

- (a) If f and g are surjective, then so is $g \circ f$;
- (b) If f and g are injective, then so is $g \circ f$;
- (c) If f and g are bijective, then so is $g \circ f$.

Proof.

- (a) Assuming f and g to be surjective, it suffices, according to Lemma 3.13, to show that $(g \circ f)(B) = Z$:

$$\begin{aligned} & (g \circ f)(B) \\ \Leftrightarrow & \quad \{ \text{Property 3.8(a)} \} \\ & g(f(B)) \\ \Leftrightarrow & \quad \{ \text{“}f \text{ is surjective”}: \text{Lemma 3.13} \} \\ & g(V) \\ \Leftrightarrow & \quad \{ \text{“}g \text{ is surjective”}: \text{Lemma 3.13} \} \\ & Z . \end{aligned}$$

⁶Instead of “surjective” and “injective” we could have used “right-total” and “right-functional”.

(b) Assuming f and g to be injective, we derive, for $a, b \in B$:

$$\begin{aligned}
 & (g \circ f)(a) = (g \circ f)(b) \\
 \Leftrightarrow & \quad \{ \text{definition of composition} \} \\
 & g(f(a)) = g(f(b)) \\
 \Rightarrow & \quad \{ \text{"}g \text{ is injective"} \} \\
 & f(a) = f(b) \\
 \Rightarrow & \quad \{ \text{"}f \text{ is injective"} \} \\
 & a = b .
 \end{aligned}$$

(c) This follows directly from (a) and (b).

□

3.20 Lemma. If functions f in $B \rightarrow V$ and g in $V \rightarrow Z$ both are bijective then

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

□

3.7 Some counting arguments

The elements of finite sets can be counted. That is, we can define a function $\#$, say, that maps every finite set to its *number of elements*. So, for finite set V its number of elements $\#V$ is a natural number. This function can be defined recursively, as follows. After all, every finite set either is \emptyset or is $U \cup \{v\}$, for some smaller finite set U and some value v not in U .

3.21 Definition. For every finite set U and value v with $v \notin U$:

$$\begin{aligned}
 \#\emptyset &= 0 \\
 \#(U \cup \{v\}) &= \#U + 1
 \end{aligned}$$

□

Without proof we state that $\#$ has the following properties, which can be proved by induction on the structure of the finite sets, if so desired.

3.22 Properties. For all finite sets U, V and value v :

- (0) $\#U = 0 \Leftrightarrow U = \emptyset$
- (1) $\#\{v\} = 1$
- (2) $0 \leq \#U$

- (3) $\#(U \cup V) = \#U + \#V - \#(U \cap V)$
 (4) $\#(U \cup V) \leq \#U + \#V$
 (5) $\#(U \cup V) = \#U + \#V$, if $U \cap V = \emptyset$
 (6) $\#V = \#U + \#(V \setminus U)$, if $U \subseteq V$
 (7) $\#U \leq \#V$, if $U \subseteq V$
 (8) $\#U = \#V \Leftrightarrow U = V$, if $U \subseteq V$

□

Notice that these properties are mutually dependent. For example, property (4) follows directly from (3) and (2), property (5) is an instance of (3), and (6) is an instance of (5), whereas (7) follows from (6) and (2).

* * *

Functions on finite sets yield finite images. In what follows function f has type $B \rightarrow V$, where B and V are finite sets.

3.23 Lemma. $\#(f(B)) \leq \#B$

Proof. By induction on the structure of B . If $B = \emptyset$ then $f(B) = \emptyset$ too, hence, in this case we even have $\#(f(B)) = \#B$. If $B = A \cup \{b\}$, for some b not in A we have:

$$\begin{aligned}
 & \#(f(A \cup \{b\})) \\
 = & \quad \{ f \text{ over } \cup \} \\
 & \#(f(A) \cup f(\{b\})) \\
 \leq & \quad \{ \text{property 3.22 (4)} \} \\
 & \#(f(A)) + \#(f(\{b\})) \\
 = & \quad \{ \text{definition of } f(\{b\}) \} \\
 & \#(f(A)) + \#(\{f(b)\}) \\
 = & \quad \{ \text{property 3.22 (1) (twice)} \} \\
 & \#(f(A)) + \#(\{b\}) \\
 \leq & \quad \{ \text{Induction Hypothesis} \} \\
 & \#A + \#(\{b\}) \\
 = & \quad \{ \text{property 3.22 (5), using } \neg(b \in A) \} \\
 & \#(A \cup \{b\})
 \end{aligned}$$

□

3.24 Lemma. “ f is injective” $\Leftrightarrow \#(f(B)) = \#B$

Proof. By mutual implication.

“ \Rightarrow ”: Very similar to the proof of Lemma 3.23, using that $\#(f(A) \cup f(\{b\}))$ now is equal to $\#(f(A)) + \#(f(\{b\}))$, because for injective f we have $\neg(f(b) \in f(A))$ if $\neg(b \in A)$.

“ \Leftarrow ”: By contraposition, that is, if f is not injective then $\#(f(B)) \neq \#B$. If f is not injective then B contains elements a, b , say, with $a \neq b$ and $f(a) = f(b)$. Now:

$$\begin{aligned} & \#(f(B)) \\ = & \quad \{ a \neq b \text{ and } f(a) = f(b), \text{ so } f(B \setminus \{a\}) = f(B) \} \\ & \#(f(B \setminus \{a\})) \\ \leq & \quad \{ \text{Lemma 3.23} \} \\ & \#(B \setminus \{a\}) \\ = & \quad \{ a \in B \} \\ & \#B - 1 \quad , \end{aligned}$$

from which we conclude that $\#(f(B)) < \#B$, hence, also $\#(f(B)) \neq \#B$.

□

3.25 Lemma. “ f is injective” $\Rightarrow \#B \leq \#V$

Proof.

$$\begin{aligned} & \#B \\ = & \quad \{ f \text{ is injective: Lemma 3.24} \} \\ & \#(f(B)) \\ \leq & \quad \{ f(B) \subseteq V: \text{property 3.22(7)} \} \\ & \#V \end{aligned}$$

□

3.26 Lemma. “ f is surjective” $\Rightarrow \#V \leq \#B$

Proof.

$$\begin{aligned} & \#V \\ = & \quad \{ f \text{ is surjective: } f(B) = V \} \\ & \#(f(B)) \\ \leq & \quad \{ \text{Lemma 3.23} \} \\ & \#B \end{aligned}$$

□

Corollary: “ f is bijective” $\Rightarrow \#V = \#B$

□

Aside: Without further proofs we observe that the converses to the two latter lemmata hold as well.

Properties.

(0) $\#B \leq \#V \Rightarrow (\exists f: f \in B \rightarrow V: \text{“}f \text{ is injective”})$

(1) $\#V \leq \#B \Rightarrow (\exists f: f \in B \rightarrow V: \text{“}f \text{ is surjective”})$

□

Now we are ready for the main theorem of this subsection.

3.27 Theorem. [Pigeonhole Principle] For f a function of type $B \rightarrow V$, for finite sets B, V , and if $\#B = \#V$, we have:

“ f is injective” \Leftrightarrow “ f is surjective”

Proof.

“ f is injective”

\Leftrightarrow { Lemma 3.24 }

$\#(f(B)) = \#B$

\Leftrightarrow { $\#B = \#V$ }

$\#(f(B)) = \#V$

\Leftrightarrow { $f(B) \subseteq V$: property 3.22 (8) }

$f(B) = V$

\Leftrightarrow { Lemma 3.13 }

“ f is surjective”

□

3.28 Remark. The above result is called the pigeonhole principle because of the following. If one has n pigeons (the set A) and the same number of holes (the set B), then one pigeonhole is empty if and only if one of the other holes contains at least two pigeons.

3.29 Example. Suppose p and q are two different prime numbers. We consider the function φ in $[0..p) \rightarrow [0..p)$, defined by $\varphi(x) = (x * q) \bmod p$, for all $x \in [0..p)$.

We prove that φ is a bijection. By the Pigeonhole Principle it suffices to show that φ is injective. To this end we derive, for $x, y \in [0..p)$:

$\varphi(x) = \varphi(y)$

\Leftrightarrow { definition of φ }

$$\begin{aligned}
& (x * q) \bmod p = (y * q) \bmod p \\
\Leftrightarrow & \quad \{ \text{property of mod} \} \\
& ((x-y) * q) \bmod p = 0 \\
\Leftrightarrow & \quad \{ \text{property of mod and } | \} \\
& p \mid ((x-y) * q) \\
\Leftrightarrow & \quad \{ \text{"p is prime": } (p \mid) \text{ distributes over } * \} \\
& p \mid (x-y) \vee p \mid q \\
\Leftrightarrow & \quad \{ \text{"p and q are prime" and } p \neq q, \text{ hence: } \neg(p \mid q) \} \\
& p \mid (x-y) \\
\Leftrightarrow & \quad \{ x, y \in [0..p), \text{ hence } -p < x-y < p \} \\
& x = y \quad .
\end{aligned}$$

□

3.8 Of finite and infinite

Every now and then we encounter in our discussions the notions of *finite* or *infinite* sets. If we wish our proofs to be really formal, we need a formal definition of what constitutes a finite or infinite set.

Fortunately, this is not a study on the foundation of Mathematics; therefore, we can get away with a very pragmatic attitude: without much ado, we consider the natural numbers as (the prototype of) an infinite set and we consider any set infinite in which the natural numbers can be embedded. Formally, this means that an *injective* function from the naturals into that set exist. In addition, we define a set to be finite if it is not infinite⁷.

3.30 Definition. Set V is infinite if and only if a function f in $\text{Nat} \rightarrow V$ exists such that all function values are *different*, that is, such that:

$$(\forall i, j: 0 \leq i < j: f_i \neq f_j)$$

3.31 Definition. Set V is finite if and only if it is not infinite. By application of the rules of The Morgan to the definition of “infinite” we obtain that V is finite if and only if *for all* functions f in $\text{Nat} \rightarrow V$ we have:

$$(\exists i, j: 0 \leq i < j: f_i = f_j)$$

□

Functions on the naturals are also called “infinite sequences”. The above definition of infinity can, therefore, also be remembered as follows: a set is infinite if it contains

⁷The notions of finite and infinite go hand in hand: it is virtually impossible to define the one without the other.

an infinite sequence all whose elements are different. And, conversely, a set is finite if every infinite sequence in the set contains at least one value twice.

A different way to define infinity is: a set is infinite if every finite subset of it is not the whole set. An apparent advantage of this definition is that it does not seem to refer to the natural numbers, but that is not true: this definition refers to “finite subsets”, the definition of which still involves, one way or the other, the naturals.

Nevertheless, this latter definition of infinity sometimes is useful too⁸, but the above definition in terms of functions generally is more practical.

3.32 Example. A rather trivial example is: the natural numbers themselves are infinite. Just take for function f the identity function, as defined by $(\forall i: 0 \leq i: f_i = i)$. Then, obviously, all function values are different.
□

3.33 Properties. If set V is infinite and if $V \subseteq W$ then set W is infinite too. This proposition is equivalent to the proposition that every subset of a finite set is finite as well.

3.9 A Useful Classification

3.9.1 Several kinds of relations

Throughout this little essay we study relations from a given set B to a given set V , so we study subsets of the cartesian product $B \times V$. For R such a relation we write bRv as an abbreviation of $(b, v) \in R$, as usual.

For R such a relation we define four properties that R may or may not have:

- (2) “ R is L-total” $\Leftrightarrow (\forall b: b \in B: (\exists v: v \in V: bRv))$
- (3) “ R is L-functional” $\Leftrightarrow (\forall b, u, v: b \in B \wedge u, v \in V: bRu \wedge bRv \Rightarrow u = v)$
- (4) “ R is R-total” $\Leftrightarrow (\forall v: v \in V: (\exists b: b \in B: bRv))$
- (5) “ R is R-functional” $\Leftrightarrow (\forall b, c, v: b, c \in B \wedge v \in V: bRv \wedge cRv \Rightarrow b = c)$

Note that “ R is L-total” expresses that every element in B is related (by R) to *at least one* element in V , whereas “ R is L-functional” expresses that every element in B is related to *at most one* element in V . Hence, together – in conjunction – they express that every element in B is related to *exactly one* element in V .

Similarly, “ R is R-total” expresses that every element in V is related (by R) to *at least one* element in B , whereas “ R is R-functional” expresses that every element in V is related to *at most one* element in B . So, together they express that every element in V is related to *exactly one* element in B .

Also note the symmetry between (2) and (3) on the one hand and (4) and (5) on the other hand: the two pairs are transformed into one another under relation transposition. That is, for every relation R we have:

⁸Example: the well-known proof that the set of prime numbers is infinite.

$$\begin{aligned}
\text{“}R \text{ is L-total”} &\Leftrightarrow \text{“}R^T \text{ is R-total”} \text{ , and:} \\
\text{“}R \text{ is L-functional”} &\Leftrightarrow \text{“}R^T \text{ is R-functional”} \text{ .}
\end{aligned}$$

As far as I know, Netty van Gasteren, in her PhD-thesis, was the first to emphasize the importance of these four properties (although she used different names).

3.9.2 Several kinds of functions

Using the terminology introduced above we can now classify several kinds of functions, in the following way. Again, R is a relation from B to V .

$$\begin{aligned}
\text{“}R \text{ is a (partial) function (from } B \text{ to } V\text{)”} &\Leftrightarrow \\
\text{“}R \text{ is L-functional”} &\text{ ,} \\
\text{“}R \text{ is a (total) function”} &\Leftrightarrow \\
\text{“}R \text{ is L-total”} \wedge \text{“}R \text{ is L-functional”} &\text{ ,} \\
\text{“}R \text{ is a surjective function”} &\Leftrightarrow \\
\text{“}R \text{ is L-total”} \wedge \text{“}R \text{ is L-functional”} \wedge \text{“}R \text{ is R-total”} &\text{ ,} \\
\text{“}R \text{ is an injective function”} &\Leftrightarrow \\
\text{“}R \text{ is L-total”} \wedge \text{“}R \text{ is L-functional”} \wedge \text{“}R \text{ is R-functional”} &\text{ ,} \\
\text{“}R \text{ is a bijection”} &\Leftrightarrow \\
\text{“}R \text{ is L-total”} \wedge \text{“}R \text{ is L-functional”} \wedge \text{“}R \text{ is R-total”} \wedge \text{“}R \text{ is R-functional”} &\text{ .}
\end{aligned}$$

In particular this classification clarifies the relation between functions and their inverses (if they exist). We have, for instance, that relation R is a bijection if and only if both R and R^T are (total) functions, which then are each other’s inverses, as we will show.

3.9.3 Applications

To demonstrate the usefulness of the above concepts we formulate and prove two lemmata.

3.34 Lemma. Every bijection R satisfies: $R;R^T = I_B$, where “ $;$ ” denotes relation composition and where I_B denotes the identity relation (on $B \times B$). Similarly, $R^T;R = I_V$.

Proof. We prove $R;R^T = I_B$ element-wise, that is, we prove $(\forall b, c: b, c \in B: b(R;R^T)c \Leftrightarrow b = c)$, as follows, for any $b, c \in B$:

$$\begin{aligned}
&b(R;R^T)c \\
\Leftrightarrow &\quad \{ \text{definition of } ; \} \\
&(\exists v: v \in V: bRv \wedge vR^Tc) \\
\Leftrightarrow &\quad \{ \text{definition of } \tau \}
\end{aligned}$$

$$\begin{aligned}
& (\exists v : v \in V : bRv \wedge cRv) \\
\Leftrightarrow & \quad \{ R \text{ is R-functional: (5) } \} \\
& (\exists v : v \in V : bRv \wedge cRv \wedge b=c) \\
\Leftrightarrow & \quad \{ \text{Leibniz, to simplify} \} \\
& (\exists v : v \in V : bRv \wedge b=c) \\
\Leftrightarrow & \quad \{ \wedge \text{ distributes over } \exists \} \\
& (\exists v : v \in V : bRv) \wedge b=c \\
\Leftrightarrow & \quad \{ R \text{ is L-total: (2) } \} \\
& b=c .
\end{aligned}$$

Notice that, in this proof, we only have used that R is R-functional and L-total. The other two properties, that R is L-functional and R-total, are needed for the proof of the symmetric counterpart $R^T ; R = I_V$.

□

The next lemma actually used to be an exercise which students tend to consider a difficult one.

3.35 Lemma. For R a relation from B to V and S a relation from V to B we have: if $R ; S = I_B$ and $S ; R = I_V$ then both R and S are bijections.

Proof. The problem is entirely symmetric in R and S , in that it is invariant under the substitution $B, V, R, S := V, B, S, R$. Hence, it suffices to prove that R is a bijection. As we have seen, this means that R has the four properties (2) through (5). The problem also is symmetric in that it is invariant under relation transposition. Hence, as we have seen, it suffices to prove that R is L-total and L-functional only. To be able to do so, however, we must rephrase our premisses $R ; S = I_B$ and $S ; R = I_V$ in a way that does not involve relation composition anymore, because properties (2) through (5) do not contain compositions.

Rewriting $R ; S = I_B$ as the conjunction of $R ; S \subseteq I_B$ and $I_B \subseteq R ; S$ and $S ; R = I_V$ as the conjunction of $S ; R \subseteq I_V$ and $I_V \subseteq S ; R$ respectively, and eliminating the compositions and identities – by applying their definitions – our premisses boil down to the following four given properties:

$$\begin{aligned}
(6) \quad & (\forall b, c : b, c \in B : (\exists v : v \in V : bRv \wedge vSc) \Rightarrow b=c) \\
(7) \quad & (\forall b : b \in B : (\exists v : v \in V : bRv \wedge vSb)) \\
(8) \quad & (\forall u, v : u, v \in V : (\exists b : b \in B : uSb \wedge bRv) \Rightarrow u=v) \\
(9) \quad & (\forall v : v \in V : (\exists b : b \in B : vSb \wedge bRv))
\end{aligned}$$

Now we prove that R is L-total, as follows:

$$\begin{aligned}
& R \text{ is L-total} \\
\Leftrightarrow & \quad \{ (2) \}
\end{aligned}$$

$$\begin{aligned}
& (\forall b: b \in B: (\exists v: v \in V: bRv)) \\
\Leftarrow & \quad \{ \text{strengthening, in view of (7)} \} \\
& (\forall b: b \in B: (\exists v: v \in V: bRv \wedge vSb)) \\
\Leftarrow & \quad \{ (7) \} \\
& \text{true} ,
\end{aligned}$$

and, hence, by the aforementioned symmetries, we conclude that S is R-total too. Now, we prove that R is L-functional, as follows, using definition (3), for all $b \in B$ and $u, v \in V$:

$$\begin{aligned}
& bRu \wedge bRv \\
\Leftarrow & \quad \{ S \text{ is R-total} \} \\
& (\exists w: w \in V: wSb) \wedge bRu \wedge bRv \\
\Leftarrow & \quad \{ \wedge \text{ distributes over } \exists \} \\
& (\exists w: w \in V: wSb \wedge bRu) \wedge bRv \\
\Leftarrow & \quad \{ (8), \text{ with } u, v := w, u \} \\
& (\exists w: w \in V: wSb \wedge bRu \wedge w = u) \wedge bRv \\
\Leftarrow & \quad \{ \text{one-point rule} \} \\
& uSb \wedge bRu \wedge bRv \\
\Rightarrow & \quad \{ \text{weakening} \} \\
& uSb \wedge bRv \\
\Rightarrow & \quad \{ \exists\text{-introduction} \} \\
& (\exists b: b \in B: uSb \wedge bRv) \\
\Rightarrow & \quad \{ (8) \} \\
& u = v ,
\end{aligned}$$

as required. Notice that a step like the first one in this calculation is unavoidable: the proof obligation is about R only, so we must introduce, one way or another, relation S into the game, in order to be able to use the premisses (6) through (9), all of which involve both R and S .

□

3.9.4 A more abstract and more algebraic approach

In the proof of Lemma 1 I have observed that the premisses of the lemma contain relation composition, whereas the proof obligations, amounting to (2) through (5), do not contain relation composition. Therefore, so I decided, we must eliminate relation composition from the premisses, which gave rise to formulae (6) through (9).

When writing down this observation, however, it suddenly dawned upon me that there is an alternative approach: instead of *eliminating* relation composition from the

premisses we might as well *introduce* it into formulae (2) through (5), thus lifting the whole discussion to a more abstract, element-free, level.

To be able to do so, I need a general (logical) property which I will refer to either as “ \exists -elimination” or as “ \exists -introduction”. In a natural-deduction style of logic this rule, when applied from left to right, is usually called “ \exists -elimination”. Note, however, that the rule actually expresses an equivalence which, therefore, may be applied from right to left as well, as a way of “ \exists -introduction”.

\exists -elimination: For any predicate $P(x)$, possibly containing x as a free variable, and for any predicate Q , in which x does *not* occur as free variable, we have:

$$(\exists x :: P(x)) \Rightarrow Q \quad \Leftrightarrow \quad (\forall x :: P(x) \Rightarrow Q) \quad .$$

proof: By straightforward calculation.

□

Now we can translate the classification from Section 3.9.1 into relational form, as follows:

$$\begin{aligned} & R \text{ is L-total} \\ \Leftrightarrow & \quad \{ (2) \} \\ & (\forall b : b \in B : (\exists v : v \in V : bRv)) \\ \Leftrightarrow & \quad \{ \text{idempotence of } \wedge \text{ and transposition} \} \\ & (\forall b : b \in B : (\exists v : v \in V : bRv \wedge vR^T b)) \\ \Leftrightarrow & \quad \{ \text{definition of } ; \} \\ & (\forall b : b \in B : b(R; R^T)b) \\ \Leftrightarrow & \quad \{ \text{one-point rule} \} \\ & (\forall b, c : b, c \in B : b = c \Rightarrow b(R; R^T)c) \\ \Leftrightarrow & \quad \{ \text{definition of } I_B \} \\ & (\forall b, c : b, c \in B : b(I_B)c \Rightarrow b(R; R^T)c) \\ \Leftrightarrow & \quad \{ \text{definition of } \subseteq \} \\ & I_B \subseteq R; R^T \quad . \end{aligned}$$

So, relation R is L-total if and only if $I_B \subseteq R; R^T$, and this may be taken as an alternative definition. Notice, however, that the second step in the above derivation – introduction of R^T – is rather arbitrary. This is a true design decision, but not the only possibility. Instead, for example, we could also have rewritten bRv to $bRv \wedge v \top b$, where \top denotes the maximal relation, that is the whole set $V \times B$. This gives rise to yet another alternative definition: R is L-total if and only if $I_B \subseteq R; \top$. As a matter of fact, there are good reasons to retain both variants, and performing the same exercises with R-totality we obtain:

$$(10) \quad \text{“}R \text{ is L-total”} \quad \Leftrightarrow \quad I_B \subseteq R; R^T$$

- (11) “ R is L-total” $\Leftrightarrow I_B \subseteq R; \top$
 (12) “ R is R-total” $\Leftrightarrow I_V \subseteq R^T; R$
 (13) “ R is R-total” $\Leftrightarrow I_V \subseteq \top; R$

Note that (11) is weaker than (10) and that (13) is weaker than (12): when we need to *prove* totality the weaker forms are to be preferred, but when we *use* totality the stronger forms are more attractive.

* * *

As to functionality we proceed as follows:

$$\begin{aligned}
 & R \text{ is L-functional} \\
 \Leftrightarrow & \quad \{ (3) \} \\
 & (\forall b, u, v : b \in B \wedge u, v \in V : bRu \wedge bRv \Rightarrow u = v) \\
 \Leftrightarrow & \quad \{ \text{dummy } b \text{ does not occur in } u = v : \text{nesting} \} \\
 & (\forall u, v : u, v \in V : (\forall b : b \in B : bRu \wedge bRv \Rightarrow u = v)) \\
 \Leftrightarrow & \quad \{ \exists\text{-introduction} \} \\
 & (\forall u, v : u, v \in V : (\exists b : b \in B : bRu \wedge bRv) \Rightarrow u = v) \\
 \Leftrightarrow & \quad \{ \text{transposition} \} \\
 & (\forall u, v : u, v \in V : (\exists b : b \in B : uR^T b \wedge bRv) \Rightarrow u = v) \\
 \Leftrightarrow & \quad \{ \text{definition of } ; \text{ and } I_V \} \\
 & (\forall u, v : u, v \in V : u(R^T; R)v \Rightarrow u(I_V)v) \\
 \Leftrightarrow & \quad \{ \text{definition of } \subseteq \} \\
 & R^T; R \subseteq I_V .
 \end{aligned}$$

R-functionality can, of course, be rewritten in exactly the same way. Thus, we obtain:

- (14) “ R is L-functional” $\Leftrightarrow R^T; R \subseteq I_V$
 (15) “ R is R-functional” $\Leftrightarrow R; R^T \subseteq I_B$

* * *

Now let us see how the two lemmata from the previous section can be proved in terms of the above definitions.

Lemma 3.34: Every bijection R satisfies $R; R^T = I_B$ and $R^T; R = I_V$.

proof: This is rather trivial now:

$$\begin{aligned}
 & R; R^T = I_B \\
 \Leftrightarrow & \quad \{ \text{set equality via mutual inclusion} \}
 \end{aligned}$$

$$\begin{aligned}
& R;R^T \subseteq I_B \wedge I_B \subseteq R;R^T \\
\Leftrightarrow & \{ R \text{ is a bijection, so } R \text{ is R-functional, (15), and } R \text{ is L-total, (10) } \} \\
& \text{true .}
\end{aligned}$$

□

Lemma 3.35: For R a relation from B to V and S a relation from V to B we have: if $R;S = I_B$ and $S;R = I_V$ then both R and S are bijections.

proof: As we have observed already in the previous proof, it suffices to prove that R is L-total and L-functional; obviously, here we use the weaker form of totality:

$$\begin{aligned}
& R \text{ is L-total} \\
\Leftrightarrow & \{ (11) \} \\
& I_B \subseteq R;\top \\
\Leftrightarrow & \{ R;S = I_B \} \\
& R;S \subseteq R;\top \\
\Leftarrow & \{ ; \text{ is monotonic } \} \\
& S \subseteq \top \\
\Leftrightarrow & \{ \text{definition of } \top \} \\
& \text{true .}
\end{aligned}$$

Notice that use of the weaker definition, (11), really helps here: this leads to $S \subseteq \top$, which is easy, whereas the stronger definition would have given rise to $S \subseteq R^T$, which requires more work.

Furthermore:

$$\begin{aligned}
& R \text{ is L-functional} \\
\Leftrightarrow & \{ (14) \} \\
& R^T;R \subseteq I_V \\
\Leftrightarrow & \{ S;R = I_V \} \\
& R^T;R \subseteq S;R \\
\Leftarrow & \{ ; \text{ is monotonic } \} \\
& R^T \subseteq S \\
\Leftrightarrow & \{ \text{transposition} \} \\
& R \subseteq S^T \\
\Leftrightarrow & \{ I_V \text{ is identity of } ; \} \\
& R \subseteq S^T;I_V \\
\Leftrightarrow & \{ S;R = I_V \}
\end{aligned}$$

$$\begin{aligned}
& R \subseteq S^T ; S ; R \\
\Leftrightarrow & \quad \{ I_B \text{ is identity of } ; \} \\
& I_B ; R \subseteq S^T ; S ; R \\
\Leftarrow & \quad \{ ; \text{ is monotonic } \} \\
& I_B \subseteq S^T ; S \\
\Leftrightarrow & \quad \{ S \text{ is R-total: (12), with } B, V, R := V, B, S \} \\
& \text{true .}
\end{aligned}$$

From a heuristic point of view, the step labelled “transposition” is the most difficult one in this proof. It is needed to prepare for what follows, were we need to obtain I_B to the left of a set inclusion. Honesty forces me to admit here that I could only do this because I already knew that in the element-wise proof (in the previous section) I have used that S is R-total and that, therefore, it is most likely needed here as well.

□

It is somewhat amazing that, once we have formulated definitions (10) through (15) we have used no other properties of relation composition than that it has an identity element, that it is associative – which we have used implicitly –, and that it is monotonic (with respect to \subseteq).

3.9.5 Afterthoughts

More generally, the notion *exactly one* is an awkward one, mathematically speaking. The best way to treat it is to view it as the conjunction of the two notions *at least one* and *at most one*, which can be treated in relative isolation. In my experience, the notion *at least one* – non-emptiness – really differs from the notion *at most one* – uniqueness –.

* * *

The purely relational characterisations – formulae (10) through (15) – in terms of relation composition really are useful. They show that totality and functionality actually are faces of the same coin, and they allow for much more concise proofs. Relational algebra, with composition as an important operator, still is undervalued!

The transition even renders Lemma 0 completely trivial. Lemma 1, on the other hand, is not trivial at all; in this respect it is not strange that students find this a difficult exercise.

3.10 Exercises

1. Set A is given by $A = \{1, 2, 3, 4\}$. Which of the following relations are functions from A to A ?

(a) $\{(1, 3), (2, 4), (3, 1), (4, 2)\}$;

- (b) $\{(1, 3)(2, 4)\}$;
 (c) $\{(1, 1), (2, 2), (3, 3), (4, 4), (1, 3), (2, 4), (3, 1), (4, 2)\}$;
 (d) $\{(1, 1), (2, 2), (3, 3), (4, 4)\}$.
2. Suppose f and g are functions from \mathbb{R} to \mathbb{R} defined by $f(x) = x^2$ and $g(x) = x+1$ for all $x \in \mathbb{R}$. What is $g \circ f$ and what is $f \circ g$?
3. Which of the following functions is injective, surjective and/or bijective?
- (a) $f : \mathbb{R} \rightarrow \mathbb{R}, f(x) = x^2$ for all $x \in \mathbb{R}$.
 (b) $f : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}, f(x) = x^2$ for all $x \in \mathbb{R}$.
 (c) $f : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}, f(x) = x^2$ for all $x \in \mathbb{R}$.
4. Suppose R and S are relations on a set V with $R;S = I$ and $S;R = I$. Prove that both R and S are bijective functions.
5. Let R be a finite relation with adjacency matrix A . Prove the following statements:
- (a) If every row of A contains one non-zero entry, then R is a function.
 (b) If, in addition, every column of A contains at most one entry, then function R is injective.
 (c) If every row and column of A contain exactly one 1, then R is a bijection. What is the adjacency matrix of the inverse function?
6. Let B and V be sets and let R be a relation from B to V . Then, for every $v \in V$, we have defined⁹ the *pre-image* of v as the set ${}_R[v]$, thus:
- $${}_R[v] = \{b \in B \mid bRv\} ,$$
- which, obviously, is a subset of B .
- (a) Prove that the relation $\{(v, {}_R[v]) \mid v \in V\}$ is a function from V to $\mathcal{P}(B)$ (the set of all subsets of B).
 (b) Prove that, if F is a function in $V \rightarrow \mathcal{P}(B)$, then the set R_F defined by $R_F = \{(b, v) \mid b \in F(v)\}$ is a relation from B to V , with ${}_{R_F}[v] = F(v)$, for all $v \in V$.
7. Prove Lemma 3.5.
 8. Prove Properties 3.6.
 9. Prove Properties 3.7.
 10. Prove Properties 3.8.

⁹See the chapter on relations.

11. Prove Theorem 3.10.
12. Prove Lemma 3.13.
13. Prove Lemma 3.16.
14. Prove Lemma 3.18.
15. Given are two bijective functions f , of type $U \rightarrow V$, and g , of type $V \rightarrow W$. Prove that: $(g \circ f)^{-1} = f^{-1} \circ g^{-1}$.
16. We consider functions f, g, h , of type $U \rightarrow U$, such that both $g \circ f$ and $h \circ g$ are bijective. Prove that $h \circ g \circ f$ is bijective as well.
17. Prove that an injective function f in $B \rightarrow V$ has a *partial inverse*, that is, a partial function g from V to B such that $g \circ f = I_B$. What is the domain of g ?
18. Prove Lemma 3.20.
19. Prove Properties 3.22.

4 Posets and lattices

4.1 Partial orders

4.1 Definition. An (endo)relation \sqsubseteq (“under”) on a set P is called a *partial order* if it is reflexive, antisymmetric, and transitive. We recall that this means that, for all $x, y, z \in P$, we have:

- $x \sqsubseteq x$;
- $x \sqsubseteq y \wedge y \sqsubseteq x \Rightarrow x = y$;
- $x \sqsubseteq y \wedge y \sqsubseteq z \Rightarrow x \sqsubseteq z$.

The pair (P, \sqsubseteq) is called a *partially ordered set* or, for short, a *poset*.

Two elements x and y in a poset (P, \sqsubseteq) are called *comparable* if $x \sqsubseteq y$ or $y \sqsubseteq x$, otherwise they are called *incomparable*, that is, if $\neg(x \sqsubseteq y)$ and $\neg(y \sqsubseteq x)$.

A partial order is a *total order*, also called *linear order*, if every two elements are comparable.

□

4.2 Lemma. For every poset (P, \sqsubseteq) and for every subset X of P , the pair (X, \sqsubseteq) is a poset too.

□

4.3 Example.

- On every set, the identity relation I is a partial order. It is the *smallest* possible partial order relation on that set.
- On the real numbers \mathbb{R} the relation \leq is a total order: every two numbers $x, y \in \mathbb{R}$ satisfy $x \leq y$ or $y \leq x$. Restriction of \leq to any subset of \mathbb{R} –for example, restriction to $\mathbb{Q}, \mathbb{Z}, \mathbb{N}$ – also yields a total order on that subset.
- The power set $\mathcal{P}(V)$ of a set V , that is, the set of all subsets of V , with relation \subseteq (subset inclusion), is a poset. This P contains a smallest element, namely \emptyset , and a largest element, namely V itself.
- The relation $|$ (“divides”) is a partial order on the positive naturals. This example can be considered an instance of the previous one, as follows. For each $n \in \mathbb{N}^+$ we define $D(n)$ as the set of all divisors of n . Then we have that

$$m | n \Leftrightarrow D(m) \subseteq D(n) \text{ , for all } m, n \in \mathbb{N}^+ \text{ .}$$

Thus, relation $|$ on \mathbb{N}^+ is “equivalent”¹⁰ to relation \subseteq on set $\{D(n) \mid n \in \mathbb{N}^+\}$, which is a subset of $\mathcal{P}(\mathbb{N}^+)$.

¹⁰Actually, the proper word to be used here is “isomorphic”, to be introduced later.

- Let P be the set of all partitions of a set V . We define the relation “refines” as follows. Partition Π_1 refines partition Π_2 if and only if each $X \in \Pi_1$ is contained in some $Y \in \Pi_2$, that is:

$$(\forall X : X \in \Pi_1 : (\exists Y : Y \in \Pi_2 : X \subseteq Y)) .$$

The relation “refines” is a partial order on P . For the corresponding equivalence relations R_{Π_1} and R_{Π_2} we have that Π_1 refines Π_2 if and only if $R_{\Pi_1} \subseteq R_{\Pi_2}$.

□

- 4.4 Definition.** The *irreflexive part* of a partial order relation \sqsubseteq is often denoted by \sqsubset (“strictly under”) and is defined by, for all $x, y \in P$:

$$x \sqsubset y \Leftrightarrow x \sqsubseteq y \wedge x \neq y .$$

□

4.5 Properties.

- Relation \sqsubset , as defined above, is irreflexive, antisymmetric, and transitive.
- Conversely, we have, for all $x, y \in P$: $x \sqsubseteq y \Leftrightarrow x \sqsubset y \vee x = y$.

□

- 4.6 Lemma.** If (P, \sqsubseteq) is a poset, then the corresponding directed graph, with vertex set P and arrows (x, y) whenever $x \sqsubset y$, is *acyclic*.

□

If we want to draw a picture of the poset, we usually do not draw the whole digraph. Instead we only draw an edge from x to y from P with $x \sqsubseteq y$ if there is no z , distinct from both x and y , for which we have $x \sqsubseteq z$ and $z \sqsubseteq y$. This digraph is called the *Hasse diagram* for (P, \sqsubseteq) , named after the German mathematician Helmut Hasse (1898-1979). Usually pictures of Hasse diagrams are drawn in such a way that two vertices x and y with $x \sqsubseteq y$ are connected by an edge going upwards. For example the Hasse diagram for the poset $(\mathcal{P}(\{a, b, c\}), \sqsubseteq)$ is drawn as in Figure 23.

* * *

There are various ways of constructing new posets out of old ones. We discuss some of them. In the sequel both (P, \sqsubseteq) and (Q, \sqsubseteq) are posets. Notice that we use the *same* symbol, \sqsubseteq , for the two *different* partial order relations on sets P and Q . If confusion may arise we distinguish the two relations by using \sqsubseteq_P and \sqsubseteq_Q , respectively.

- As already stated in Lemma 4.2, for every subset X of P the pair (X, \sqsubseteq) is a poset, with relation \sqsubseteq restricted to X . Thus restricted \sqsubseteq is called the *induced* order on X .

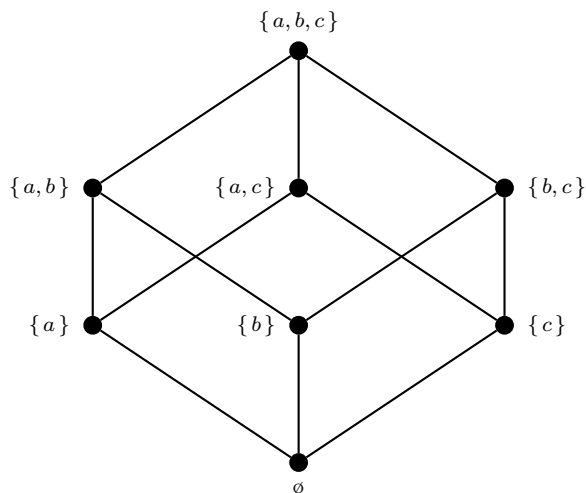


Figure 23: A Hasse diagram of $(\mathcal{P}(\{a, b, c\}), \subseteq)$

- Relation \supseteq (“above”), defined by, for all $x, y \in P$, $x \supseteq y \Leftrightarrow y \subseteq x$ is a partial order too, called the *dual order to* \subseteq ; hence, (P, \supseteq) also is a poset.
- Let V be a set. On the set $V \rightarrow P$ of functions from V to P we can define a partial order \subseteq_{VP} , say, as follows, for all $f, g \in V \rightarrow P$:

$$f \subseteq_{VP} g \Leftrightarrow (\forall v: v \in V: f(v) \subseteq_P g(v)) \ .$$

Then $(V \rightarrow P, \subseteq_{VP})$ is a poset.

- On the cartesian product $P \times Q$ we can define a partial order as follows. For $(p, q), (x, y) \in P \times Q$ we define:

$$(p, q) \subseteq (x, y) \Leftrightarrow p \subseteq_P x \wedge q \subseteq_Q y \ .$$

Thus defined relation \subseteq is a partial order, called the *product order* on $P \times Q$.

- On the cartesian product $P \times Q$ we also can define a partial order as follows. For $(p, q), (x, y) \in P \times Q$ we define:

$$(p, q) \subseteq (x, y) \Leftrightarrow (p \neq x \wedge p \subseteq_P x) \vee (p = x \wedge q \subseteq_Q y) \ .$$

This relation \subseteq is a partial order too, called the *lexicographic order* on $P \times Q$.

The notions of product order and lexicographic order can be extended to (finite) products of more than two sets.

4.2 Indirect Equality

A poset's structure is determined by its partial order relation, as denoted by \sqsubseteq in the previous section. Sometimes we wish to prove equalities in a poset, and then it can be convenient if we can reformulate an equality in terms of the partial order relation: that enables us to reason in terms of the partial order.

The following lemma provides such a connection between equality and the partial order, which turns out to be quite useful. In the course of this text we will refer to it as “Indirect Equality”.

4.7 Lemma. In every poset (P, \sqsubseteq) we have, for all $x, y \in P$:

$$(a) \quad x = y \Leftrightarrow (\forall z: z \in P: x \sqsubseteq z \Leftrightarrow y \sqsubseteq z) \quad ;$$

$$(b) \quad x = y \Leftrightarrow (\forall z: z \in P: z \sqsubseteq x \Leftrightarrow z \sqsubseteq y) \quad .$$

Proof. We prove (a) only; the proof of (b) follows, mutatis mutandis, the same pattern. We do so by mutual implication.

“ \Rightarrow ”: Assuming $x = y$ we may substitute x for y and vice versa wherever we like; in particular, if we substitute x for y in our demonstrandum $(\forall z: z \in P: x \sqsubseteq z \Leftrightarrow y \sqsubseteq z)$, we obtain $(\forall z: z \in P: y \sqsubseteq z \Leftrightarrow y \sqsubseteq z)$, which is true because of the reflexivity of \Leftrightarrow .

“ \Leftarrow ”: Assuming $(\forall z: z \in P: x \sqsubseteq z \Leftrightarrow y \sqsubseteq z)$, by the instantiation $z := x$ we obtain $x \sqsubseteq x \Leftrightarrow y \sqsubseteq x$, which is equivalent to $y \sqsubseteq x$, because \sqsubseteq is reflexive. Moreover, by the instantiation $z := y$ we obtain $x \sqsubseteq y \Leftrightarrow y \sqsubseteq y$, which is equivalent to $x \sqsubseteq y$. By the antisymmetry of \sqsubseteq we conclude $x = y$, as required.

□

4.8 Exercise. Prove the rules of “Indirect Inequality” in a poset (P, \sqsubseteq) , for all $x, y \in P$:

$$(a) \quad x \sqsubseteq y \Leftrightarrow (\forall z: z \in P: x \sqsubseteq z \Leftarrow y \sqsubseteq z) \quad ;$$

$$(b) \quad x \sqsubseteq y \Leftrightarrow (\forall z: z \in P: z \sqsubseteq x \Rightarrow z \sqsubseteq y) \quad .$$

□

4.3 Extreme elements

4.9 Definition. For any subset X of a poset (P, \sqsubseteq) we define, for all $m \in P$:

(a) m is X 's *maximum*, or *largest element* if:

$$m \in X \quad \wedge \quad (\forall x: x \in X: x \sqsubseteq m) \quad .$$

(b) m is X 's *minimum*, or *least element* if:

$$m \in X \quad \wedge \quad (\forall x: x \in X: m \sqsubseteq x) \quad .$$

(c) m is a *maximal* element of X if:

$$m \in X \quad \wedge \quad (\forall x: x \in X: \neg(m \sqsubset x)) \quad .$$

(d) m is a *minimal* element of X if:

$$m \in X \wedge (\forall x: x \in X: \neg(x \sqsubset m)) .$$

□

Notice the difference between the notions “maximum” and “maximal”. A value $m \in P$ is X ’s maximum if $m \in X$ and all elements in X are under m , whereas m is a maximal element of X if $m \in X$ and X contains no elements strictly above m .

A subset of a partially ordered set does not necessarily contain such extreme elements. If it exists, however, the maximum of a subset is unique, and so is the minimum, if it exists. If subset X has a maximum we denote it by $\max X$, and its minimum we denote by $\min X$.

4.10 Lemma. Let X be a subset of a poset (P, \sqsubseteq) . If m and n both are X ’s maximum then $m = n$. If m and n both are X ’s minimum then $m = n$.

□

A subset of a partially ordered set does not necessarily contain maximal or minimal elements. Such elements are not unique either: a subset may have many maximal or minimal elements. As a rather trivial example, consider poset (P, I_P) , where the identity relation I_P is the smallest possible partial order on P : $x I_P y \Leftrightarrow x = y$. With this particular order all elements of a subset $X \subseteq P$ both are maximal and minimal.

4.11 Definition. Let (P, \sqsubseteq) be a poset. If the whole set P has a minimum this is often denoted by \perp (“bottom”), and if P has a maximum this is often denoted by \top (“top”). If P has a minimum the minimal elements of $P \setminus \{\perp\}$ are called P ’s *atoms*.

4.12 Example.

- If we consider the poset of all subsets of a set V , then the empty set \emptyset is the minimum of the poset, whereas the whole set V is the maximum. The atoms are the subsets of V containing just a single element.
- In the poset $(\mathbb{N}^+, |)$ the whole set, \mathbb{N}^+ , has no maximum. Its minimum, however, equals 1. The atoms are the prime numbers.
- If (P, \sqsubseteq) is totally ordered then subset $\{x, y\}$ has a maximum and a minimum, for all $x, y \in P$.

□

4.13 Lemma. In a poset (P, \sqsubseteq) for every subset X its maximum, if it exists, is a maximal element of X ; also, its minimum, if it exists, is a minimal element of X .

4.14 Lemma. If a poset (P, \sqsubseteq) is a total order then every subset $X \subseteq P$ has at most one maximal element, which then also is its maximum. Also, X has at most one minimal element, which then also is its minimum.

4.15 Lemma. Let (P, \sqsubseteq) be a *nonempty* and *finite* poset. Then P contains a maximal element and a minimal element.

Proof. This is equivalent to the proposition that every finite and acyclic graph contains both a vertex without incoming arrows – a minimal element – and a vertex without outgoing arrows – a maximal element –.

□

4.16 Algorithm. [Topological sorting] For any finite poset (P, \sqsubseteq) the elements of P can be ordered totally in such a way that the partial order is “respected”. That is, on P a total order \preceq , say, exists satisfying $x \sqsubseteq y \Rightarrow x \preceq y$, for all $x, y \in P$.

For finite set P of size N , so $N = \#P$, such a total order can be conveniently represented by a finite sequence $[s_0, \dots, s_{N-1}]$, containing all elements of P . Formally, this means that s is a surjective function of type $[0..N) \rightarrow P$. The total order of P ’s elements then is given by the order in which they occur in s . This order “respects” partial order \sqsubseteq if and only if:

$$(16) \quad s_i \sqsubseteq s_j \Rightarrow i \leq j, \text{ for all } i, j: 0 \leq i, j < N.$$

The process of constructing such a sequence s is called “topological sorting”. This process can be formulated recursively, in a simple way. We define a sequence $[Q_0, \dots, Q_N]$ of subsets of P , as follows. These subsets will be such that $\#Q_n = n$, for all n , $0 \leq n \leq N$; so, $Q_0 = \emptyset$ and all others will be nonempty. We define $Q_N = P$, and for every n , $0 \leq n < N$, we observe that, because (Q_{n+1}, \sqsubseteq) is a nonempty and finite poset too, by Lemma 4.15, set Q_{n+1} contains at least one maximal element s_n , say. Now we define $Q_n = Q_{n+1} \setminus \{s_n\}$. Thus defined, indeed we have $\#Q_n = n$, for all n , so Q_{n+1} is nonempty, as required.

In addition we now also have defined a function s in $[0..N) \rightarrow P$, and by construction it satisfies $Q_n = \{s_j \mid 0 \leq j < n\}$, for all n . Function s is surjective, hence, because $N = \#P$, s also is injective. By definition, for any n , $0 \leq n < N$, we have that s_n is a maximal element of Q_{n+1} . By Definition 4.9 this means that:

$$\begin{aligned} & (\forall x: x \in Q_{n+1} : \neg(s_n \sqsubset x)) \\ \Rightarrow & \quad \{ Q_n \subseteq Q_{n+1} \} \\ & (\forall x: x \in Q_n : \neg(s_n \sqsubset x)) \\ \Leftrightarrow & \quad \{ \text{dummy transformation } x := s_j \} \\ & (\forall j: 0 \leq j < n : \neg(s_n \sqsubset s_j)) \\ \Leftrightarrow & \quad \{ \text{range-term trading, and contraposition} \} \\ & (\forall j: 0 \leq j < N : s_n \sqsubset s_j \Rightarrow \neg(j < n)) \\ \Leftrightarrow & \quad \{ \text{property of } < \text{ and } \leq \} \\ & (\forall j: 0 \leq j < N : s_n \sqsubset s_j \Rightarrow n \leq j) \\ \Leftrightarrow & \quad \{ s \text{ is injective, so } s_n = s_j \Rightarrow n \leq j \} \\ & (\forall j: 0 \leq j < N : s_n \sqsubseteq s_j \Rightarrow n \leq j) \end{aligned}$$

from which we conclude that s satisfies property (16), as required.

□

4.17 Example. Topological sorting has various applications. For example consider a (so-called) *spreadsheet*. In a spreadsheet the values in various *cells* depend on each other, but, in a correct spreadsheet, in an acyclic way only. The value in any particular cell in the spreadsheet can only be computed if the values in all cells on which this particular cell depends have been computed already. Therefore, an efficient implementation of these computations requires that they are performed in the “right” order. This gives rise to a partial order on the set of cells within a spreadsheet. By topological sorting the set of cells can be linearized in such a way that every cell precedes all cells depending on it; thus the computations of the values in the cells can be performed in a linear order.

4.4 Upper and lower bounds

4.18 Definition. For any subset X of a poset (P, \sqsubseteq) we define, for all $m \in P$:

- (a) m is an *upper bound* of X if: $(\forall x: x \in X: x \sqsubseteq m)$
- (b) m is a *lower bound* of X if: $(\forall x: x \in X: m \sqsubseteq x)$

□

4.19 Properties.

- (a) If P has a maximum \top then \top is an upper bound of every subset of P .
- (b) If P has a minimum \perp then \perp is a lower bound of every subset of P .
- (c) Every element in P is an upper bound and a lower bound of \emptyset .
- (d) If it exists $\max X$ is an upper bound of X , for all $X \subseteq P$.
- (e) If it exists $\min X$ is a lower bound of X , for all $X \subseteq P$.

□

For any subset $X \subseteq P$ we can consider the set $\{m \in P \mid (\forall x: x \in X: x \sqsubseteq m)\}$ of *all* upper bounds of X . This set may or may not have a minimum. If it has a minimum, this minimum is called the *supremum* of X , notation $\sup X$. Alternatively, it is sometimes also called X 's *least upper bound*, notation $\text{lub } X$.

Similarly, we can consider the set $\{m \in P \mid (\forall x: x \in X: m \sqsubseteq x)\}$ of *all* lower bounds of X . This set may or may not have a maximum. If it has a maximum, this maximum is called the *infimum* of X , notation $\inf X$. Alternatively, it is sometimes also called X 's *greatest lower bound*, notation $\text{glb } X$.

By combination of the definitions of maximum/minimum and of upper/lower bounds we can define supremum and infimum in a more direct way.

4.20 Definition. For any subset X of a poset (P, \sqsubseteq) we define, for all $m \in P$:

(a) m is X 's *supremum* if both:

$$(\forall x : x \in X : x \sqsubseteq m) \quad , \text{ and:}$$

$$(\forall x : x \in X : x \sqsubseteq z) \Rightarrow m \sqsubseteq z \quad , \text{ for all } z \in P \quad .$$

Notice that the first requirement expresses that m is an upper bound of X , and that the second one expresses that m is under all upper bounds of X .

(b) m is X 's *infimum* if both:

$$(\forall x : x \in X : m \sqsubseteq x) \quad , \text{ and:}$$

$$(\forall x : x \in X : z \sqsubseteq x) \Rightarrow z \sqsubseteq m \quad , \text{ for all } z \in P \quad .$$

□

4.21 Example.

- For a set V its power set $\mathcal{P}(V)$ – the set of all subsets of V – with relation \sqsubseteq is a poset, and any subset X of $\mathcal{P}(V)$ has a supremum, namely $(\bigcup_{U:U \in X} U)$, and an infimum, namely $(\bigcap_{U:U \in X} U)$.
- The set \mathbb{N}^+ of positive natural numbers with relation $|$ (“divides”) is a poset. The supremum of two elements $a, b \in \mathbb{N}^+$ is their *least common multiple*, that is, the *smallest* of all positive naturals m satisfying $a|m$ and $b|m$; usually, this value is denoted by $lcm(a, b)$.

Similarly, the greatest common divisor of a and b , denoted by $gcd(a, b)$, is the infimum of $\{a, b\}$.

□

4.22 Lemma. For poset (P, \sqsubseteq) and for $p \in P$ we have: $\sup\{p\} = p$ and $\inf\{p\} = p$.

Proof. By Definition 4.20, to prove $\sup\{p\} = p$ we must prove:

$$\begin{aligned} & (\forall x : x \in \{p\} : x \sqsubseteq p) \\ \Leftrightarrow & \quad \{ \text{definition of } \{p\} \} \\ & (\forall x : x = p : x \sqsubseteq p) \\ \Leftrightarrow & \quad \{ \text{1-pt. rule} \} \\ & p \sqsubseteq p \\ \Leftrightarrow & \quad \{ \sqsubseteq \text{ is reflexive} \} \\ & \text{true} \quad , \end{aligned}$$

and, for all $z \in P$:

$$\begin{aligned}
& (\forall x : x \in \{p\} : x \sqsubseteq z) \\
\Leftrightarrow & \quad \{ \text{same steps as above} \} \\
& p \sqsubseteq z \quad ,
\end{aligned}$$

which is the desired result.

□

4.23 Lemma. In a poset (P, \sqsubseteq) any subset $X \subseteq P$ for which $\sup X$ exists satisfies, for all $m \in P$:

$$m = \max X \Leftrightarrow m = \sup X \wedge m \in X \quad ,$$

similarly, if $\inf X$ exists then, for all $m \in P$:

$$m = \min X \Leftrightarrow m = \inf X \wedge m \in X \quad .$$

Proof. By direct application of the definitions of \max and \sup , and of \min and \inf respectively.

□

Corollary: If poset (P, \sqsubseteq) is such that $\sup P$ exists then $\sup P = \max P$, and if $\inf P$ exists then $\inf P = \min P$.

□

An important difference between supremum and maximum of a set is that a set's supremum may or may not be an element of that set, whereas a set's maximum always is an element of that set. The above lemma states, however, that if a set's supremum is in that set, then this supremum also is the set's maximum.

The following lemma provides a different but equivalent characterization of supremum and infimum that occasionally turns out to be very useful.

4.24 Lemma. Let X be a subset of a poset (P, \sqsubseteq) . For $m \in P$ we have:

(a) m is X 's *supremum* if and only if:

$$(\forall z : z \in P : m \sqsubseteq z \Leftrightarrow (\forall x : x \in X : x \sqsubseteq z)) \quad .$$

(b) m is X 's *infimum* if and only if:

$$(\forall z : z \in P : z \sqsubseteq m \Leftrightarrow (\forall x : x \in X : z \sqsubseteq x)) \quad .$$

Proof. We prove (a) only; the proof for (b) follows, mutatis mutandis, the same pattern. Firstly, we let m be X 's supremum according to Definition 4.20. Then, for $z \in P$ we prove the equivalence of $m \sqsubseteq z$ and $(\forall x : x \in X : x \sqsubseteq z)$, by "cyclic implication":

$$\begin{aligned}
& m \sqsubseteq z \\
\Leftrightarrow & \quad \{ \text{Definition 4.20: } m \text{ is an upper bound of } X \} \\
& (\forall x : x \in X : x \sqsubseteq m) \wedge m \sqsubseteq z \\
\Rightarrow & \quad \{ \forall \text{ introduction } \} \\
& (\forall x : x \in X : x \sqsubseteq m) \wedge (\forall x : x \in X : m \sqsubseteq z) \\
\Leftrightarrow & \quad \{ \text{combining terms } \} \\
& (\forall x : x \in X : x \sqsubseteq m \wedge m \sqsubseteq z) \\
\Rightarrow & \quad \{ \sqsubseteq \text{ is transitive } \} \\
& (\forall x : x \in X : x \sqsubseteq z) \\
\Rightarrow & \quad \{ \text{Definition 4.20: } m \text{ is under all upper bounds } \} \\
& m \sqsubseteq z \quad .
\end{aligned}$$

Secondly, let m satisfy:

$$(17) \quad (\forall z : z \in P : m \sqsubseteq z \Leftrightarrow (\forall x : x \in X : x \sqsubseteq z)) \quad .$$

Then we must prove that m is X 's supremum. Well, m is an upperbound:

$$\begin{aligned}
& (\forall x : x \in X : x \sqsubseteq m) \\
\Leftrightarrow & \quad \{ (17) , \text{ with } z := m \} \\
& m \sqsubseteq m \\
\Leftrightarrow & \quad \{ \sqsubseteq \text{ is reflexive } \} \\
& \text{true} \quad ,
\end{aligned}$$

and that m is under all upper bounds follows directly from (17), because \Leftrightarrow is stronger than \Rightarrow .

□

4.25 Properties.

- (a) If P has a maximum \top then $\top = \sup P$ and $\top = \inf \emptyset$.
- (b) If P has a minimum \perp then $\perp = \inf P$ and $\perp = \sup \emptyset$.

□

4.5 Lattices

4.5.1 Definition

In the previous section we have introduced the notions of *supremum* –least upper bound– and *infimum* –greatest lower bound– of subsets of a poset. Generally, such a subset does not have a supremum or an infimum, just as, generally, not every

subset has a maximum or a minimum. (Recall Lemma 4.23, for the relation between supremum and maximum, and between infimum and minimum, respectively.)

Partially ordered sets in which particular subsets do have suprema and/or infima are of interest. We will study three types of such posets, called “lattices”, “complete lattices”, and “complete partial orders”.

4.26 Definition. A poset (P, \sqsubseteq) is a *lattice*, if for all $x, y \in P$ the subset $\{x, y\}$ has a supremum and an infimum. Because this pertains to two-element sets, it is customary to use infix-notation to denote their suprema and infima. For this purposes binary operators \sqcup (“cup”) and \sqcap (“cap”) are used: the supremum of $\{x, y\}$ then is written as $x \sqcup y$ and its infimum as $x \sqcap y$.

□

4.27 Example. Here are some examples of lattices we already encountered before.

- (\mathbb{R}, \leq) is a lattice. For $x, y \in \mathbb{R}$ we have: $x \sqcup y = x \max y$ and $x \sqcap y = x \min y$.
- For a set V the poset $(\mathcal{P}(V), \subseteq)$ is a lattice, with $\sqcup = \cup$ and $\sqcap = \cap$.
- The poset $(\mathbb{N}^+, |)$ is a lattice, with $a \sqcup b = lcm(a, b)$ and $a \sqcap b = gcd(a, b)$.

□

The following lemma actually is a special case of Lemma 4.24, namely for non-empty subsets with *at most two* elements; that is, this is Lemma 4.24 with $X := \{x, y\}$.

4.28 Lemma. In every lattice (P, \sqsubseteq) we have, for all $x, y, z \in P$:

- (a) $x \sqcup y \sqsubseteq z \Leftrightarrow x \sqsubseteq z \wedge y \sqsubseteq z$;
- (b) $z \sqsubseteq x \sqcap y \Leftrightarrow z \sqsubseteq x \wedge z \sqsubseteq y$.

□

Actually, this lemma can serve as an *alternative definition* of \sqcup and \sqcap , as the original definition follows from it. As a special case, we obtain the following lemma, expressing that $x \sqcup y$ is an upper bound and that $x \sqcap y$ is a lower bound.

4.29 Lemma. In every lattice (P, \sqsubseteq) we have, for all $x, y \in P$:

- (a) $x \sqsubseteq x \sqcup y \wedge y \sqsubseteq x \sqcup y$;
- (b) $x \sqcap y \sqsubseteq x \wedge x \sqcap y \sqsubseteq y$.

Proof. Instantiate Lemma 4.28 with $z := x \sqcup y$ and $z := x \sqcap y$, respectively.

□

4.30 Lemma. In every lattice (P, \sqsubseteq) we have, for all $x, y \in P$:

- (a) $x \sqsubseteq y \Leftrightarrow x \sqcup y = y$;
 (b) $x \sqsubseteq y \Leftrightarrow x \sqcap y = x$.

Proof. We prove (a) only, by calculation:

$$\begin{aligned}
 & x \sqcup y = y \\
 \Leftrightarrow & \quad \{ \sqsubseteq \text{ is reflexive and antisymmetric} \} \\
 & x \sqcup y \sqsubseteq y \wedge y \sqsubseteq x \sqcup y \\
 \Leftrightarrow & \quad \{ x \sqcup y \text{ is an upperbound of } y \} \\
 & x \sqcup y \sqsubseteq y \\
 \Leftrightarrow & \quad \{ \text{Lemma 4.28, with } z := y \} \\
 & x \sqsubseteq y \wedge y \sqsubseteq y \\
 \Leftrightarrow & \quad \{ \sqsubseteq \text{ is reflexive} \} \\
 & x \sqsubseteq y
 \end{aligned}$$

□

4.5.2 Algebraic properties

The lattice operators have interesting algebraic properties, as reflected by the following theorem.

4.31 Theorem. Let (P, \sqsubseteq) be a lattice. Then for all $x, y, z \in P$ we have:

- (a) $x \sqcup x = x$ and $x \sqcap x = x$: \sqcup and \sqcap are *idempotent*;
 (b) $x \sqcup y = y \sqcup x$ and $x \sqcap y = y \sqcap x$: \sqcup and \sqcap are *commutative*;
 (c) $x \sqcup (y \sqcup z) = (x \sqcup y) \sqcup z$ and $x \sqcap (y \sqcap z) = (x \sqcap y) \sqcap z$: \sqcup and \sqcap are *associative*;
 (d) $x \sqcup (x \sqcap y) = x$ and $x \sqcap (x \sqcup y) = x$: *absorption*.

Proof.

- (a) See Lemma 4.22.
 (b) By symmetry: set $\{x, y\}$ equals set $\{y, x\}$.
 (c) By Indirect Equality (Lemma 4.7), for all $w \in P$:

$$\begin{aligned}
 & x \sqcup (y \sqcup z) \sqsubseteq w \\
 \Leftrightarrow & \quad \{ \text{Lemma 4.28} \} \\
 & x \sqsubseteq w \wedge y \sqcup z \sqsubseteq w \\
 \Leftrightarrow & \quad \{ \text{Lemma 4.28} \}
 \end{aligned}$$

$$\begin{aligned}
& x \sqsubseteq w \wedge y \sqsubseteq w \wedge z \sqsubseteq w \\
\Leftrightarrow & \quad \{ \text{Lemma 4.28} \} \\
& x \sqcup y \sqsubseteq w \wedge z \sqsubseteq w \\
\Leftrightarrow & \quad \{ \text{Lemma 4.28} \} \\
& (x \sqcup y) \sqcup z \sqsubseteq w .
\end{aligned}$$

(d) We prove $x \sqcup (x \sqcap y) = x$ only, by calculation:

$$\begin{aligned}
& x \sqcup (x \sqcap y) = x \\
\Leftrightarrow & \quad \{ \text{Lemma 4.30} \} \\
& x \sqcup (x \sqcap y) \sqsubseteq x \\
\Leftrightarrow & \quad \{ \text{Lemma 4.28} \} \\
& x \sqsubseteq x \wedge x \sqcap y \sqsubseteq x \\
\Leftrightarrow & \quad \{ \sqsubseteq \text{ is reflexive, and } x \sqcap y \text{ is a lower bound of } x \} \\
& \text{true} .
\end{aligned}$$

□

Conversely, the following theorem expresses that every structure with operators having the above algebraic properties “is” – can be extended into – a lattice.

4.32 Theorem. Let P be a set with two binary operators \sqcup and \sqcap ; that is, these operators have type $P \times P \rightarrow P$. Let these operators have algebraic properties (a) through (d), as in the previous theorem. Then the relation \sqsubseteq , on P and defined by $x \sqsubseteq y \Leftrightarrow x \sqcup y = y$ for all $x, y \in P$, is a partial order, and (P, \sqsubseteq) is a lattice.

□

A direct consequence of Theorem 4.31, particularly of the associativity of the operators, is that every *finite* and *non-empty* subset of a lattice has a supremum and an infimum. Notice that the requirement “non-empty” is essential here: in a lattice the empty set may not have a supremum or infimum.

4.33 Theorem. Let (P, \sqsubseteq) be a lattice. Then every finite and non-empty subset of P has a supremum and an infimum.

Proof. By Mathematical Induction on the size of the subsets, and using Theorem 4.31.

□

4.5.3 Distributive lattices

The prototype example of a lattice is the poset of all subsets of a set V , with \subseteq as the partial order relation. As already mentioned, in this lattice set union, \cup , and intersection, \cap , are the binary lattice operators. In this particular lattice, the operators have an additional algebraic property, namely (mutual) *distributivity*; that is, “ \cup distributes over \cap ” and “ \cap distributes over \cup ”, respectively:

$$X \cup (Y \cap Z) = (X \cup Y) \cap (X \cup Z) \quad , \text{ for all } X, Y, Z \subseteq V \quad ;$$

$$X \cap (Y \cup Z) = (X \cap Y) \cup (X \cap Z) \quad , \text{ for all } X, Y, Z \subseteq V \quad .$$

Generally, lattices do *not* have these properties. They do satisfy, however, a weaker version of these properties, namely “in one direction” only.

4.34 Theorem. Let (P, \subseteq) be a lattice. Then for all $x, y, z \in P$ we have:

- (a) $x \sqcup (y \sqcap z) \subseteq (x \sqcup y) \sqcap (x \sqcup z)$;
- (b) $(x \sqcap y) \sqcup (x \sqcap z) \subseteq x \sqcap (y \sqcup z)$.

Proof.

- (a) By calculation:

$$\begin{aligned}
 & x \sqcup (y \sqcap z) \subseteq (x \sqcup y) \sqcap (x \sqcup z) \\
 \Leftrightarrow & \quad \{ \text{Lemma 4.28(b)} \} \\
 & x \sqcup (y \sqcap z) \subseteq x \sqcup y \quad \wedge \quad x \sqcup (y \sqcap z) \subseteq x \sqcup z \\
 \Leftrightarrow & \quad \{ \text{Lemma 4.28(a) (twice)} \} \\
 & x \subseteq x \sqcup y \quad \wedge \quad y \sqcap z \subseteq x \sqcup y \quad \wedge \quad x \subseteq x \sqcup z \quad \wedge \quad y \sqcap z \subseteq x \sqcup z \\
 \Leftrightarrow & \quad \{ \text{Lemma 4.29(a) (twice)} \} \\
 & y \sqcap z \subseteq x \sqcup y \quad \wedge \quad y \sqcap z \subseteq x \sqcup z \\
 \Leftarrow & \quad \{ \subseteq \text{ is transitive (twice)} \} \\
 & y \sqcap z \subseteq y \quad \wedge \quad y \subseteq x \sqcup y \quad \wedge \quad y \sqcap z \subseteq z \quad \wedge \quad z \subseteq x \sqcup z \\
 \Leftrightarrow & \quad \{ \text{Lemma 4.29(a) (twice) and Lemma 4.29(b) (twice)} \} \\
 & \text{true} \quad .
 \end{aligned}$$

- (b) By duality.

□

As we have seen, some lattices –like $(\mathcal{P}(V), \subseteq)$ – do satisfy the distribution properties. Such lattices are called “distributive”.

4.35 Definition. A *distributive lattice* is a lattice (P, \sqsubseteq) in which \sqcup and \sqcap distribute over each other; that is, for all $x, y, z \in P$:

$$(a) \quad x \sqcup (y \sqcap z) = (x \sqcup y) \sqcap (x \sqcup z) ;$$

$$(b) \quad x \sqcap (y \sqcup z) = (x \sqcap y) \sqcup (x \sqcap z) .$$

□

4.36 Example. Not every lattice is distributive. The smallest example illustrating this has 5 elements \top, a, b, c, \perp , say, with this partial order: \perp is under all elements, a, b, c are under \top and are mutually incomparable. (See the Hasse diagram in Figure 21.) In this lattice $a \sqcup (b \sqcap c) = a$ whereas $(a \sqcup b) \sqcap (a \sqcup c) = \top$.

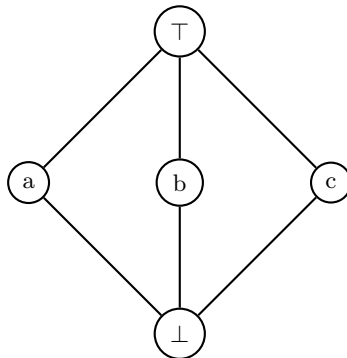


Figure 24: The smallest non-distributive lattice

4.5.4 Complete lattices

4.37 Definition. A *complete lattice* is a partial order in which *every* subset has a supremum and an infimum. In particular, the whole set has a supremum and an infimum, usually denoted by \top and \perp , respectively. Notice that \top and \perp also are the lattice's maximum and minimum. (Recall Property 4.25.)

□

4.38 Lemma. Every *finite* lattice is complete.

Proof. Let (P, \sqsubseteq) be a finite lattice. By Theorem 4.33 this lattice is *almost* complete already: every non-empty subset has a supremum and an infimum, so all we must do is prove that \emptyset has a supremum and an infimum. From Property 4.25 we know that $\sup \emptyset = \perp$ and $\inf \emptyset = \top$.

□

4.39 Example. The power set $(\mathcal{P}(V), \subseteq)$ of all subsets of a set V is a complete lattice. The supremum of a set X of subsets of V is the union of all sets in X , which is $(\bigcup_{U:U \in X} U)$; its infimum is the intersection of all its elements, which is $(\bigcap_{U:U \in X} U)$.

□

4.40 Example. The poset (\mathbb{R}, \leq) is a lattice, but it is not complete, but every *closed* interval $[a, b]$, with $a \leq b$ and with the same partial order \leq , is a complete (sub)lattice.

□

Completeness is a strong property, so strong, actually, that we only have to prove half of it: if every subset has an infimum that it also has a supremum, so if every subset has an infimum the partial order is a complete lattice.

4.41 Theorem. Let (P, \subseteq) be a poset. Then:

(a) “Every subset of P has an infimum” \Rightarrow “ (P, \subseteq) is a complete lattice” ;

(b) “Every subset of P has a supremum” \Rightarrow “ (P, \subseteq) is a complete lattice” .

Proof. We prove (a) only. To prove this we assume that every subset of P has an infimum. Then, to prove that (P, \subseteq) is complete we only must prove that every subset of P has a supremum. So, let X be a subset of P . We define a subset Y by $Y = \{ y \in P \mid (\forall x : x \in X : x \subseteq y) \}$, that is, Y is the set of all upper bounds of X . Now let $m = \inf Y$. We prove that this m is X 's supremum.

“ m is an upper bound of X ”:

$$\begin{aligned}
 & (\forall x : x \in X : x \subseteq m) \\
 \Leftrightarrow & \quad \{ m \text{ is } Y\text{'s } \textit{greatest} \text{ lower bound} \} \\
 & (\forall x : x \in X : (\forall y : y \in Y : x \subseteq y)) \\
 \Leftrightarrow & \quad \{ \text{exchanging dummies} \} \\
 & (\forall y : y \in Y : (\forall x : x \in X : x \subseteq y)) \\
 \Leftrightarrow & \quad \{ \text{definition of } Y \} \\
 & (\forall y : y \in Y : y \in Y) \\
 \Leftrightarrow & \quad \{ \text{predicate calculus} \} \\
 & \text{true} .
 \end{aligned}$$

“ m is under all upper bounds of X ”: For any $y \in P$ we derive:

$$\begin{aligned}
 & (\forall x : x \in X : x \subseteq y) \\
 \Leftrightarrow & \quad \{ \text{definition of } Y \} \\
 & y \in Y \\
 \Rightarrow & \quad \{ m \text{ is a lower bound of } Y \} \\
 & m \subseteq y .
 \end{aligned}$$

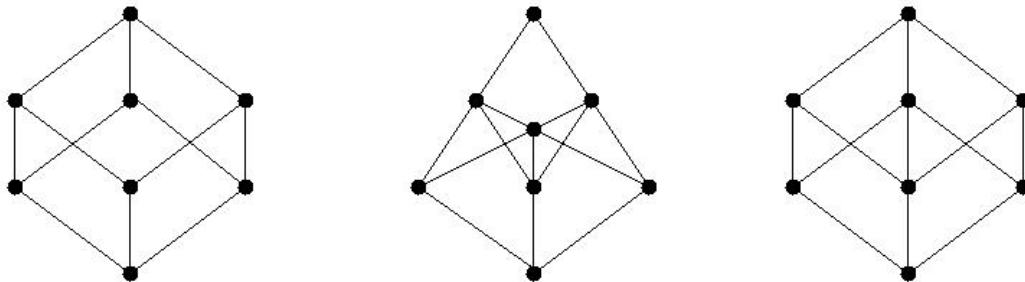
□

remark: In the proof of this theorem we introduced set Y as the set of *all* upper bounds of X . It is possible that X has no upper bounds, in which case Y is *empty*. This is harmless, because if *every* subset has an infimum then so has the empty set. Thus, this proof crucially depends on the property that also the empty set has an infimum. Recall that in a non-complete lattice the empty set does not need to have an infimum or supremum. (Also see the discussion preceding Theorem 4.33.)

□

4.6 Exercises

1. Let (P, \sqsubseteq) be a poset and X a subset of P . Prove that an element $m \in X$ is maximal if and only if for all $x \in X$ we have $m \sqsubseteq x \Rightarrow m = x$.
2. Let (P, \sqsubseteq) be a lattice and $x, y, z \in P$. Prove that:
 - (a) $x \sqsubseteq y \Rightarrow x \sqcup z \sqsubseteq y \sqcup z$;
 - (b) $x \sqsubseteq z \Rightarrow z \sqcup (x \sqcap y) = z$.
3. Prove Lemma 4.13
4. Prove Lemma 4.14
5. (a) Prove that the set $\{x \in \mathbb{Q} \mid x < 1\}$ has no maximum.
 (b) Prove that this set has a supremum.
6. We consider a poset (P, \sqsubseteq) ; let X and Y be subsets of P such that $\sup X$ and $\sup Y$ both exist. Prove that if $\sup X \in Y$ then $\sup Y$ is an upperbound of X .
7. We consider a poset (P, \sqsubseteq) ; let X and Y be subsets of P such that $\sup X$ and $\inf Y$ both exist. In addition it is given that $x \sqsubseteq y$, for all $x \in X$ and $y \in Y$. Prove that $\sup X \sqsubseteq \inf Y$.
8. We consider a poset (P, \sqsubseteq) ; let X be a subset of P .
 - (a) Prove that if X contains two (different) maximal elements then X has no maximum.
 - (b) Prove that if $(\forall x, y: x, y \in X: x \sqsubseteq y \vee y \sqsubseteq x)$ and if X contains a maximal element then X has a maximum.
9. In the figure below you see three diagrams. Which of these diagrams are Hasse diagrams? Which of these diagrams represents a lattice?
10. Show that every lattice with *at most* 4 elements is distributive.



11. Is the poset $(\mathbb{N}^+, |)$ a complete lattice? How about $(\mathbb{N}, |)$?
12. Suppose (P, \sqsubseteq) is a lattice and $a, b \in P$ with $a \sqsubseteq b$. Prove that $([a, b], \sqsubseteq)$ is a lattice too. Here $[a, b]$ denotes the interval from (and including) a upto (and including) b ; how would you define this interval?
13. Let (P, \sqsubseteq) be a lattice. Prove that if \sqcup distributes over \sqcap then \sqcap also distributes over \sqcup .
14. Prove that, in a complete lattice (P, \sqsubseteq) , the extreme element \top is the identity element of \sqcap and a zero element of \sqcup . Similarly, show that \perp is the identity of \sqcup and a zero of \sqcap .

4.7 Monotonic and continuous functions

In this subsection we study some classes of functions in the context of partial orders and lattices. Apart from having some general relevance, such functions play an important role in the fixed-point theorems in the following section.

4.42 Definition. Let (P, \sqsubseteq_P) and (Q, \sqsubseteq_Q) be posets. A function f from P to Q is *monotonic* if f preserves the partial orders, that is, if, for all $x, y \in P$:

$$x \sqsubseteq_P y \Rightarrow f(x) \sqsubseteq_Q f(y) .$$

As a special case, an (endo)function from poset (P, \sqsubseteq) to itself is monotonic if, for all $x, y \in P$:

$$x \sqsubseteq y \Rightarrow f(x) \sqsubseteq f(y) .$$

□

4.43 Example.

- (a) In a lattice (P, \sqsubseteq) the operator \sqcup is monotonic in both arguments; that is, for fixed $z \in P$ functions f and g , of type $P \rightarrow P$, defined by $f(x) = x \sqcup z$ and $g(y) = z \sqcup y$, both are monotonic. Similarly, \sqcap is monotonic too.

- (b) For sets B and V and for every function f from B to V the lifted function f , from $\mathcal{P}(B)$ to $\mathcal{P}(V)$, is monotonic with respect to the partial order \subseteq .
- (c) For endofunctions f on \mathbb{R} – or any subset thereof like $\mathbb{N}, \mathbb{Z}, \mathbb{Q}$ –, with partial order \leq , monotonicity means, for all $x, y \in \mathbb{R}$:

$$x \leq y \Rightarrow f(x) \leq f(y) .$$

Examples of monotonic functions on \mathbb{R} are $x \mapsto x+1$ and $x \mapsto e^x$. Other functions, like $x \mapsto x^2$ and $x \mapsto \sqrt{x}$, only are monotonic on a subset of \mathbb{R} , like $\mathbb{R}_{\geq 0}$.

□

4.44 Definition. In a partial order (P, \sqsubseteq) a function f is *continuous* if for every non-empty subset $X \subseteq P$ we have:

$$\text{“sup } X \text{ exists”} \Rightarrow \text{“sup } (f(X)) \text{ exists”} \wedge f(\text{sup } X) = \text{sup } (f(X)) .$$

□

4.45 Example. In a partial order (P, \sqsubseteq) , for every $b \in P$, the constant function with value b , that is, function f defined by $f(x) = b$ for all $x \in P$, is continuous. The identity function I_P is continuous too.

□

4.46 Lemma. In a lattice (P, \sqsubseteq) every continuous function f in $P \rightarrow P$ distributes over \sqcup , that is, for all $x, y \in P$:

$$f(x \sqcup y) = f(x) \sqcup f(y) .$$

As a result, in view of Theorem 4.33, a continuous function distributes over suprema of all finite and non-empty subsets of the lattice; that is, if f is continuous then $f(\text{sup } X) = \text{sup } (f(X))$, for all finite and non-empty $X \subseteq P$.

□

4.47 Lemma. In a lattice (P, \sqsubseteq) every continuous function is monotonic.

Proof. Let f in $P \rightarrow P$ be continuous; then, for $x, y \in P$ we calculate:

$$\begin{aligned} & f(x) \sqsubseteq f(y) \\ \Leftrightarrow & \quad \{ \text{Lemma 4.30} \} \\ & f(x) \sqcup f(y) = f(y) \\ \Leftrightarrow & \quad \{ f \text{ is continuous: Lemma 4.46} \} \\ & f(x \sqcup y) = f(y) \\ \Leftarrow & \quad \{ \text{Leibniz} \} \\ & x \sqcup y = y \\ \Leftrightarrow & \quad \{ \text{Lemma 4.30} \} \\ & x \sqsubseteq y , \end{aligned}$$

which proves that f is monotonic, as required.

□

* * *

Generally, continuous functions are *not* required to distribute over the supremum of the empty set. If function f would distribute over the supremum of the empty set, we would have:

$$\begin{aligned} f(\sup \emptyset) &= \sup(f(\emptyset)) \\ \Leftrightarrow \quad \{ f(\emptyset) = \emptyset \} \\ f(\sup \emptyset) &= \sup \emptyset \\ \Leftrightarrow \quad \{ \sup \emptyset = \perp \} \\ f(\perp) &= \perp , \end{aligned}$$

which is a rather strong property: very often we are interested in functions *not* having this property. Functions f that do satisfy $f(\perp) = \perp$ sometimes are called *strict* functions.

In a complete lattice existence of suprema is guaranteed; therefore, for complete lattices the parts “ $\sup X$ exists” and “ $\sup(f(X))$ exists” in the definition of continuity can be omitted. This is reflected by the following lemma.

4.48 Lemma. In a complete lattice (P, \sqsubseteq) a function f in $P \rightarrow P$ is continuous if and only if f distributes over suprema of all non-empty subsets of P , that is, $f(\sup X) = \sup(f(X))$, for all non-empty $X \subseteq P$.

□

4.49 Example. For a set U we recall that the poset $(\mathcal{P}(U), \subseteq)$ is a complete lattice. For every function f in $U \rightarrow U$ its lifted version, of type $\mathcal{P}(U) \rightarrow \mathcal{P}(U)$, is continuous.

□

4.50 Example. In the complete lattice $(\mathcal{P}(U \times U), \subseteq)$ of all (endo)relations on set U , relation composition is continuous (in both arguments), cf. Theorem 1.29.

□

4.8 Fixed point theorems

4.8.1 Least solutions

4.51 Definition. For any poset (P, \sqsubseteq) and for any predicate R on P , an element $p \in P$ is a *least solution* of the equation $x : R(x)$ if and only if p satisfies:

$$(23) \quad R(p) \text{ , and:}$$

$$(24) \quad (\forall x : x \in P : R(x) \Rightarrow p \sqsubseteq x)$$

Notice that condition (23) expresses that p is a *solution* of $x:R(x)$ and that (24) expresses that p is the *smallest* of all such solutions. In conjunction, they express that $p = \min\{x \in P \mid R(x)\}$.

□

In a poset not every equation $x:R(x)$ has a least solution, but *if* a least solution exists *then* it is *unique*. This follows directly from Lemma 4.10.

4.8.2 Fixed points and prefix points

4.52 Definition. For poset (P, \sqsubseteq) and for function f in $P \rightarrow P$, a value $x \in P$ is called a *fixed point* of f if $f(x) = x$, and x is called a *prefix point* of f if $f(x) \sqsubseteq x$.

□

So, by this definition, a fixed point of function f is a solution of the equation $x: f(x) = x$, whereas prefix points are solutions of the equation $x: f(x) \sqsubseteq x$. In Computer Science we are particularly interested in *least* fixed points and *least* prefix points. It is a direct corollary of Lemma 4.10 that, if they exist, such least fixed points and prefix points are unique, without further assumptions on poset (P, \sqsubseteq) or function f . Existence of such points, however, requires a little more.

4.53 Definition. For poset (P, \sqsubseteq) and for function f in $P \rightarrow P$ the *least fixed point* of f is the (if it exists, unique) value $p \in P$ satisfying:

- (2) $f(p) = p$, and:
- (3) $(\forall x: x \in P: f(x) = x \Rightarrow p \sqsubseteq x)$

The *least prefix point* of f is the (if it exists, unique) value $q \in P$ satisfying:

- (4) $f(q) \sqsubseteq q$, and:
- (5) $(\forall x: x \in P: f(x) \sqsubseteq x \Rightarrow q \sqsubseteq x)$

We denote f 's least fixed point by $fix(f)$ and its least prefix point by $pre(f)$.

□

4.54 Example. Let $b \in P$ in a poset (P, \sqsubseteq) . We consider the constant function with value b , that is, function f defined by $f(x) = b$, for all $x \in P$. Function f 's *only* fixed point, so also its least fixed point, is b : the only x satisfying $f(x) = x$ is b . The prefix points of f are all upperbounds of b , that is, all $x \in P$ satisfying $b \sqsubseteq x$. Obviously, f 's least prefix point is b .

For the identity function I_P , *every* $x \in P$ is a fixed point and, hence, a prefix point too. Hence, I_P 's least fixed point and prefix point is \perp , provided, of course, that \perp exists, that is, that P has a minimum.

□

Functions may or may not have least fixed points and/or least prefix points but *monotonic* functions have the property that, *provided* they exist, $fix(f) = pre(f)$.

4.55 Theorem. Any poset (P, \sqsubseteq) and any *monotonic* function f in $P \rightarrow P$ satisfy: $fix(f) = pre(f)$, provided $fix(f)$ and $pre(f)$ both exist.

Proof. Let p and q be f 's least fixed point and least prefix point, respectively; so, p and q satisfy (2) through (5) above. Now we prove $p = q$ by, using the antisymmetry of \sqsubseteq , proving $p \sqsubseteq q$ and $q \sqsubseteq p$ separately:

$$\begin{aligned}
 & p \sqsubseteq q \\
 \Leftarrow & \quad \{ (3) , \text{ with } x := q \} \\
 & f(q) = q \\
 \Leftarrow & \quad \{ \sqsubseteq \text{ is antisymmetric } \} \\
 & f(q) \sqsubseteq q \wedge q \sqsubseteq f(q) \\
 \Leftrightarrow & \quad \{ (4) \} \\
 & q \sqsubseteq f(q) \\
 \Leftarrow & \quad \{ (5) , \text{ with } x := f(q) \} \\
 & f(f(q)) \sqsubseteq f(q) \\
 \Leftarrow & \quad \{ f \text{ is monotonic } \} \\
 & f(q) \sqsubseteq q \\
 \Leftrightarrow & \quad \{ (4) \} \\
 & \text{true} ,
 \end{aligned}$$

and:

$$\begin{aligned}
 & q \sqsubseteq p \\
 \Leftarrow & \quad \{ (5) , \text{ with } x := p \} \\
 & f(p) \sqsubseteq p \\
 \Leftarrow & \quad \{ \sqsubseteq \text{ is reflexive } \} \\
 & f(p) = p \\
 \Leftrightarrow & \quad \{ (2) \} \\
 & \text{true} .
 \end{aligned}$$

Notice that each of the properties (2) through (5) of p and q have been used at least once.

□

Theorem 4.55 tells us that the least fixed point and the least prefix point of a monotonic function are equal, if they exist. Fortunately, the proof of this theorem gives us a little more.

4.56 Theorem. Any poset (P, \sqsubseteq) and any *monotonic* function f in $P \rightarrow P$ satisfy: if $pre(f)$ exists then so does $fix(f)$.

Proof. Let $q = \text{pre}(f)$, so q satisfies (4) and (5) in Definition 4.53. By Theorem 4.55, if $\text{fix}(f)$ exists it is equal to q , so we must prove that q is f 's least fixed point. This means that q satisfies (2) and (3) but with $p := q$. As to (2), that is, $f(q) = q$, we observe that the proof of $p \sqsubseteq q$ in Theorem 4.55 already implies this, because it is based on (4) and (5) only. As to (3) (with $p := q$) we observe, for all $x \in P$:

$$\begin{aligned} & f(x) = x \\ \Rightarrow & \quad \{ \sqsubseteq \text{ is reflexive } \} \\ & f(x) \sqsubseteq x \\ \Rightarrow & \quad \{ (5) \} \\ & q \sqsubseteq x \end{aligned}$$

□

4.8.3 Existence of fixed points

We have seen that if a monotonic function has a least prefix point then this also is its least fixed point. Monotonicity of functions all by itself, however, is not enough to guarantee the existence of least prefix points. In complete lattices, however, the situation is better.

4.57 Theorem. [Prefix Point Theorem] In a complete lattice (P, \sqsubseteq) every monotonic function f in $P \rightarrow P$ has a least prefix point.

Proof. Let X be the set of f 's prefix points, that is, $X = \{x \mid x \in P \wedge f(x) \sqsubseteq x\}$. Proving that f has a *least* prefix point then is proving that X has a *minimum*. Because (P, \sqsubseteq) is a complete lattice we do know that X has an *infimum* q , say, so: $q = \inf X$. By Lemma 4.23, to conclude that $q = \min X$ it suffices to prove that $q \in X$, that is, $f(q) \sqsubseteq q$. We recall that “ q is X 's infimum” means two things:

- (6) $(\forall x : x \in P : f(x) \sqsubseteq x \Rightarrow q \sqsubseteq x)$, and:
 (7) $(\forall x : x \in P : f(x) \sqsubseteq x \Rightarrow y \sqsubseteq x) \Rightarrow y \sqsubseteq q$, for all $y \in P$.

Now we are ready to prove:

$$\begin{aligned} & f(q) \sqsubseteq q \\ \Leftarrow & \quad \{ (7) , \text{ with } y := f(q) \} \\ & (\forall x : x \in P : f(x) \sqsubseteq x \Rightarrow f(q) \sqsubseteq x) \\ \Leftarrow & \quad \{ \sqsubseteq \text{ is transitive } \} \\ & (\forall x : x \in P : f(x) \sqsubseteq x \Rightarrow f(q) \sqsubseteq f(x)) \\ \Leftarrow & \quad \{ f \text{ is monotonic } \} \\ & (\forall x : x \in P : f(x) \sqsubseteq x \Rightarrow q \sqsubseteq x) \\ \Leftrightarrow & \quad \{ (6) \} \end{aligned}$$

true

□

4.58 Theorem. [Fixed Point Theorem] In a complete lattice (P, \sqsubseteq) every monotonic function in $P \rightarrow P$ has a least fixed point.

Proof. By the Prefix Point Theorem every monotonic function has a least prefix point. By Theorem 4.56, therefore, every monotonic function has a least fixed point as well.

□

* * *

The Prefix Point Theorem gives the least prefix (and fixed) point of a monotonic function as the infimum of the set of all prefix points. For continuous functions, however, a different characterization exists: the least fixed point of a continuous function is the supremum of all “approximations from below”.

4.59 Theorem. [Limit Theorem] In a complete lattice (P, \sqsubseteq) and for continuous function f in $P \rightarrow P$ we have $fix(f) = \sup \{ f^i(\perp) \mid 0 \leq i \}$.

Proof. We define function s , of type $\mathbb{N} \rightarrow P$, by $s_i = f^i(\perp)$, for all $i \in \mathbb{N}$. A recursive definition of s then is, for all $i, 0 \leq i$:

$$(8) \quad s_0 = \perp \wedge s_{i+1} = f(s_i) .$$

Now we have that $\{ f^i(\perp) \mid 0 \leq i \}$ equals $\{ s_i \mid 0 \leq i \}$; hence, also $\sup \{ f^i(\perp) \mid 0 \leq i \}$ equals $\sup \{ s_i \mid 0 \leq i \}$. By Theorems 4.56 and 4.55 it is sufficient to prove that $\sup \{ s_i \mid 0 \leq i \}$ is f 's least prefix point.

Firstly, we prove that $\sup \{ s_i \mid 0 \leq i \}$ is a prefix point of f :

$$\begin{aligned} & f(\sup \{ s_i \mid 0 \leq i \}) \\ = & \quad \{ f \text{ is continuous} \} \\ & \sup (f(\{ s_i \mid 0 \leq i \})) \\ = & \quad \{ \text{definition of lifted function} \} \\ & \sup (\{ f(s_i) \mid 0 \leq i \}) \\ = & \quad \{ \text{recursive property (8) of } s \} \\ & \sup \{ s_{i+1} \mid 0 \leq i \} \\ \sqsubseteq & \quad \{ \text{monotonicity of } \sup \} \\ & \sup \{ s_i \mid 0 \leq i \} . \end{aligned}$$

Secondly, we prove that $\sup \{ s_i \mid 0 \leq i \}$ is under every prefix point z , say, of f ; that is, assuming $f(z) \sqsubseteq z$ we calculate:

$$\begin{aligned}
& \sup\{s_i \mid 0 \leq i\} \sqsubseteq z \\
\Leftrightarrow & \quad \{ \text{Lemma 4.24} \} \\
& (\forall i: 0 \leq i: s_i \sqsubseteq z) \\
\Leftrightarrow & \quad \{ \text{Mathematical Induction} \} \\
& s_0 \sqsubseteq z \wedge (\forall i: 0 \leq i: s_i \sqsubseteq z \Rightarrow s_{i+1} \sqsubseteq z) \\
\Leftarrow & \quad \{ \text{recursive property (8) of } s, \text{ and } f(z) \sqsubseteq z, \text{ using transitivity of } \sqsubseteq \} \\
& \perp \sqsubseteq z \wedge (\forall i: 0 \leq i: s_i \sqsubseteq z \Rightarrow f(s_i) \sqsubseteq f(z)) \\
\Leftarrow & \quad \{ \text{property of } \perp \text{ and } f \text{ is monotonic} \} \\
& \text{true} ,
\end{aligned}$$

which concludes the proof.

□

Notice that in the above proof we have defined $s_0 = \perp$, but all we really needed about s_0 is $s_0 \sqsubseteq z$ for all prefix points z of f : any value for s_0 meeting this requirement will do.

Although we have formulated this theorem for continuous functions f , all we have *actually used* about f is $f(\sup\{s_i \mid 0 \leq i\}) = \sup(f(\{s_i \mid 0 \leq i\}))$ and that f is monotonic. Infinite sequence s in this proof is *ascending* –see the exercises– and ascending infinite sequences also are called *chains*. Thus, although continuity means that f distributes over the suprema of all non-empty subsets, all we need here is that f is monotonic and that f distributes over the suprema of chains. This entails a weaker form of continuity that is also called *chain continuity*.

Chain continuity very much resembles the common notion of continuity in traditional analysis: there, the supremum of a chain equals the *limit* of that chain, and there continuity means that if a chain has a limit then the chain of function values has a limit as well, and this limit equals the value of the function applied to the limit of the chain; that is, the function distributes over limits. (And, one of the many characteristic properties of \mathbb{R} is that every bounded chain has a limit.)

4.9 Complete Partial Orders

Theorem 4.59 has another interesting consequence. In a complete lattice *every* subset has a supremum and, hence, also an infimum. In the proof of Theorem 4.59, however, we only have used that a *chain* –that is, an ascending infinite sequence– has a supremum, and that \perp exists. The requirement that a partial order contains a minimum, \perp , and that all chains have suprema is weaker than the requirements of a complete lattice. Such a structure is called a *Complete Partial Order*, or CPO for short, and Theorem 4.59 then boils down to the proposition that in a CPO every chain-continuous function f has $\sup\{f^i(\perp) \mid 0 \leq i\}$ as its least fixed point.

Occasionally, CPOs are useful because, in a way, they are more “open-ended” than complete lattices: the whole set does not need to have a maximum, as is the

case with complete lattices, and subsets do not need to have infima. In addition, the requirement of the existence of \perp can be relaxed, as we already noticed right after Theorem 4.59.

4.10 Boolean algebras

In Section 4.5 we have seen that a lattice is a poset with two additional binary operators, \sqcup and \sqcap , that happen to have the following algebraic properties, for all $x, y, z \in P$:

- $x \sqcup x = x$ and $x \sqcap x = x$ (*idempotence*);
- $x \sqcup y = y \sqcup x$ and $x \sqcap y = y \sqcap x$ (*commutativity*);
- $x \sqcup (y \sqcup z) = (x \sqcup y) \sqcup z$ and $x \sqcap (y \sqcap z) = (x \sqcap y) \sqcap z$ (*associativity*);
- $x \sqcup (x \sqcap y) = x$ and $x \sqcap (x \sqcup y) = x$ (*absorption*).

We have also seen –Theorem 4.32– that, conversely, every set in which two binary operators have these algebraic properties is a lattice, if the partial order relation is defined in terms of the operators in the correct way.

In addition, a distributive lattice is a lattice in which the operators distribute over each other, that is, for all $x, y, z \in P$:

- $x \sqcup (y \sqcap z) = (x \sqcup y) \sqcap (x \sqcup z)$ (\sqcup *distributes over* \sqcap);
- $x \sqcap (y \sqcup z) = (x \sqcap y) \sqcup (x \sqcap z)$ (\sqcap *distributes over* \sqcup).

A particularly interesting special class of distributive lattices is formed by the so-called *Boolean algebras*, which are distributive lattices with a few additional algebraic properties. A Boolean algebra has extreme elements \top and \perp , and an additional unary operator called “complement” or “negation”.

Remark: Boolean algebras are named after the British mathematician George Boole (1815-1864), who was the first to investigate the algebraic structure of the operations \cup and \cap on sets.

□

4.60 Definition. A Boolean algebra is a set V , say, containing two elements \top (“top”) and \perp (“bottom”), two binary operators \sqcup (“cup”) and \sqcap (“cap”), and a unary operator \neg (“not”), with the following algebraic properties:

- (a) \sqcup and \sqcap are idempotent, commutative, and associative;
- (b) \sqcup and \sqcap distribute over each other;
- (c) *absorption*: $x \sqcup (x \sqcap y) = x$ and $x \sqcap (x \sqcup y) = x$, for all $x, y \in V$;
- (d) *zero elements*: $x \sqcup \top = \top$ and $x \sqcap \perp = \perp$, for all $x \in V$;

- (e) *identity elements*: $x \sqcup \perp = x$ and $x \sqcap \top = x$, for all $x \in V$;
- (f) *complement properties*: $x \sqcup \neg x = \top$ and $x \sqcap \neg x = \perp$, for all $x \in V$.

□

The list of algebraic properties in this definition is not *minimal*: the absorption rule (c) can be proved from (a) through (e) without (c), and, also, rule (e) for identity elements can be proved from (a) through (d). So, either rule (c) or rule (e) (but not both!) could be omitted without affecting the definition.

Notice that only rule (f) seemingly gives only little, rather implicit, information about the complement operator \neg . Nevertheless, this rule, together with the other rules, gives enough information, as various other properties can be derived.

4.61 Lemma. In a Boolean algebra $(V, \top, \perp, \sqcup, \sqcap, \neg)$ the complement of every element $x \in V$ is *unique*. That is, for every $x, y, z \in V$ we have:

$$x \sqcup y = \top \wedge x \sqcap y = \perp \wedge x \sqcup z = \top \wedge x \sqcap z = \perp \Rightarrow y = z .$$

Proof. Assuming the left-hand side of the implication we derive:

$$\begin{aligned} & y \\ = & \quad \{ \text{identity elements, to introduce } \perp \} \\ & y \sqcup \perp \\ = & \quad \{ x \sqcap z = \perp \} \\ & y \sqcup (x \sqcap z) \\ = & \quad \{ \sqcup \text{ distributes over } \sqcap \} \\ & (y \sqcup x) \sqcap (y \sqcup z) \\ = & \quad \{ \text{commutativity and } x \sqcup y = \top \} \\ & \top \sqcap (y \sqcup z) \\ = & \quad \{ \text{commutativity and identity elements} \} \\ & y \sqcup z , \end{aligned}$$

from which we conclude that $y = y \sqcup z$; because the situation is symmetric in y and z we also conclude $z = y \sqcup z$; from these two we obtain $y = z$, as required.

□

In the proof of this lemma we have invoked “commutativity” twice. Usually, however, commutativity and associativity are considered so elementary that they are taken for granted and used without explicitly mentioning.

This lemma is important: if some $x, y \in V$ satisfy $x \sqcup y = \top$ and $x \sqcap y = \perp$ then we may conclude, also using complement properties (f) that $y = \neg x$. Hence, to prove $y = \neg x$ it suffices to show that $x \sqcup y = \top$ and $x \sqcap y = \perp$. We will call this proof obligation “satisfying the complement properties”.

4.62 Exercise. Prove that:

- (a) $\neg \perp = \top$ and $\neg \top = \perp$;
- (b) $\neg(\neg x) = x$, for all $x \in V$.

□

4.63 Lemma. [de Morgan's rules] In a Boolean algebra $(V, \top, \perp, \sqcup, \sqcap, \neg)$ we have:

$$\neg(x \sqcup y) = \neg x \sqcap \neg y \text{ and } \neg(x \sqcap y) = \neg x \sqcup \neg y, \text{ for all } x, y \in V.$$

Proof. We prove the first conjunct only. By Lemma 4.61, to prove this it suffices to prove that $\neg x \sqcap \neg y$ satisfies the complement properties:

$$\begin{aligned} & (x \sqcup y) \sqcup (\neg x \sqcap \neg y) \\ = & \quad \{ \sqcup \text{ distributes over } \sqcap \} \\ & ((x \sqcup y) \sqcup \neg x) \sqcap ((x \sqcup y) \sqcup \neg y) \\ = & \quad \{ x \sqcup \neg x = \top \text{ and } y \sqcup \neg y = \top \} \\ & (y \sqcup \top) \sqcap (x \sqcup \top) \\ = & \quad \{ \text{zero elements (twice)} \} \\ & \top \sqcap \top \\ = & \quad \{ \text{idempotence} \} \\ & \top, \end{aligned}$$

and:

$$\begin{aligned} & (x \sqcup y) \sqcap (\neg x \sqcap \neg y) \\ = & \quad \{ \sqcap \text{ distributes over } \sqcup \} \\ & (x \sqcap (\neg x \sqcap \neg y)) \sqcup (y \sqcap (\neg x \sqcap \neg y)) \\ = & \quad \{ x \sqcap \neg x = \perp \text{ and } y \sqcap \neg y = \perp \} \\ & (\perp \sqcap \neg y) \sqcup (\perp \sqcap \neg x) \\ = & \quad \{ \text{zero elements (twice)} \} \\ & \perp \sqcap \perp \\ = & \quad \{ \text{idempotence} \} \\ & \perp. \end{aligned}$$

□

4.64 Example. For every set V its power set $\mathcal{P}(V)$ is a Boolean algebra, with V , \emptyset , \cup , and \cap for \top , \perp , \sqcup , and \sqcap , respectively, and with complement C – defined by $X^C = V \setminus X$ – for \neg .

□

4.65 Example. We consider set \mathbb{B} with $\mathbb{B} = \{1, 0\}$. On \mathbb{B} we define operators \vee , \wedge , and \neg as follows, for all $x, y \in \mathbb{B}$: $x \vee y = x \max y$, $x \wedge y = x \min y$, $\neg 1 = 0$, and $\neg 0 = 1$. Then $(\mathbb{B}, 1, 0, \vee, \wedge, \neg)$ is a Boolean algebra. This is the *smallest possible* Boolean algebra.

□

4.66 Example. We consider a set V and the set $V \rightarrow \mathbb{B}$ of all functions from V to \mathbb{B} . On this set we define operators \vee (“lifted or”), \wedge (“lifted and”), and \neg (“lifted negation”) as follows, for all functions $f, g \in V \rightarrow \mathbb{B}$. Firstly, $f \vee g$ is the function h in $V \rightarrow \mathbb{B}$ defined by $h(v) = f(v) \vee g(v)$, for all $v \in V$; secondly, $f \wedge g$ is the function h in $V \rightarrow \mathbb{B}$ defined by $h(v) = f(v) \wedge g(v)$, for all $v \in V$; thirdly, $\neg f$ is the function h in $V \rightarrow \mathbb{B}$ defined by $h(v) = \neg(f(v))$, for all $v \in V$; Finally, we introduce \top in $V \rightarrow \mathbb{B}$ as the constant function whose value is everywhere 1 and \perp in $V \rightarrow \mathbb{B}$ as the constant function whose value is everywhere 0. Then $((V \rightarrow \mathbb{B}), \top, \perp, \vee, \wedge, \neg)$ is a Boolean algebra.

□

4.11 Applications

The fixed point theorems play a role in defining the meaning of recursive definitions. Very often sets, and data structures in particular, and functions are defined recursively. This raises the question, however, what exactly is defined by such a recursive definition. The fixed point theorems provide the answer to this question. We illustrate this with some examples.

4.11.1 Recursively defined sets

Sometimes it is convenient to define sets recursively. As an example we consider the (so-called) “Hamming set” which is the set of all positive natural numbers that are divisible by 2, 3, and 5 only. Calling this set H it can be defined recursively as follows:

$$(9) \quad 1 \in H ;$$

$$(10) \quad 2*n \in H \wedge 3*n \in H \wedge 5*n \in H , \text{ for all } n \in H ;$$

$$(11) \quad \text{The only elements of } H \text{ are those required by (9) and (10) .}$$

This “definition” can be, and usually is, interpreted as follows. Apparently, this definition is intended to define a subset of \mathbb{N} , or even of \mathbb{N}^+ if you like. Rule (9) can be rewritten to the following equivalent form:

$$(12) \quad \{1\} \subseteq H .$$

Furthermore, we introduce functions $f2$, $f3$, and $f5$, say, defined by, for all $n \in \mathbb{N}$:

$$f2(n) = 2*n \wedge f3(n) = 3*n \wedge f5(n) = 5*n .$$

Using these functions lifted, we can reformulate rule (10) equivalently as:

$$(13) \quad f2(H) \subseteq H \wedge f3(H) \subseteq H \wedge f5(H) \subseteq H .$$

Using that, for all sets, $A \subseteq C \wedge B \subseteq C$ is equivalent to $A \cup B \subseteq C$, we can now combine (12) and (13) into:

$$(14) \quad \{1\} \cup f2(H) \cup f3(H) \cup f5(H) \subseteq H .$$

Thus, the single rule (14) is equivalent to the conjunction of rules (9) and (10), and rule (11) now can be taken into account by defining H as the *smallest* of all sets satisfying (14). It now also is clear that this latter requirement is essential to make such a definition meaningful: \mathbb{N} is a solution to (14) too, but not a very interesting one. (Actually, of course, \mathbb{N} is the *largest* of all subsets of \mathbb{N} satisfying (14).)

Formula (14) can be abbreviated further by the introduction of a function F , of type $\mathcal{P}(\mathbb{N}) \rightarrow \mathcal{P}(\mathbb{N})$, and defined by, for all $X \subseteq \mathbb{N}$:

$$(15) \quad F(X) = \{1\} \cup f2(X) \cup f3(X) \cup f5(X) .$$

In terms of this function formula (14) just expresses that H is a prefix point of F :

$$F(H) \subseteq H ,$$

and that H is the smallest of all sets satisfying (13) now boils down to the requirement that H is the *least* prefix point of F . Because, by Theorem 4.55, the least prefix point of any monotonic function equals its least fixed point, we conclude that set H also satisfies:

$$F(H) = H .$$

As we know, the poset $(\mathcal{P}(\mathbb{N}), \subseteq)$ is a complete lattice, and in this structure function F is continuous. Thus, by the fixed point theorems, F has a least prefix point, which is equal to its least fixed point, and which is given by the explicit formula:

$$\left(\bigcup_{i:0 \leq i} F^i(\emptyset) \right) .$$

4.11.2 The natural numbers and Mathematical Induction

In the previous example we have defined the Hamming set as a subset of the natural numbers. This enabled us to define the Hamming set as a least prefix point of a continuous function of type $\mathcal{P}(\mathbb{N}) \rightarrow \mathcal{P}(\mathbb{N})$, in the complete lattice $(\mathcal{P}(\mathbb{N}), \subseteq)$. But what about the natural numbers themselves, that is, how are these defined?

A very common way to define the set, \mathbb{N} , of natural numbers is by means of the axioms by Peano. These involve a “constant” 0 (“zero”) and a “function” *succ* (“successor”), in terms of which \mathbb{N} is defined in the following recursive way:

$$(16) \quad 0 \in \mathbb{N} ;$$

$$(17) \quad succ(n) \in \mathbb{N} , \quad \text{for all } n \in \mathbb{N} ;$$

$$(18) \quad \text{The only elements of } \mathbb{N} \text{ are those required by (16) and (17) .}$$

This resembles the previous example very much, but there is an important difference: the Hamming set can be defined as a subset of \mathbb{N} but we cannot, of course, define \mathbb{N} as a subset of itself. In addition, what should be the types of constant 0 and function $succ$? From (16) and (17) we infer that 0 and $succ$ are likely to have types \mathbb{N} and $\mathbb{N} \rightarrow \mathbb{N}$, respectively. Unfortunately, 0 and $succ$ are the constructors *used* to define \mathbb{N} , so stating their types in terms of the set to be defined is unsatisfactory: we should rather consider the constructors, and their types, as given a priori, independently of the set yet to be defined, so as to avoid circular dependencies.

These issues are resolved if we postulate the existence of a set Ω , the “universe of values”, that provides the set of which \mathbb{N} , and all other datatypes for that matter, will be subsets: then we can define \mathbb{N} as a subset of Ω , and then the poset $(\mathcal{P}(\Omega), \subseteq)$ is the complete lattice needed for our purpose¹¹.

In terms of set Ω we now require that $0 \in \Omega$ and that $succ$ has type $\Omega \rightarrow \Omega$. This brings us onto familiar grounds again. Just as in the previous example we rewrite (16) and (17), by lifting and combining, into the following equivalent form:

$$(19) \quad \{0\} \cup succ(\mathbb{N}) \subseteq \mathbb{N} .$$

This shows that we can define \mathbb{N} as the least prefix point of a function F , of type $\mathcal{P}(\Omega) \rightarrow \mathcal{P}(\Omega)$, and defined by, for all $X \subseteq \Omega$:

$$F(X) = \{0\} \cup succ(X) .$$

As before, function F is continuous and F 's least prefix/fixed point is the set $(\bigcup_{i:0 \leq i} F^i(\emptyset))$, which boils down to $\{succ^i(0) \mid 0 \leq i\}$.

* * *

So, set \mathbb{N} is the smallest solution of the equation, in unknown X :

$$(20) \quad \{0\} \cup succ(X) \subseteq X .$$

We recall that this means that

$$(21) \quad \{0\} \cup succ(\mathbb{N}) \subseteq \mathbb{N} ,$$

which expresses that \mathbb{N} is *just a* solution, and

$$(22) \quad \{0\} \cup succ(X) \subseteq X \Rightarrow \mathbb{N} \subseteq X , \text{ for all } X \subseteq \Omega ,$$

which expresses that \mathbb{N} is a subset of all solutions, that is, \mathbb{N} is the *smallest* solution.

Now let R be a predicate on \mathbb{N} . That R is universally true, that is, is true for all elements of \mathbb{N} , is usually formulated as $(\forall n : n \in \mathbb{N} : R(n))$, but in more set theoretic terms this can also be formulated as: the subset of \mathbb{N} of all natural numbers for which R is true, is equal to the whole \mathbb{N} . Formally, this boils down to $\{n \in \mathbb{N} \mid R(n)\} = \mathbb{N}$. Now, how would we prove such a proposition? One possible approach now is:

¹¹Demonstrating the existence of such a universe is beyond the scope of this text. Actually, we seem to have a chicken-and-egg problem here: to define an infinite set like \mathbb{N} recursively we need an infinite universe of values like Ω , but defining infinite sets always seems to require recursion!

$$\begin{aligned}
& \{n \in \mathbb{N} \mid R(n)\} = \mathbb{N} \\
\Leftrightarrow & \quad \{ \text{mutual set inclusion} \} \\
& \{n \in \mathbb{N} \mid R(n)\} \subseteq \mathbb{N} \quad \wedge \quad \mathbb{N} \subseteq \{n \in \mathbb{N} \mid R(n)\} \\
\Leftrightarrow & \quad \{ \text{by definition, set } \{n \in \mathbb{N} \mid R(n)\} \text{ is a subset of } \mathbb{N} \} \\
& \mathbb{N} \subseteq \{n \in \mathbb{N} \mid R(n)\} \\
\Leftarrow & \quad \{ (22), \text{ with } X := \{n \in \mathbb{N} \mid R(n)\} \} \\
& \{0\} \cup \text{succ}(\{n \in \mathbb{N} \mid R(n)\}) \subseteq \{n \in \mathbb{N} \mid R(n)\} \\
\Leftrightarrow & \quad \{ \cup/\subseteq\text{-connection} \} \\
& \{0\} \subseteq \{n \in \mathbb{N} \mid R(n)\} \quad \wedge \quad \text{succ}(\{n \in \mathbb{N} \mid R(n)\}) \subseteq \{n \in \mathbb{N} \mid R(n)\} \\
\Leftrightarrow & \quad \{ \text{definition of } \subseteq; \text{Property 3.7(f) (of lifted functions)} \} \\
& 0 \in \{n \in \mathbb{N} \mid R(n)\} \quad \wedge \quad (\forall n: n \in \mathbb{N} \wedge R(n) : \text{succ}(n) \in \{n \in \mathbb{N} \mid R(n)\}) \\
\Leftrightarrow & \quad \{ \text{“definition” of } \in \text{ (twice), using } 0 \in \mathbb{N} \text{ and } \text{succ}(n) \in \mathbb{N} \} \\
& R(0) \quad \wedge \quad (\forall n: n \in \mathbb{N} \wedge R(n) : R(\text{succ}(n))) \\
\Leftrightarrow & \quad \{ \text{range-term trading} \} \\
& R(0) \quad \wedge \quad (\forall n: n \in \mathbb{N} : R(n) \Rightarrow R(\text{succ}(n))) \quad .
\end{aligned}$$

Thus, we have proved that $(\forall n: n \in \mathbb{N} : R(n))$ is equivalent to $R(0) \wedge (\forall n: n \in \mathbb{N} : R(n) \Rightarrow R(\text{succ}(n)))$, so to prove the former it is sufficient to prove the latter, which is formally weaker and, therefore, usually easier. This is the well-known Principle of Mathematical Induction: we see that, in the above lattice-theoretical setting, the validity of this principle can be *proved*. Notice that in this proof we have used property (22), and we cannot do without it: the Principle of Mathematical Induction is a direct consequence of the part that defines \mathbb{N} as the *smallest* solution of equation (20).

4.11.3 Finite lists

A well-known datatype, particularly in Functional Programming, is the datatype of *finite lists*. Usually, all elements of a (finite) list have the same type, called the *element type*. The datatype of lists with elements of type B , say, then is defined informally in the following way, by means of two, so-called, *constructors*, here denoted by $[]$ (“empty”) and \triangleright (“cons”):

- (23) $[]$ is a list ;
- (24) $b \triangleright s$ is a list, for all $b \in B$ and for all lists s ;
- (25) The only lists are those generated by (23) and (24) .

Formally, a datatype is just a set. The datatype of lists with elements of type B is a set $\mathcal{L}_*(B)$, say, and the above informal rules can then be rendered more formally, as follows:

$$(23) \quad [] \in \mathcal{L}_*(B) \ ;$$

$$(24) \quad b \triangleright s \in \mathcal{L}_*(B) \ , \text{ for all } b \in B \text{ and } s \in \mathcal{L}_*(B) \ ;$$

$$(25) \quad \text{The only elements of } \mathcal{L}_*(B) \text{ are those required by (23) and (24) .}$$

Thus, a recursively defined datatype just is a recursively defined set, and to interpret a recursive definition like the above one we can apply the same techniques as in the previous subsections. Apparently, again, $\mathcal{L}_*(B)$ is the *smallest* of all sets satisfying (23) and (24).

If, however, we wish to define the datatype as a least prefix point we do need a function in a complete lattice. So, this raises the question: what lattice? In addition, this raises the question: what are the types of the constructors $[]$ and \triangleright ?

As in the previous subsection, we resolve these issues by postulating the existence of a set Ω , the “universe of values”, of which our datatype will be a subset: now we can define $\mathcal{L}_*(B)$ as a subset of Ω , and the poset $(\mathcal{P}(\Omega), \subseteq)$ is the complete lattice needed for our purpose.

In terms of set Ω we now require that $[] \in \Omega$ and that \triangleright has type $B \times \Omega \rightarrow \Omega$. We define datatype $\mathcal{L}_*(B)$ as the *smallest* subset of Ω satisfying (23) and (24). Thus, we are on familiar grounds again. For every $b \in B$ we define an associated function f_b , of type $\Omega \rightarrow \Omega$, by:

$$f_b(x) = b \triangleright x \ , \text{ for all } x \in \Omega \ .$$

Now rules (23) and (24) above can be rewritten into:

$$(26) \quad \{ [] \} \subseteq \mathcal{L}_*(B) \ ;$$

$$(27) \quad f_b(\mathcal{L}_*(B)) \subseteq \mathcal{L}_*(B) \ , \text{ for all } b \in B \ .$$

We now define a function F , of type $\mathcal{P}(\Omega) \rightarrow \mathcal{P}(\Omega)$, by:

$$(28) \quad F(X) = \{ [] \} \cup \left(\bigcup_{b \in B} f_b(X) \right) \ .$$

Function F is continuous and $\mathcal{L}_*(B)$ appears as the least prefix point, hence also as the least fixed point, of F . Notice that $(\bigcup_{b \in B} f_b(X))$ just denotes the set $\{ b \triangleright x \mid b \in B \wedge x \in X \}$.

Remark: Although the above reasoning is valid for all constructors $[]$ and \triangleright having the right types, some additional requirements are needed to yield a *useful* datatype. For example, the above does allow the possibility that $b \triangleright [] = []$, for all $b \in B$. In that case the datatype collapses: then we would have $\mathcal{L}_*(B) = \{ [] \}$, which is not particularly interesting. To avoid this, it is usual to require that $b \triangleright x \neq []$, for all b and x , and it also is usual to require that function \triangleright is injective. The latter is achieved by postulating the existence of two additional functions hd , of type $\Omega \rightarrow B$, and tl , of type $\Omega \rightarrow \Omega$, say, satisfying $hd(b \triangleright x) = b$ and $tl(b \triangleright x) = b$.

These additional requirements are needed to obtain a practically useful datatype of finite lists, but they are not needed for the definition *per se* of the datatype as a least prefix point.

□

4.11.4 Closures of relations

We recall that an (endo)relation on a set U is a subset of $U \times U$. In Chapter 1 we have defined several closures of relations as smallest solutions of equations. All of these equations have the following general shape, where R is a given relation on U and where X is the unknown of the equation:

$$(29) \quad F_R(X) \subseteq X \quad ,$$

in which F_R is a function of type $U \times U \rightarrow U \times U$, so F_R maps relations on U to relations on U ; generally F_R depends on R .

In all cases, the closure of R then is defined as the smallest solution of (29), that is, as the least prefix point of F_R in the complete lattice $(U \times U, \subseteq)$. We now know, on account of the Prefix Point Theorem, that monotonicity of function F_R is sufficient to guarantee the existence of F_R 's least prefix point, and that it is equal to F_R 's least fixed point.

In all cases of relation closure, function F_R is even continuous, hence its least prefix point is given explicitly by the Limit Theorem.

For example, for the reflexive closure of R function F_R can be defined by:

$$F_R(X) = R \cup I \quad , \text{ for all } X \subseteq U \times U \quad .$$

This example is not very exciting, because F_R is constant, so its prefix and fixed points are its constant value, that is, $R \cup I$.

Slightly more interesting is the symmetric closure of relation R . In this case function F_R can be defined by:

$$F_R(X) = R \cup X^T \quad , \text{ for all } X \subseteq U \times U \quad .$$

Because $F_R(F_R(X)) = R \cup R^T \cup X$, it is not difficult to show that, in this case, F_R 's least prefix point equals $R \cup R^T$.

The most interesting are the cases of the transitive closure and the reflexive-transitive closures. We leave the analysis of the latter to the reader, and for the transitive closure F_R can be defined by:

$$F_R(X) = R \cup X;X \quad , \text{ for all } X \subseteq U \times U \quad .$$

Now the theorems in this chapter can be used to prove theorems like Theorem 1.31 and Theorem 1.32 in Chapter 1. Proving Theorem 1.31, for example, requires defining a function G_R by:

$$G_R(X) = R \cup X;R \quad , \text{ for all } X \subseteq U \times U \quad ,$$

and, then, proving that G_R and F_R have the same least prefix points.

4.11.5 Grammars and languages

We confine ourselves to a simple example, but the pattern of reasoning is applicable to all context-free grammars.

We consider the following context-free grammar, with non-terminal symbol S and terminal symbols a and b :

$$S := \varepsilon \mid aSb .$$

This grammar represents a “language” L_S , in the following way. Let U be the set of all finite sequences containing a ’s and b ’s only. The empty such sequence is denoted by ε and we use simple juxtaposition for sequence extensions; for example, for $s \in U$ we write as for the sequence obtained from s by prefixing it with a and we write sb for the sequence obtained by postfixing s with b . Sequence extension is associative: as $(as)b$ and $a(sb)$ are the same sequence we can omit the parantheses and simply write asb .

The language L_S represented by the above grammar now is defined as follows. It is a subset of U , and it satisfies:

$$(30) \quad \varepsilon \in L_S ;$$

$$(31) \quad asb \in L_S , \text{ for all } s \in L_S ;$$

$$(32) \quad L_S \text{ contains no other elements than as required by (30) and (31) .}$$

A more formal way to define what this means is as follows. First, we define function f in $U \rightarrow U$ by:

$$f(s) = asb , \text{ for all } s \in U .$$

Now rules (a) and (b) above can be reformulated and combined as follows; notice that, again, we use the lifted version of f , of type $\mathcal{P}(U) \rightarrow \mathcal{P}(U)$, here:

$$(33) \quad \{\varepsilon\} \cup f(L_S) \subseteq L_S .$$

Rule (c) above now means that L_S is the *smallest* of all possible sets satisfying (33), that is, L_S is the smallest solution of the equation, with unknown $X \subseteq U$:

$$\{\varepsilon\} \cup f(X) \subseteq X .$$

Again, the poset $(\mathcal{P}(U), \subseteq)$ is a complete lattice, and in this lattice L_S now is the least prefix/fixed point of function F , of type $\mathcal{P}(U) \rightarrow \mathcal{P}(U)$, and defined by, for all $X \subseteq U$:

$$F(X) = \{\varepsilon\} \cup f(X) .$$

4.12 Exercises

1. We consider a complete lattice (P, \sqsubseteq) . So, every subset $X \subseteq P$ has a (unique) supremum $\sup X$. Therefore, \sup is a function from poset $(\mathcal{P}(P), \subseteq)$ to poset (P, \sqsubseteq) .
 - (a) Prove that function \sup is monotonic, that is: $X \subseteq Y \Rightarrow \sup X \sqsubseteq \sup Y$, for all $X, Y \in \mathcal{P}(P)$.
 - (b) In which way can, similarly, function \inf be considered monotonic?

2. Assuming that a given subset $X \subseteq P$ in a poset (P, \sqsubseteq) has the following properties: $\sup X$ exists and $b \sqsubseteq c$ for some $b \in P$ and for some $c \in X$. Prove that $\sup(X \cup \{b\}) = \sup X$.
3. Prove that infinite sequence s , as recursively defined by (8) in the proof of Theorem 4.59, is *ascending*, that is: $(\forall i: 0 \leq i: s_i \sqsubseteq s_{i+1})$.
4. We consider a set V defined by $V = \mathbb{N} \cup \{y, z\}$, where y and z are two elements *not* in \mathbb{N} ; so, V is the naturals *extended with* two new elements, here called y and z , with $y \neq z$. On V we define a partial order \sqsubseteq , as follows:

$$\begin{aligned} m \sqsubseteq n &\Leftrightarrow m \leq n, \text{ for all } m, n \in \mathbb{N}; \\ m \sqsubseteq y \wedge m \sqsubseteq z &, \text{ for all } m \in \mathbb{N}; \\ y \sqsubseteq y \wedge y \sqsubseteq z \wedge z \sqsubseteq z &. \end{aligned}$$

In $V \rightarrow V$ we define a function f by, for all $m \in \mathbb{N}$:

$$f(m) = m+1 \wedge f(y) = z \wedge f(z) = z.$$

- (a) Prove that (V, \sqsubseteq) is a complete lattice.
 - (b) Prove that function f is monotonic and that f distributes over \sqcup .
 - (c) Show that f is *not* continuous.
 - (d) What is $\text{fix}(f)$ and what is $\sup\{f^i(\perp) \mid 0 \leq i\}$?
5. We consider function F , as defined by formula (15) in Subsection 4.11.1.
 - (a) Prove that, for all $n \in \mathbb{N}$: $F^n(\emptyset) = \{2^i * 3^j * 5^k \mid i, j, k \in \mathbb{N} \wedge i+j+k < n\}$.
 - (b) Prove that F 's least fixed point equals $\{2^i * 3^j * 5^k \mid i, j, k \in \mathbb{N}\}$.
 6. We consider function F , as defined by (28) in Subsection 4.11.3. How can F^n , for $n \in \mathbb{N}$, be interpreted?
 7. The (so-called) "Fibonacci sequence" is the function fib , of type $\mathbb{N} \rightarrow \mathbb{N}$, and defined recursively by, for all $n \in \mathbb{N}$:

$$\text{fib}_0 = 0 \wedge \text{fib}_1 = 1 \wedge \text{fib}_{n+2} = \text{fib}_n + \text{fib}_{n+1}.$$

Prove that thus defined fib is a function in $\mathbb{N} \rightarrow \mathbb{N}$ indeed.

5 Monoids and Groups

5.1 Operators and their properties

We consider a set V . A *binary operator on V* is a function of type $V \times V \rightarrow V$, so such a function maps pairs of elements of V to elements of V . Very often, applications of a binary operator are written in *infix-notation*, that is, the function's name is written *in between* the arguments.

In this chapter we use $*$ to denote any operator on a set V . Using infix-notation, we write the application of $*$ to pair $(x, y) \in V \times V$ as $x * y$ (instead of the more standard *prefix-notation* $*(x, y)$).

Sometimes binary operators have special properties, which deserve to be named.

5.1 Definition. Let $*$ be a binary operator on a set V . Then $*$ is called:

- *idempotent*, if for all $x \in V$ we have: $x * x = x$;
- *commutative*, if for all $x, y \in V$ we have: $x * y = y * x$;
- *associative*, if for all $x, y, z \in V$ we have: $(x * y) * z = x * (y * z)$;

□

5.2 Examples.

- (a) On \mathbb{N} addition, $+$, and multiplication, $*$, are binary operators; both are commutative and associative but not idempotent.
- (b) On \mathbb{Z} addition, $+$, and subtraction, $-$, are binary operators. Subtraction is neither idempotent, nor commutative, nor associative.
- (c) On \mathbb{Z} maximum, \max , and minimum, \min , are binary operators; both are idempotent, commutative, and associative.
- (d) On the set of finite lists concatenation, $++$, is a binary operator; it is neither idempotent nor commutative but it is associative.
- (e) On the set of all relations relation composition, $;$, is an associative binary operator; as a special case, so is function composition, \circ , on the set of all functions.

□

For an associative operator $*$, the expressions $(x * y) * z$ and $x * (y * z)$ are equal and, therefore, we may safely omit the parentheses and write $x * y * z$ instead. Thus we do not only save a little writing, we also avoid the choice between two forms that are equivalent. Therefore, particularly with associative operators it pays to use infix-notation. With prefix-notation no parentheses can be omitted and we are always forced to decide whether to write $*(x, *(y, z))$ or $*(*(x, y), z)$: a rather irrelevant choice!

More important than the possibility to omit parentheses, however, is that associativity offers a *manipulative* opportunity in proofs: if we have a formula of the shape $(x * y) * z$ and if $*$ is associative then we may reposition the parentheses and obtain $x * (y * z)$ (and vice versa). So, even if we have omitted the parentheses we better stay aware of their (hidden) presence and of the opportunity to reposition them. We will see examples of this.

5.3 Definition. Let $*$ be a binary operator on a set V . An element $I \in V$ is called $*$'s *identity (element)* if it satisfies, for all $x \in V$:

$$x * I = x \wedge I * x = x .$$

□

Not every binary operator has an identity element but every operator has *at most one* identity element; that is, if it exists an operator's identity element is unique.

5.4 Lemma. Let I and J both be identity elements of binary operator $*$ on a set V . Then $I = J$.

Proof. By calculation:

$$\begin{aligned} & I \\ = & \quad \{ J \text{ is identity element, with } x := I \} \\ & I * J \\ = & \quad \{ I \text{ is identity element, with } x := J \} \\ & J . \end{aligned}$$

□

5.5 Examples.

- (a) On \mathbb{N} and \mathbb{Z} the identity element of $+$ is 0, and the identity of $*$ is 1.
- (b) On \mathbb{N}^+ operator $+$ has no identity element.
- (c) On \mathbb{Z} operator $-$ has no identity element.
- (d) On \mathbb{Z} operators \max and \min have no identity elements; on \mathbb{N} , however, the identity element of \max is 0 whereas \min still has no identity element.
- (e) On the set of finite lists the identity element of $++$ is $[\]$ – the *empty* list –.
- (f) The identity element of both relation composition and function composition is the identity relation/function, I .

□

Sometimes we are interested in the relation between two (or even more) binary operators. Although we do not elaborate this in this text, we mention one property that already has been used extensively in previous chapters.

5.6 Definition. Let $*$ and $+$ (say) be binary operators on a set V . Then we say that $*$ *distributes (from the left) over* $+$ if, for all $x, y, z \in V$:

$$x * (y + z) = (x * y) + (x * z) .$$

Similarly, $*$ *distributes (from the right) over* $+$ if, for all $x, y, z \in V$:

$$(y + z) * x = (y * x) + (z * x) .$$

Of course, if $*$ is commutative the distinction between “left” and “right” disappears and we just say “distributes over”. (This is the case for almost all examples we have seen, with relation composition as the most notable exception.)

□

5.7 Examples.

- (a) On \mathbb{N} and \mathbb{Z} multiplication, $*$, distributes over addition, $+$, but addition does not distribute over multiplication.
- (b) On \mathbb{Z} operator \max distributes over \min and \min distributes over \max .
- (c) On \mathbb{Z} operator $+$ distributes over both \max and \min .
- (d) On \mathbb{N} operator $*$ distributes over both \max and \min , whereas this is not true on \mathbb{Z} .
- (e) Set union, \cup , distributes over set intersection, \cap , and vice versa.
- (f) Relation composition, $;$, distributes, from the left and from the right, over union of relations, \cup . (Recall that composition is not commutative.)

□

5.2 Semigroups and monoids

So-called *algebraic structures* are sets with operators having particular properties. The simplest such structure is called a *semigroup*, which is just a set with an associative operator.

5.8 Definition. Let $*$ be a binary operator on a set V . The pair $(V, *)$ is a *semigroup* if $*$ is associative.

□

If the operator in a semigroup has an identity element the structure already becomes a little more interesting.

5.9 Definition. A *monoid* is a triple $(V, *, I)$, where $*$ is an associative binary operator on set V , so $(V, *)$ is a semigroup, and $I, I \in V$, is the identity element of $*$.

□

5.10 Examples.

- (a) $(\mathbb{N}, +, 0)$ is a monoid, whereas $(\mathbb{N}^+, +)$ is a semigroup but not a monoid.
- (b) both $(\mathbb{N}, *, 1)$ and $(\mathbb{N}^+, *, 1)$ are monoids.
- (c) Let \mathcal{L}_* denote the set of all finite lists, and let \mathcal{L}_+ denote the set of all non-empty finite lists. Then $(\mathcal{L}_*, ++, [])$ is a monoid, whereas $(\mathcal{L}_+, ++)$ only is a semigroup.
- (d) All relations on a set with operator $;$ and identity relation I form a monoid. Similarly, all functions on a set, with function composition and the identity function form a monoid too.

□

5.11 Definition. Let $(V, *, I)$ be a monoid. For all $x \in V$ and for all $n \in \mathbb{N}$ we define x^n recursively, as follows:

$$x^0 = I \wedge x^{n+1} = x * x^n .$$

□

5.12 Lemma. Let $(V, *, I)$ be a monoid. For every $x \in V$ and for all $m, n \in \mathbb{N}$ we have:

$$x^{m+n} = x^m * x^n .$$

□

5.13 Definition. Let $(V, *, I)$ be a monoid. For every $x \in V$ an element $y \in V$ is called an *inverse* of x (with respect to $*$) if and only if: $x * y = I \wedge y * x = I$.

□

An element of a monoid does not necessarily have inverses, but if it has an inverse, it is unique. This is stated by the following lemma.

5.14 Lemma. Let $(V, *, I)$ be a monoid. Let $x, y, z \in V$ satisfy:

$$x * y = I \wedge y * x = I \wedge x * z = I \wedge z * x = I .$$

Then $y = z$.

Proof. By calculation:

$$\begin{aligned}
& y \\
= & \{ I \text{ is identity of } * \} \\
& y * I \\
= & \{ x * z = I \} \\
& y * (x * z) \\
= & \{ * \text{ is associative} \} \\
& (y * x) * z \\
= & \{ y * x = I \} \\
& I * z \\
= & \{ I \text{ is identity of } * \} \\
& z .
\end{aligned}$$

□

Notice that in the proof of this lemma we have only used $x * z = I$ – z is a *right-inverse* of x – and $y * x = I$ – y is a *left-inverse* of x –, so, actually, we have proved that all left-inverses are equal to all right-inverses.

5.3 Groups

Algebraically, life becomes really interesting with *groups*. A group is a monoid in which every element has an inverse (with respect to the binary operator).

5.15 Definition. A *group* is a monoid $(V, *, I)$ satisfying, for all $x \in V$:

$$(34) \quad (\exists y : y \in V : x * y = I \wedge y * x = I) .$$

An *Abelian group* is a group in which, in addition, operator $*$ is commutative.

□

As a matter of fact requirement (34) is stronger than strictly necessary: either of the two conjuncts, $x * y = I$ or $y * x = I$, can be dropped without affecting the definition. It is only for practical reasons, and because we do not wish to destroy the symmetry, that we have included both.

5.16 Lemma. Let $(V, *, I)$ be a monoid satisfying, for all $x \in V$:

$$(35) \quad (\exists y : y \in V : x * y = I) .$$

Then $(V, *, I)$ is a group.

Proof. To prove that $(V, *, I)$ is a group we must prove (34) for all $x \in V$, while using (35) for all $x \in V$. So, let $x \in V$ and let, using (35), element $y \in V$ satisfy $x * y = I$. Now to prove (34) for this particular x it suffices to show that y also satisfies $y * x = I$. Let, using (35) once more but with $x := y$, element $z \in V$ satisfy $y * z = I$. Now we calculate:

$$\begin{aligned}
& x \\
= & \quad \{ I \text{ is identity of } * \} \\
& x * I \\
= & \quad \{ y * z = I \} \\
& x * (y * z) \\
= & \quad \{ * \text{ is associative} \} \\
& (x * y) * z \\
= & \quad \{ x * y = I \} \\
& I * z \\
= & \quad \{ I \text{ is identity of } * \} \\
& z .
\end{aligned}$$

So, we have $x = z$; now z satisfies $y * z = I$, and substituting x for z in this we obtain $y * x = I$, as required.

□

The definition of groups states that every element of the set has an inverse. We have already seen that, if an element has an inverse, this inverse is unique. The inverse of an element, of course, depends on that element. Therefore, from now onwards we denote the inverse of every element $x \in V$ by x^{-1} .

5.17 Definition. Let $(V, *, I)$ be a group. For every $x \in V$ its inverse, x^{-1} , satisfies:

$$x * x^{-1} = I \wedge x^{-1} * x = I .$$

□

In the proof of Lemma 5.16 we have introduced y as the inverse of x , and z as the inverse of y , and then we have proved $z = x$. So, as an additional result, we obtain that the inverse of the inverse of an element is that element itself.

5.18 Lemma. Let $(V, *, I)$ be a group. Every $x \in V$ satisfies: $(x^{-1})^{-1} = x$.

□

5.19 Lemma. Let $(V, *, I)$ be a group. All $x, y \in V$ satisfy: $(x * y)^{-1} = y^{-1} * x^{-1}$.

□

5.20 Examples.

- (a) Let $V = \{ i \}$ and let binary operator $*$ on V be defined by $i * i = i$. Then $(V, *, i)$ is a group; this is the *smallest possible* group.
- (b) $(\mathbb{Z}, +, 0)$ is an (Abelian) group.

- (c) $(\mathbb{Q}^+, *, 1)$ and $(\mathbb{Q} \setminus \{0\}, *, 1)$ are (Abelian) groups.
- (d) For any set V we consider the set of all *bijections* from V to V , here denoted by $V \leftrightarrow V$. Then $(V \leftrightarrow V, \circ, I)$ is a group; it is not Abelian. If V is *finite* the bijections in $V \leftrightarrow V$ are also called *permutations* and the group is called a *permutation group*.
- (e) For a fixed positive natural number n we define operator \oplus , of type $\mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{Z}$, by $x \oplus y = (x+y) \bmod n$, for all $x, y \in \mathbb{Z}$. Then this operator also has type $[0..n) \times [0..n) \rightarrow [0..n)$, and $([0..n), \oplus, 0)$ is an Abelian group.

□

A group $(V, *, I)$ has the property that, for all $a, b \in V$, equations of the shape $x: a * x = b$ can be solved. The solution of such an equation even is unique:

$$\begin{aligned}
 & a * x = b \\
 \Rightarrow & \quad \{ \text{Leibniz} \} \\
 & a^{-1} * (a * x) = a^{-1} * b \\
 \Leftrightarrow & \quad \{ * \text{ is associative} \} \\
 & (a^{-1} * a) * x = a^{-1} * b \\
 \Leftrightarrow & \quad \{ a^{-1} \text{ is } a\text{'s inverse} \} \\
 & I * x = a^{-1} * b \\
 \Leftrightarrow & \quad \{ I \text{ is identity of } * \} \\
 & x = a^{-1} * b ,
 \end{aligned}$$

which shows that every solution to the equation is equal to $a^{-1} * b$. Conversely, it also is easy to show that $a^{-1} * b$ is a solution indeed, because $a * (a^{-1} * b)$ is, indeed, equal to b .

This is the characteristic property of groups: a group is the simplest possible structure in which all equations of the shape $x: a * x = b$ can be solved.

5.21 Lemma. Let $(V, *, I)$ be a group. For all $x \in V$ and $n \in \mathbb{N}$ we have:

$$(x^{-1})^n = (x^n)^{-1}$$

□

5.22 Definition. Let $(V, *, I)$ be a group. For all $x \in V$ and $n \in \mathbb{N}$ we define x^{-n} by:

$$x^{-n} = (x^{-1})^n$$

□

5.23 Lemma. Let $(V, *, I)$ be a group. For every $x \in V$ and for all $m, n \in \mathbb{Z}$ we have:

$$x^{m+n} = x^m * x^n$$

□

5.4 Subgroups

5.24 Definition. Let $(V, *, I)$ be a group and let U be a subset of V . If $(U, *, I)$ is a group this is called a *subgroup* of $(V, *, I)$.

□

To verify that $(U, *, I)$ is a subgroup we do not have to verify that $*$ is associative, that I is the identity element, and that group elements have inverses: these properties remain valid. But, we do have to verify that subset U is *closed* under the group operations, that is, to prove that $(U, *, I)$ is a subgroup we must prove the following three properties:

$$(\forall x, y : x, y \in U : x * y \in U) \quad , \text{ and:}$$

$$I \in U \quad , \text{ and:}$$

$$(\forall x : x \in U : x^{-1} \in U) \quad .$$

5.25 Examples.

(a) In $(\mathbb{Z}, +, 0)$ the subset of the *even* integers, with $+$ and 0 , form a subgroup. More generally, for any natural number n the subset of all *multiples of n* , with $+$ and 0 , form a subgroup.

(b) $(\mathbb{Q}^+, *, 1)$ is a subgroup of $(\mathbb{Q} \setminus \{0\}, *, 1)$.

(c) For any group $(V, *, I)$ and for any fixed $a \in V$ we can define a subset U by $U = \{a^i \mid i \in \mathbb{Z}\}$. Then $(U, *, I)$ is a subgroup of $(V, *, I)$, called the subgroup *generated by a* .

(d) Actually, in $(\mathbb{Z}, +, 0)$ the subgroup of all multiples of n , for some natural n , is the subgroup generated by n .

□

5.26 Definition. A group $(V, *, I)$ is called *cyclic* if V contains an element a , say, such that the subgroup generated by a is the whole group, that is, $\{a^i \mid i \in \mathbb{Z}\} = V$.

□

5.27 Definition. A group $(V, *, I)$ is finite if its set V of elements is finite. For a finite group $(V, *, I)$ the *order* of the group $(V, *, I)$ is $\#V$.

□

5.28 Examples.

(a) Let group $(V, *, I)$ be finite of order N and let this group be cyclic. Then V contains an element a , say, such that $V = \{a^i \mid 0 \leq i < N\}$ and $a^N = I$.

(b) For positive natural n , the group $([0..n), \oplus, 0)$, with \oplus as defined in Example 5.20 (e), has order n . This group is cyclic, as it is generated by 1 .

□

5.5 Cosets and Lagrange's Theorem

5.29 Definition. Let $(V, *, I)$ be a group and let $(U, *, I)$ be a subgroup. Then for every $a \in V$ the *left coset* of a is the subset $\{a * y \mid y \in U\}$ and the *right coset* of a is the subset $\{x * a \mid x \in U\}$. The left and right cosets of a are denoted by $a * U$ and $U * a$, respectively.

□

Notice that U is a (left and right) coset too, because $I * U = U$ and $U * I = U$.

If $(V, *, I)$ is a group with subgroup $(U, *, I)$ and for fixed $a \in V$, we can define a function $\varphi: U \rightarrow V$ by $\varphi(x) = a * x$, for all $x \in U$. Then, the left coset $a * U$ just is the image of U under function φ , that is, in terms of lifted functions, we have $a * U = \varphi(U)$.

5.30 Lemma. Let $(V, *, I)$ be a group and let $(U, *, I)$ be a subgroup. For fixed $a \in V$ the function $\varphi: U \rightarrow a * U$, defined by $\varphi(x) = a * x$, for all $x \in U$, is bijective.

Proof. Because $a * U = \varphi(U)$ function φ is surjective. That φ is injective as well follows from, for all $x, y \in U$:

$$\begin{aligned} & \varphi(x) = \varphi(y) \\ \Leftrightarrow & \quad \{ \text{definition of } \varphi \} \\ & a * x = a * y \\ \Rightarrow & \quad \{ \text{Leibniz} \} \\ & a^{-1} * (a * x) = a^{-1} * (a * y) \\ \Leftrightarrow & \quad \{ * \text{ is associative; definition of inverse; identity element} \} \\ & x = y \quad . \end{aligned}$$

□

All subsets of a finite set are finite as well. Therefore, in a finite group $(V, *, I)$ every subgroup $(U, *, I)$ is finite too, and so are all (left and right) cosets of this subgroup. In this case we arrive at an important consequence of the above lemma.

Corollary: In a finite group $(V, *, I)$ with subgroup $(U, *, I)$ we have, for all $a \in V$, that $\#(a * U) = \#U$ and $\#(U * a) = \#U$. In words: in a finite group with a subgroup all cosets have the same size as the subgroup from which they are derived.

□

Because $I \in U$ we have $a \in a * U$, for all $a \in V$. For $a, b \in V$ one may well wonder how the cosets $a * U$ and $b * U$ are related. By careful analysis we can derive that these cosets either are disjoint or are the same, and it so happens that $a * U = b * U$ if and only if $a^{-1} * b \in U$. This gives rise to the following lemma.

5.31 Lemma. Let $(V, *, I)$ be a group and let $(U, *, I)$ be a subgroup. On V we define a relation \sim by, for all $a, b \in V$: $a \sim b \Leftrightarrow a^{-1} * b \in U$. Then:

- (a) \sim is an equivalence relation;
- (b) The left cosets $a * U$, for all $a \in V$, are the equivalence classes of \sim .

□

Now we are ready for our final theorem, which is due to the famous mathematician Joseph Louis Lagrange.

5.32 Theorem. [Lagrange] The order of every subgroup of a finite group is a divisor of the order of the whole group.

Proof. Let M be the order of the group and let N be the order of the subgroup. The equivalence classes of relation \sim , as defined in Lemma 5.31, form a partitioning of V . Each of these classes – the left cosets – has size N , and because V is finite there are only finitely many such cosets: let K be the number of left cosets. Because these sets are mutually disjoint and because their union equals V we have $M = K * N$, hence N is a divisor of M .

□

5.6 Permutation Groups

5.6.1 Function restriction and extension

In this section we will be studying functions on intervals of the shape $[0..n)$, for positive naturals n . If f is a function on the interval $[0..n)$, then we wish to speak of the *restriction of f* to the interval $[0..m)$, for any naturals $m, n: 1 \leq m \leq n$; this is a function with domain $[0..m)$, and on this domain it has the same values as f . We denote this restriction as $f[m$ – “ f take m ” –.

5.33 Definition. For function f on $[0..n)$ and for m with $1 \leq m \leq n$ the function $f[m$, on $[0..m)$, is defined by:

$$(f[m)(i) = f(i) \text{ , for all } i: 0 \leq i < m \text{ .}$$

□

Property: If f , on $[0..n)$, is injective then so is $f[m$, for all $m, n: 1 \leq m \leq n$.

□

* * *

As a converse to function restriction we also have need of the possibility of *function extension*. If f is a function on $[0..n)$ then we wish to define a new function, on $[0..n+1)$, that coincides with f on $[0..n)$ and for which the value in n equals a pre-specified value v , say. We denote this function as $f \triangleleft v$ – “ f snoc v ” –.

5.34 Definition. For function f on $[0..n)$ and for any value v , the function $f \triangleleft v$, on $[0..n+1)$, is defined by:

$$\begin{aligned}(f \triangleleft v)(i) &= f(i) \text{ , for all } i: 0 \leq i < n \\ (f \triangleleft v)(n) &= v\end{aligned}$$

□

5.35 Lemma. Function f , on $[0..n+1)$, satisfies:

$$f = (f[n) \triangleleft f(n) \text{ ,}$$

and function f , on $[0..n)$, and value v satisfy:

$$(f \triangleleft v)[n = f \text{ .}$$

□

5.6.2 Continued Compositions

For any finite list fs of functions, all of the same type $V \rightarrow V$, we define the *continued composition* of (the functions in) list fs . Informally, if list fs has length k then the continued composition of fs is:

$$fs_0 \circ fs_1 \circ \cdots \circ fs_{k-1} \text{ .}$$

Formally, the continued composition of (finite) lists of functions can be defined as a function \mathcal{C} , say, such that $\mathcal{C}(fs)$ is the composition of the functions in fs . Function \mathcal{C} can be defined recursively as follows, for all lists fs, gs of functions and for function f , all of type $V \rightarrow V$.

5.36 Definition.

$$\begin{aligned}\mathcal{C}([]) &= I_V \\ \mathcal{C}([f]) &= f \\ (36) \quad \mathcal{C}(fs \uparrow\uparrow gs) &= \mathcal{C}(fs) \circ \mathcal{C}(gs) \text{ .}\end{aligned}$$

□

Notice that we have defined $\mathcal{C}([]) = I_V$ here because I_V is the identity element of function composition: thus, we guarantee that rule (36) also holds if either fs or gs equals $[]$. Also notice that rule (36) is ambiguous: the decomposition of a list of functions as a concatenation $fs \uparrow\uparrow gs$ of two lists fs and gs of functions is not unique, but, fortunately, this is harmless, because of the associativity of function composition: the result will be the same, independently of this decomposition. As a matter of fact, *lists* are the appropriate data structure here, because of this associativity and because function composition is *not* commutative.

5.6.3 Bijections

We recall that a *bijection* on a set V is a function in $V \rightarrow V$ that is both *injective* and *surjective*. Informally, this means that, for every $v \in V$, there is *exactly one* $u \in V$ satisfying $f(u) = v$: that f is injective means that, for every v , there is *at most one* such u , and that f is surjective means that, for every v , there is *at least one* such u . For any given set V , the bijections on V have the following properties:

- The identity function I_V is a bijection on V ;
- If f and g are bijections on V then so is their composition $f \circ g$;
- If f is a bijection on V , then f has an inverse, f^{-1} , and f^{-1} is a bijection on V too.

From this we conclude that the bijections on V , with \circ and I , form a group.

5.6.4 Permutations

For *finite* set V the bijections on V are also called *permutations of V* . In what follows we will restrict our attention to finite sets of a very particular shape, namely, initial segments of the natural numbers. That is, we consider nonempty intervals of the shape $[0..n)$, for $n: 1 \leq n$. In this section we denote the identity permutation on $[0..n)$ as I_n . Notice that there is only *one* permutation of $[0..1)$, namely I_1 .

* * *

We will use \mathcal{P}_n to denote the set of all permutations of $[0..n)$, for $n: 1 \leq n$. So, \mathcal{P}_n is the subset of those functions in $[0..n) \rightarrow [0..n)$ that are bijections. In what follows, the requirement $1 \leq n$ is left implicit. Hence, as permutations are bijections, we now have that $(\mathcal{P}_n, \circ, I_n)$ is a group, for every n .

* * *

Every permutation in \mathcal{P}_n can be represented compactly by enumerating its values in a list of length n . That is, if $s \in \mathcal{P}_n$ then it is represented by the list $[s_0, s_1, \dots, s_{n-2}, s_{n-1}]$. Notice that, because every permutation is a bijection, this list contains each of the naturals $i: 0 \leq i < n$ exactly once. For example $[0, 1, 2, 3]$ is the identity permutation of $[0..4)$, and $[3, 0, 1, 2]$ is the permutation that “rotates the elements of $[0..4)$ one place to the left”.

If $s \in \mathcal{P}_{n+1}$ then s is a permutation of $[0..n+1)$; so, s is injective and, therefore, $s \upharpoonright n$ also is injective on $[0..n)$. Function s also is surjective and if, in addition, $s_n = n$, then $s \upharpoonright n$ also is surjective in $[0..n) \rightarrow [0..n)$. Hence, if (and only if) $s_n = n$ then $s \upharpoonright n$ is a permutation in \mathcal{P}_n as well. This is expressed by the following lemma.

5.37 Lemma. $(\forall s: s \in \mathcal{P}_{n+1}: s_n = n \Rightarrow s \upharpoonright n \in \mathcal{P}_n)$.

□

Conversely, every permutation in \mathcal{P}_n can be extended to a permutation in \mathcal{P}_{n+1} in a simple way, namely by extending the function with value n .

5.38 Lemma. $(\forall s: s \in \mathcal{P}_n : s \triangleleft n \in \mathcal{P}_{n+1})$.

□

corollary: As a result of these lemmata and of Lemma (5.35) we have:

$$(\forall s: s \in \mathcal{P}_{n+1} : s_n = n \Rightarrow s = (s \lceil n) \triangleleft n) \text{ ,}$$

and:

$$(\forall s: s \in \mathcal{P}_n : s = (s \triangleleft n) \lceil n) \text{ .}$$

□

This shows that the subset of those permutations in \mathcal{P}_{n+1} that map n to n is *isomorphic* to \mathcal{P}_n : the functions $(\lceil n)$ and $(\triangleleft n)$ are the bijections from that subset to \mathcal{P}_n and back, respectively.

In what follows, therefore, we will identify the subset of the permutations in \mathcal{P}_{n+1} that map n to n and \mathcal{P}_n , that is, we will leave the application of the bijections implicit. Thus, for every permutation $s \in \mathcal{P}_{n+1}$ with $s_n = n$ will also say that $s \in \mathcal{P}_n$ and, conversely, we consider every permutation in \mathcal{P}_n to be a permutation in \mathcal{P}_{n+1} as well.

5.6.5 Swaps

For $p, q \in [0..n)$ we define the permutation $p \leftrightarrow q$ – “ p swap q ” – as follows:

$$\begin{aligned} (p \leftrightarrow q)(p) &= q \text{ ,} \\ (p \leftrightarrow q)(q) &= p \text{ ,} \\ (p \leftrightarrow q)(i) &= i \text{ , for all } i: i \in [0..n) \wedge i \neq p \wedge i \neq q \text{ .} \end{aligned}$$

So, $p \leftrightarrow q$ is the permutation that interchanges p and q and leaves everything else in place. Obviously, if $p = q$ then $p \leftrightarrow q = I_n$, so, if $p = q$ then $p \leftrightarrow q$ is not a true “swap”, but it is a permutation nevertheless: it equals I_n . We are mainly interested in *proper* swaps, which are swaps $p \leftrightarrow q$ with $p \neq q$. In the literature proper swaps are also known as “transpositions”.

convention: Because of the isomorphism discussed in the previous section we consider $p \leftrightarrow q$ to be a permutation in \mathcal{P}_n , for every n satisfying $p, q \in [0..n)$, without distinguishing these swaps notationally.

□

Property: Swapping p and q is symmetric in p and q , that is: $(p \leftrightarrow q) = (q \leftrightarrow p)$, for all p, q . Usually, we will, however, represent swaps uniquely, by confining ourselves to swaps $(p \leftrightarrow q)$ with $p < q$.

□

5.39 Lemma. Every swap is its own inverse; that is, for all $p, q \in [0..n)$ we have:

$$(p \leftrightarrow q) \circ (p \leftrightarrow q) = I_n .$$

proof: Directly from the definitions of \circ and \leftrightarrow .

□

The following property shows the effect of the composition of a permutation and a swap: the permutation $s \circ (p \leftrightarrow q)$ differs from s only in that s_p and s_q are interchanged.

Property (37): For $s \in \mathcal{P}_n$ and for $p, q \in [0..n)$ we have:

$$\begin{aligned} (s \circ (p \leftrightarrow q))(p) &= s_q , \\ (s \circ (p \leftrightarrow q))(q) &= s_p , \\ (s \circ (p \leftrightarrow q))(i) &= s_i , \text{ for all } i: i \in [0..n) \wedge i \neq p \wedge i \neq q . \end{aligned}$$

proof: Directly from the definitions of \circ and \leftrightarrow .

□

5.40 Lemma. Every permutation is the continued composition of a finite sequence of swaps.

proof: The lemma states that, for every $n: 1 \leq n$ and for every permutation in \mathcal{P}_n , there exists a finite list of swaps, the continued composition of which equals s . We prove this by Mathematical Induction on n .

base: The only permutation in \mathcal{P}_1 is I_1 , and I_1 , being the identity element of function composition, is the continued composition of [] .

step: Let $s \in \mathcal{P}_{n+1}$. Let $p, 0 \leq p < n$, be such that $s_p = n$. Then we have, by property (37), that $(s \circ (p \leftrightarrow n))(n) = n$, from which we conclude, using Lemma (5.37), that $s \circ (p \leftrightarrow n) \in \mathcal{P}_n$. By Induction Hypothesis, let ss be a (finite) list of swaps the continued composition of which equals $s \circ (p \leftrightarrow n)$; so, we have: $\mathcal{C}(ss) = s \circ (p \leftrightarrow n)$. Now we derive:

$$\begin{aligned} &s \\ = &\quad \{ I_n \text{ is identity of } \circ; \text{ Lemma (5.39) } \} \\ &s \circ (p \leftrightarrow n) \circ (p \leftrightarrow n) \\ = &\quad \{ \text{definition of } ss \} \\ &\mathcal{C}(ss) \circ (p \leftrightarrow n) \\ = &\quad \{ \text{definition of } \mathcal{C} \} \\ &\mathcal{C}(ss) \circ \mathcal{C}([(p \leftrightarrow n)]) \\ = &\quad \{ \text{definition of } \mathcal{C} \} \\ &\mathcal{C}(ss ++ [(p \leftrightarrow n)]) , \end{aligned}$$

from which we conclude that permutation s is the continued composition of the list

$ss \uparrow [(p \leftrightarrow n)]$ of swaps.

□

The proof of this lemma also provides some information on the *length* of the list of swaps. The permutation I_1 is the continued composition of $[\]$, which has length 0, that is, $1-1$. If we now, by Induction Hypothesis, assume that list ss has length $n-1$, then list $ss \uparrow [(p \leftrightarrow n)]$ has length $(n+1)-1$. Thus, we conclude that every permutation in \mathcal{P}_n is the composition of $n-1$ swaps.

Notice that, in the above proof, we have *not* distinguished the cases $p=n$ and $p \neq n$, as this is unnecessary: the given proof is valid for either case. If, however, $p=n$ then $(p \leftrightarrow n)$ equals the identity and can, therefore, be omitted. As a result, we conclude that every permutation in \mathcal{P}_n is the composition of *at most* $n-1$ swaps.

Finally, we note that the representation of a permutation by a list of swaps is not *unique*: every finite list of swaps represents some permutation, and one and the same permutation may be represented by very many different lists of swaps. In the following section, however, we will prove quite a surprising result: the permutations can be partitioned into two classes, which we will call “even permutations” and “odd permutations”, and every swap changes the class of the permutation; that is, the composition of a permutation and a swap always is in the other class than the original permutation. As a consequence, every permutation is even if and only if it is the composition of an even number of swaps, *independently* of the actual composition!

5.6.6 Neighbour swaps

A special case of swaps are the, so-called, “neighbour swaps”, which are swaps of the form $(p \leftrightarrow (p+1))$. As we will see, these provide useful stepping stones in the analysis of even and odd permutations, in the next subsection.

Just as every permutation can be composed from swaps, every swap, in turn, can be composed from neighbour swaps. To prove this we need the following lemma first.

5.41 Lemma. For every p, q with $0 \leq p < q$ we have:

$$(p \leftrightarrow (q+1)) = (q \leftrightarrow (q+1)) \circ (p \leftrightarrow q) \circ (q \leftrightarrow (q+1)) \ .$$

proof: We prove this by showing that $(p \leftrightarrow (q+1))(i)$ is equal to $((q \leftrightarrow (q+1)) \circ (p \leftrightarrow q) \circ (q \leftrightarrow (q+1)))(i)$, for all $i: 0 \leq i$. This requires distinction of 4 cases: $i = p$, $i = q$, $i = q+1$, and all other values of i . We illustrate this for the case $i = q+1$; the other cases can be verified similarly:

$$\begin{aligned} & ((q \leftrightarrow (q+1)) \circ (p \leftrightarrow q) \circ (q \leftrightarrow (q+1)))(q+1) \\ = & \quad \{ \text{Property (37)} \} \\ & ((q \leftrightarrow (q+1)) \circ (p \leftrightarrow q))(q) \\ = & \quad \{ \text{Property (37)} \} \\ & ((q \leftrightarrow (q+1)))(p) \\ = & \quad \{ \text{definition of } \leftrightarrow, \text{ using } p < q, \text{ so } p \neq q \text{ and } p \neq q+1 \} \end{aligned}$$

$$\begin{aligned}
& p \\
= & \{ \text{definition of } \leftrightarrow \} \\
& ((p \leftrightarrow (q+1))(q+1)) .
\end{aligned}$$

□

This lemma shows that the swap $(p \leftrightarrow (q+1))$ can be defined recursively as the composition of $(p \leftrightarrow q)$ and 2 neighbour swaps $(q \leftrightarrow (q+1))$. As a result we obtain the following lemma, which turns out useful in the next subsection.

5.42 Lemma. Every swap $(p \leftrightarrow q)$, for p, q with $0 \leq p < q$, is the continued composition of exactly $2 * k + 1$ neighbour swaps, where $k = q - 1 - p$. Notice that the number $2 * k + 1$ is *odd*.

proof: We prove this by Mathematical Induction on the value of k . The basis of the induction, of course, is the swap $(p \leftrightarrow (p+1))$, so $k = 0$, which all by itself is a *single* neighbour swap. For p, q with $0 \leq p < q$, the swap $(p \leftrightarrow (q+1))$ is, by Lemma (5.41), the composition of $(p \leftrightarrow q)$ and 2 neighbour swaps $(q \leftrightarrow (q+1))$. Hence, if, by Induction Hypothesis, $(p \leftrightarrow q)$ is the composition of $2 * k + 1$ neighbour swaps then $(p \leftrightarrow (q+1))$ is the composition of $2 * (k+1) + 1$ neighbour swaps.

□

Esthetical aside: Because $(q \leftrightarrow (q+1)) = ((q+1) \leftrightarrow q)$, the formula in Lemma (5.41) can be rendered in a more *symmetric* way as:

$$(p \leftrightarrow (q+1)) = ((q+1) \leftrightarrow q) \circ (p \leftrightarrow q) \circ (q \leftrightarrow (q+1)) .$$

For p, q with $0 \leq p < q$, Lemma (5.42) can now be rendered, informally, as:

$$\begin{aligned}
(p \leftrightarrow q) = & (q \leftrightarrow (q-1)) \circ ((q-1) \leftrightarrow (q-2)) \circ \dots \circ ((p+2) \leftrightarrow (p+1)) \circ \\
& (p \leftrightarrow (p+1)) \circ \\
& ((p+1) \leftrightarrow (p+2)) \circ \dots \circ ((q-2) \leftrightarrow (q-1)) \circ ((q-1) \leftrightarrow q) .
\end{aligned}$$

□

5.6.7 Even and odd permutations

A quantity that proves to be relevant in relation to permutations is the, so-called, number of inversions of a permutation. For any permutation $s \in \mathcal{P}_n$, an *inversion* is a pair (i, j) of indices $i, j \in [0..n)$ with the property:

$$i < j \wedge s_j < s_i .$$

We now investigate *the number of inversions*, which is, for $s \in \mathcal{P}_n$, the quantity:

$$(\#i, j : 0 \leq i < j < n : s_j < s_i) .$$

For example, for the identity I_n the number of inversions is 0, whereas for s defined by $s_i = n-1-i$, that is, $s = [n-1, \dots, 1, 0]$, the number of inversions is *maximal*: it equals the number of *all* pairs i, j satisfying $0 \leq i < j < n$ which equals $n * (n-1) / 2$. For every $s \in \mathcal{P}_n$ its number of inversions is a natural number from and including 0, and upto and including $n * (n-1) / 2$.

* * *

Now let $s \in \mathcal{P}_n$ and let $p, q \in [0..n)$ such that $p < q$. We now investigate how the number of inversions in $s \circ (p \leftrightarrow q)$ depends on the number of inversions in s . This analysis can be carried out in terms of $(p \leftrightarrow q)$ directly, but this brings about a rather elaborate case analysis. We have seen that every swap is the continued composition of neighbour swaps, and it so happens that investigating how neighbour swaps affect the number of inversions in a permutation simplifies matters considerably.

So, firstly, we investigate how the number of inversions in $s \circ (q \leftrightarrow (q+1))$ depends on the number of inversions in s . Recall that, by property (37), we have that $(s \circ (q \leftrightarrow (q+1)))(q) = s_{q+1}$, $(s \circ (q \leftrightarrow (q+1)))(q+1) = s_q$, and $(s \circ (q \leftrightarrow (q+1)))(i) = s_i$, for all other indices i . The only pair of indices that is relevant here is the pair $(q, q+1)$: for all *other* pairs (i, j) of indices we have that the swap $(q \leftrightarrow (q+1))$ does not affect the inversions; that is, the pair (i, j) is an inversion in $s \circ (q \leftrightarrow (q+1))$ if and only if it is an inversion in s , for all other pairs (i, j) .

Now, if $s_{q+1} < s_q$ then the pair $(q, q+1)$ is an inversion in s and it is not an inversion in $s \circ (q \leftrightarrow (q+1))$. So, in this case the number of inversions in $s \circ (q \leftrightarrow (q+1))$ equals the number of inversions in s *minus one*. Conversely, if $s_q < s_{q+1}$ then the pair $(q, q+1)$ is not an inversion in s but it is an inversion in $s \circ (q \leftrightarrow (q+1))$. So, in this case the number of inversions in $s \circ (q \leftrightarrow (q+1))$ equals the number of inversions in s *plus one*. Hence, in either case, the number of inversions in $s \circ (q \leftrightarrow (q+1))$ differs from the number of inversions in s by either $+1$ or -1 .

Secondly, by Lemma (5.42), every swap is the composition of an *odd* number of neighbour swaps. Hence, the number of inversions in $s \circ (p \leftrightarrow q)$ differs from the number of inversions in s by an odd amount of individual contributions, each of which is either $+1$ or -1 ; the net result of this – see Exercise 5.7.19 – is *odd*. Isn't that odd?

* * *

We are now ready to harvest the results of the above labour. To start with, we observe that the actual number of inversions in a permutation does not give much information, but this number being even or odd does. This we call the *parity* of a permutation.

5.43 Definition. For $s \in \mathcal{P}_n$ the *parity* of s is:

$$(\#i, j : 0 \leq i < j < n : s_j < s_i) \bmod 2 \ .$$

□

If the parity of a permutation equals 0 we also call the permutation “even” and if its parity equals 1 we also call the permutation “odd”. Now the above analysis boils down to the following lemma.

5.44 Lemma. Composition of a permutation with a proper swap changes its parity.

□

By “repeated application” – that is, of course, by Mathematical Induction – of this lemma we obtain the following theorem.

5.45 Theorem. If a permutation equals the continued composition of a sequence of proper swaps, then its parity equals the parity of the number of swaps.

proof: By Mathematical Induction on the length of the sequences, using the previous lemma.

□

Notice that this theorem pertains to *all possible* ways in which a given permutation equals the continued composition of a sequence of proper swaps. As a consequence, if two different such sequences represent the same permutation, then their lengths have equal parities. So, an even permutation can only be composed from an even number of swaps and an odd permutation can only be composed from an odd number of swaps, independently of *how* the permutation is composed from swaps.

Having the same parity is an equivalence relation. This relation partitions \mathcal{P}_n into two equivalence classes, containing the even and the odd permutations in \mathcal{P}_n , respectively. Recalling that $(\mathcal{P}_n, \circ, I_n)$ is a group, we now also obtain the following additional result.

5.46 Theorem. In the group $(\mathcal{P}_n, \circ, I_n)$ the subset of the even permutations, with \circ and I_n , form a subgroup of $(\mathcal{P}_n, \circ, I_n)$. This means that:

- a) I_n is even;
- b) if s and t are even then so is $s \circ t$;
- c) if s is even then so is s^{-1} .

proof: Left as an exercise.

□

5.7 Exercises

1. Prove Lemma 5.12.
2. Prove that every group $(V, *, I)$ satisfies: $I^{-1} = I$.
3. Describe all groups with exactly 2 elements, with 3 elements, and with 4 elements.

4. Why is $(\mathbb{N}, +, 0)$ not a group?
5. Prove Lemma 5.21.
6. Prove that in every group $(V, *, I)$ we have, for all $x \in V$ and $n \in \mathbb{N}$:
 $x^n = I \Leftrightarrow x^{-n} = I$.
7. Prove Lemma 5.23.
8. For a fixed natural number p , $2 \leq p$, we define operator \otimes , of type $\mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{Z}$, by $x \otimes y = (x * y) \bmod p$, for all $x, y \in \mathbb{Z}$. Prove that:
 - (a) $([0..p], \otimes, 1)$ is a monoid;
 - (b) $([0..p], \otimes, 1)$ is not a group;
 - (c) $([1..p], \otimes, 1)$ is a group if and only if p is a prime number.
9. What is, in a group $(V, *, I)$, the subgroup generated by I ?
10. What is, in the group $(\mathbb{Q} \setminus \{0\}, *, 1)$, the subgroup generated by 2?
11. Why is $(\mathbb{Q}^-, *, 1)$ *not* a subgroup of $(\mathbb{Q} \setminus \{0\}, *, 1)$?
12. Show that $(\mathbb{Z}, +, 0)$ is cyclic.
13. We consider a group $(V, *, I)$. Prove that, for every non-empty subset $U \subseteq V$, the structure $(U, *, I)$ is a subgroup of $(V, *, I)$ if and only if:

$$(\forall x, y : x, y \in U : x * y^{-1} \in U) \text{ .}$$
14. We consider, for some positive natural n , the group $([0..n], \oplus, 0)$, with \oplus as defined in Example 5.20 (e).
 - (a) Prove that, for all positive natural m , the subgroup generated by m is the whole group if and only if m and n are relatively prime.
 - (b) Identify all subgroups of this group.
15. We consider a monoid $(M, *, I)$ with exactly 8 elements. M contains an element g , say, of order 7. This means that $g^7 = I$ and that $g^i \neq I$, for all $i : 0 \leq i < 7$.
 - (a) Prove that M contains an element h , say, that is not a power of g .
 - (b) Prove that such h satisfies $h * g = h$.
16. Prove Lemma 5.31.
17. Identify all subgroups of $([1..p], \otimes, 1)$, for $p \in \{2, 5, 7, 11\}$, and where operator \otimes is defined as in Exercise 8.

18. Give an example showing that composition of swaps usually is not commutative. That is, give example values for k, l, p, q such that:
$$(k \leftrightarrow l) \circ (p \leftrightarrow q) \neq (p \leftrightarrow q) \circ (k \leftrightarrow l) .$$
19. Derive under what condition, to be imposed on k, l, p , and q , we do have that:
$$(k \leftrightarrow l) \circ (p \leftrightarrow q) = (p \leftrightarrow q) \circ (k \leftrightarrow l) .$$
20. Prove Theorem (5.46) .
21. Prove, formally but as elegantly as possible, that the sum of any odd amount of odd numbers is odd.
22. We consider (in-situ) sorting algorithms in which the *only* primitive operations used to modify the array to be sorted are swaps of the shape $p \leftrightarrow (p+1)$. Prove that for every algorithm of this kind the (worst-case) time complexity is not better than $\mathcal{O}(n^2)$.

6 Combinatorics: the art of counting

6.1 Introduction

Combinatorics is the branch of Mathematics in which methods to solve counting problems are studied. Here is a list of questions that are considered combinatorial problems. As we will see, some of these problems in this list may *look* different but actually happen to be (instances of) the *same* problem.

1. What is the number of sequences of length n and constructed from the symbols 0 and 1 only, for any given natural n ?
2. What is the number of sequences of length n and constructed from the symbols 0 and 1 only, but containing no two 0's in succession, for any given natural n ?
3. What is the number of subsets of a given finite set?
4. What is the number of elements of the cartesian product of two given finite sets?
5. What is the number of possible Hungarian licence plates for cars? (A Hungarian licence plate contains three letters followed by three digits.)
6. What is the number of relations on a given finite set?
7. What is the number of sequences of length n and containing each of the numbers $0, 1, \dots, n-1$ (exactly) once, for given $n: 1 \leq n$?
8. What is the number of sequences of length n and containing different objects chosen from a given collection of m different objects, for given $n, m: 1 \leq n \leq m$?
9. What is the number of ways to select n objects from a given collection of m objects, for given $n, m: 1 \leq n \leq m$?
10. What is the number of numbers, in the decimal system, the digits of which are all different?
11. What is the number of "words" consisting of 5 letters?
12. What is the number of ways in which 5 pairs can be formed from a group of 10 persons?
13. What is the number of $n \times n$ matrices, all elements of which are 0 or 1? How many of these matrices have an odd determinant?
14. What is the number of steps needed to sort a given finite sequence?
15. What is the minimal number of "yes/no"-questions needed to determine which one out of a given (finite) set of possibilities is actually the case?

16. In a digital signalling system 8 wires are used, each of which may or may not carry a voltage. For the transmission of a, so-called, “code word” via these wires exactly 4 wires carry a voltage (and the other 4 wires carry no voltage). How many code words are thus possible?

To provide some flavour of the theory needed to answer these questions in a systematic way, we discuss some of these questions here, not necessarily to solve them already but, at least, to shed some light on them.

question 1:

What is the number of sequences of length n and constructed from the symbols 0 and 1 only, for any given natural n ? Well, a sequence of n symbols has n positions, each of which may contain a 0 or a 1. So, each position admits 2 possibilities, and the choice at each position is *independent* from the choice at every other position. For $n=2$, for instance, we have a total of $2*2$ choices, whereas for $n=3$ we have a total of $2*2*2$ possibilities. Generally, for arbitrary n the number of possibilities is 2^n .

This is about the simplest possible counting problem. A systematic way to solve it, which is also applicable to more difficult problems, is by means of induction on n . The only sequence of length 0 is the “empty” sequence, of which there is only one. So, for $n=0$ the number of sequences equals 1. (And, if one wishes to avoid the notion of the empty sequence, one starts with $n=1$ and finds the answer 2, because there are only 2 sequences of length 1, namely consisting of a single symbol 0 or 1.) Next, a sequence of length $n+1$ can be viewed as the *extension* of a sequence of length n with an additional symbol, being 0 or 1. So, such an extension is possible in two ways, each yielding a new sequence. Hence, the number of sequences of length $n+1$ is twice the number of sequences of length n . Thus, we also arrive at the conclusion that the number of sequences of length n equals 2^n : for $n=0$, the number is 2^0 , and each extension doubles the answer, with $2^n * 2 = 2^{n+1}$.

question 2:

What is the number of sequences of length n and constructed from the symbols 0 and 1 only, but containing no two 0’s in succession, for any given natural n ? This question is harder than the previous one, because of the restriction that the sequences contain no two 0’s in succession: now different positions in the sequence are not independent anymore, and the question is not as easily answered as the previous one. As an abbreviation, we call the sequences constructed from the symbols 0 and 1 only, but containing no two 0’s in succession, “admissible”.

To answer this question it pays to introduce a *variable* – a name – for the answer. Because the answer depends on variable n , the parameter of the problem, we let this variable depend on n ; that is, we make it a function on \mathbb{N} , as $n \in \mathbb{N}$. So, we say: let a_i be the number of admissible sequences of length i , for all $i \in \mathbb{N}$. Then a is the function, with type $\mathbb{N} \rightarrow \mathbb{N}$. This enables us to try to formulate an equation for a which, hopefully, we can solve.

For our current problem we reason as follows. The one and only sequence of length 0 is the empty sequence, and it is admissible: as it contains no 0's at all, it certainly contains no two 0's in succession. Thus, we decide that $a_0 = 1$. Next, for $i \in \mathbb{N}$, a sequence of length $i+1$ can now be obtained in two different ways: either by extending an admissible sequence of length i with a symbol 1, always resulting in an admissible sequence, or by extending an admissible sequence of length i with a symbol 0; this, however, yields an admissible sequence *only if* the sequence thus extended does not end with a 0 itself: if it did we would obtain two 0's in succession!

Hence, a new kind of sequences enter the game, namely "admissible sequences not ending with a symbol 0". Therefore, we introduce yet another name b , say: let b_i be the number of admissible sequences not ending with a symbol 0, for all $i \in \mathbb{N}$. Using b , we can now complete our equation for a ; we obtain: $a_{i+1} = a_i + b_i$. Notice that, in this formula, a_i is the number of admissible sequences of length $i+1$ obtained by extending an admissible sequence with a symbol 1, and that b_i is the number of admissible sequences of length $i+1$ obtained by extending an admissible sequence, not ending in a 0, with a symbol 0; their sum, then, is the total number of admissible sequences of length $i+1$.

Thus we obtain the following equation for our function a :

$$\begin{aligned} a_0 &= 1 \\ a_{i+1} &= a_i + b_i, \text{ for all } i \in \mathbb{N} \end{aligned}$$

This equation contains our new variable b and to be able to solve the equation we also need a similar equation for b . By means of precisely the same kind of reasoning we decide that the only admissible sequence, not ending in a 0, of length 0 is the empty sequence; hence, $b_0 = 1$. And, the only way to obtain an admissible sequence, not ending in a 0, of length $i+1$ is by extending an admissible sequence of length i with a symbol 1; so, $b_{i+1} = a_i$. Thus we obtain the following equation for our function b :

$$\begin{aligned} b_0 &= 1 \\ b_{i+1} &= a_i, \text{ for all } i \in \mathbb{N} \end{aligned}$$

These two equations can now be combined into one set of equations for a and b together. Actually, they constitute a *recursive definition* for a and b , because, for any given $n \in \mathbb{N}$, they can be used as rewrite rules to calculate the values of a_n and b_n in a finite number of steps. Nevertheless, we also call this an "equation" because it does not give *explicit* expressions for a and b .

$$\begin{aligned} a_0 &= 1 \\ a_{i+1} &= a_i + b_i, \text{ for all } i \in \mathbb{N} \\ b_0 &= 1 \\ b_{i+1} &= a_i, \text{ for all } i \in \mathbb{N} \end{aligned}$$

Because these equations pertain to functions on \mathbb{N} , and because of their recursive nature, they are also called *recurrence relations*. Recurrence relations of this and similar forms can be solved in a systematic way. This is the subject of a separate section.

question 3:

What is the number of subsets of a given finite set? Well, every element may or may not be an element of a subset: for every element of the given set we have a choice out of 2 possibilities, independently of (the choices for) the other elements. Therefore, if the given set has n elements, we have n independent choices out of 2 possibilities; hence, the number of subsets of a set with n elements equals 2^n .

Actually, this is the very same problem as the one in question 1. The elements of any finite set can be ordered into a finite sequence, and now every subset of it can be represented as a sequence of length n consisting of symbols 0 and 1: a 1 in a given position encodes that the corresponding – according to the order chosen – element of the given set is an element of the subset, and a 0 in a given position encodes that the corresponding element of the given set is not an element of the subset. Thus, there is a one-to-one correspondence – that is: a bijection – between the set of all subsets of a given set of size n and the set of all sequences of length n consisting of symbols 0 and 1. Thus, question 3 and question 1 are essentially the same, and so are their answers.

question 4:

What is the number of elements of the cartesian product of two given finite sets? Let V and W be finite sets with m and n elements, respectively. The cartesian product $V \times W$ is the set of all pairs (v, w) with $v \in V$ and $w \in W$. Because v can be chosen out of m elements, and because, independently, w can be chosen out of n elements, the number of possible pairs equals $m \times n$. So, we have: $\#(V \times W) = \#V * \#W$.

Similarly, the number of elements of the cartesian product of *three* finite sets is the product of the *three* sizes of these sets: $\#(U \times V \times W) = \#U * \#V * \#W$, and so on. The following two questions contain applications of this.

question 5:

What is the number of possible Hungarian licence plates for cars? (A Hungarian licence plate contains three letters followed by three digits.) Each letter is chosen from the alphabet of 26 letters, and each digit is one out of 10 decimal digits. Actually, with A for the alphabet and with D for the set of decimal digits, a Hungarian license plate number is an element of the cartesian product $A \times A \times A \times D \times D \times D$; hence, the number of possible combinations equals $\#A * \#A * \#A * \#D * \#D * \#D$, that is: $26^3 * 10^3$.

question 6:

What is the number of relations on a given finite set? A relation on a set V is a subset of the cartesian product $V \times V$, so the number of relations on V equals the number of subsets of $V \times V$. If V has n elements then $V \times V$ has n^2 elements, so the number of relations equals 2^{n^2} .

question 7:

What is the number of sequences of length n and containing each of the numbers $0, 1, \dots, n-1$ (exactly) once? Notice that the requirement “containing each of the numbers $0, 1, \dots, n-1$ (exactly) once” is rather *overspecific*: that it is about the numbers $0, 1, \dots, n-1$ is not very relevant, all that matters is that the sequence contains n given *different* objects. The element in the first position can be chosen out of n possible objects, and this leaves only $n-1$ objects for the remainder of the sequence. Hence, the second element can be chosen out of these $n-1$ objects, and this, in turn, leaves $n-2$ objects for the (next) remainder of the sequence; and so on, until we are left with exactly one object to choose from for the last element of the sequence. Hence, the number of possible such sequences equals $n * (n-1) * (n-2) * \dots * 2 * 1$; this quantity is usually denoted as $n!$ – “ n factorial” –.

Here, too, a recursive characterisation is possible. The number of sequences of length 1, and containing one given object exactly once, equals 1. And, for any $n \in \mathbb{N}$, a sequence of length $n+1$ containing each of $n+1$ different objects exactly once can be constructed, firstly, by choosing one of these objects, which can be done in $n+1$ ways, and, secondly, by constructing a sequence of length n containing each of the remaining n objects exactly once; finally, the sequence thus obtained is extended with the single object chosen in the first step. We conclude that the number of sequences of length $n+1$ equals $n+1$ times the number of sequences of length n . Thus we arrive at this, well-known, recursive definition of $n!$:

$$\begin{aligned} 1! &= 1 \\ (n+1)! &= (n+1) * n! \text{ , for all } n \in \mathbb{N}^+ \end{aligned}$$

question 8:

What is the number of sequences of length n and containing different objects chosen from a given collection of m different objects, for given $n, m: 1 \leq n \leq m$? One way to approach this problem is to observe that the m given objects can be arranged into a sequence in $m!$ ways, as we have seen in the previous question. The first n objects of such a sequence then constitute a sequence of length n and containing different objects chosen from the given collection of m objects. The order of the remaining $m-n$ objects (in the sequence of length m), however, is completely irrelevant: these remaining objects can be ordered in $(m-n)!$ ways, so there are actually $(m-n)!$ sequences of length n , that all begin with the same n objects in the same order, and that only differ in the order of their last $m-n$ elements. So, to count the number of sequences of length n we have to *divide* the number of sequences of length m by $(m-n)!$. Thus we obtain as answer to our question: $m! / (m-n)!$.

question 9:

What is the number of ways to select m objects from a given collection of n objects, for given $m, n: 1 \leq m \leq n$? One may well wonder if, and if so how, this question differs from the previous one. Well, the previous question was about *sequences*, whereas this

question is about *sets*. The question can be rephrased as: What is the number of subsets of size n of a given set of size m , for given $n, m: 1 \leq n \leq m$? Well, every sequence of length n of the kind in the previous question represents such a subset, albeit that the order in which the objects occur in the sequence now is irrelevant. That is, every two sequences, of length n , containing the *same* objects in possibly *different* orders now represent the *same* subset. Actually, for every selection of n objects there are $n!$ sequences containing these objects and all representing the same subset. Hence, the number of ways to form a subset of size m from a given set of size n equals the answer to the previous question divided by $n!$. Thus we obtain: $m! / ((m-n)! * n!)$. This quantity is called a *binomial coefficient* and it is usually denoted as $\binom{m}{n}$, which is defined by:

$$\binom{m}{n} = \frac{m!}{(m-n)! * n!} .$$

Binomial coefficients happen to have interesting properties which we will study later.

question 10:

What is the number of numbers, in the decimal system, the digits of which are all different? Obviously, such a number consists of *at most* 10 digits. It is questionable whether we allow such a number to start with the digit 0: on the one hand, “0123” does represent a number in the decimal system, on the other hand, we usually do not write down such leading zeroes because they are meaningless: “0123” represents the same number as “123”. So, let us *decide* to exclude meaningless leading zeroes. Then, the only number starting with the digit 0 is “0”.

In this case, the problem is most easily solved by distinguishing the possible numbers according to their length. We have already observed that the length of our numbers is at most 10. For numbers consisting of a single digit we have 10 possibilities, as each of the 10 decimal digits is permitted. For numbers of length $n+1$, $1 \leq n \leq 9$, we can choose the first digit in 9 ways, 0 having been excluded as the first digit. Independently of how this first digit has been chosen we still have 9 digits left for the remaining n digits of the number, because now digit 0 is included again. So, for these remaining n digits the question becomes: what is the number of sequences of n digits, chosen from a set of 9 different digits? This question, however, is an instance of problem 8, with $m := 9$; hence, the answer to this question is $9! / (9-n)!$. Thus, we conclude that the number of numbers, in the decimal system, the digits of which are all different, is 10 single-digit numbers, and $9 * 9! / (9-n)!$ numbers consisting of $n+1$ digits, for $1 \leq n \leq 9$.

question 11:

What is the number of “words” consisting of 5 letters? This depends on what one considers a “word”. If we are asking for the number of 5-letter words in a specific language the question is a linguistic question, not a mathematical one. So, let us

decide that a “word” is just an arbitrary sequence consisting of letters from the 26-letter alphabet. Then, for every position in the sequence we have an independent choice out of 26 possibilities, so the number of 5-letter words equals 26^5 .

Notice that this question is very similar to questions 1 and 4.

question 12:

What is the number of ways in which 5 pairs can be formed from a group of 10 persons? This question does not admit a simple, straightforward answer. Therefore, as in the discussion of question 2, we introduce a variable to represent it: let a_i be the number of ways in which i pairs can be formed from a group of $2*i$ persons, for all $i \in \mathbb{N}^+$. The original question, then, amounts to asking for the value of a_5 .

For a group of 2 persons, the answer is quite simple: only one pair can be formed, so $a_1 = 1$. For a group of $2*(i+1)$ persons, which is equal to $2*i+2$, one person can be matched to $2*i+1$ other persons, and the remaining $2*i$ persons can then be arranged into pairs in a_i ways. Thus we obtain: $a_{i+1} = (2*i+1) * a_i$. Combining these results we obtain the following recurrence relations for a :

$$\begin{aligned} a_0 &= 1 \\ a_{i+1} &= (2*i+1) * a_i, \text{ for all } i \in \mathbb{N}^+ \end{aligned}$$

So, informally, we have: $a_i = (2*i-1) * (2*i-3) * \dots * 3 * 1$, and the answer for a group of 10 persons is $9 * 7 * 5 * 3 * 1$. In some textbooks, the expression for a_i is denoted as $(2*i-1)!!$.

question 13:

What is the number of $n \times n$ matrices, all elements of which are 0 or 1? How many of these have an odd determinant? The first question has no relation to linear algebra. An $n \times n$ matrix has n^2 elements, each of which may be 0 or 1. Hence, the number of such matrices equals 2^{n^2} . This answer is the same as the answer to question 6, and this is no coincidence. Why?

The problem posed in the second question does belong to linear algebra. The answer, which we shall not explain here, is that the number of $n \times n$ 0/1-matrices with an odd determinant is given by this nice formula:

$$\left(\prod_{i: 0 \leq i < n} 2^n - 2^i \right) .$$

For $n=3$, for example, we thus find 168 matrices with an odd determinant, and 344, namely: $2^{3^2} - 168$, with an even determinant.

question 16:

In a digital signalling system 8 wires are used, each of which may or may not carry a voltage. For the transmission of a, so-called, “code word” via these wires exactly 4 wires carry a voltage (and the other 4 wires carry no voltage). How many code words are thus possible? Every code words corresponds to a particular selection of 4 wires

from the 8 available wires, so the number of possible code words equals the number of ways to select 4 objects out of a collection of 8. This is exactly question 9, with $n := 4$ and $m := 8$. Hence, the answer is $8!/(4!*4!)$, which equals 70.

6.2 Recurrence Relations

6.2.1 An example

In the introduction we have discussed this question: What is the number of sequences of length n and constructed from the symbols 0 and 1 only, but containing no two 0's in succession, for any given natural n ? For the sake of brevity we have called such sequences “admissible”.

To solve this problem we have introduced two functions, a and b , on \mathbb{N} , with the following interpretation, for all $i \in \mathbb{N}$:

$$\begin{aligned} a_i &= \text{“the number of admissible sequences of length } i \text{ , and:} \\ b_i &= \text{“the number of admissible sequences of length } i, \\ &\quad \text{not ending with a 0” .} \end{aligned}$$

The answer to the above question then is a_n , and we already have formulated the following set of equations for a and b , also called *recurrence relations*:

$$\begin{aligned} (0) \quad a_0 &= 1 \\ (1) \quad a_{i+1} &= a_i + b_i \text{ , for all } i \in \mathbb{N} \\ (2) \quad b_0 &= 1 \\ (3) \quad b_{i+1} &= a_i \text{ , for all } i \in \mathbb{N} \end{aligned}$$

These recurrence relations form equations for *two* unknowns, namely a and b . We can, however, use (2) and (3) to *eliminate* b from equation (1): after all, we can view (2) and (3) as a definition for b in terms of a . Because of the case distinction between b_0 and b_{i+1} we must apply a similar case distinction to (1); that is, we must split (1) into separate equations for a_1 and a_{i+2} . In the equation for a_1 we can now substitute 1 for b_0 (and 1 for a_0), and in the equation for a_{i+2} we can now substitute a_i for b_{i+1} . Thus we obtain a new set of equations for a in which b does not occur anymore:

$$\begin{aligned} (4) \quad a_0 &= 1 \\ (5) \quad a_1 &= 2 \\ (6) \quad a_{i+2} &= a_{i+1} + a_i \text{ , for all } i \in \mathbb{N} \end{aligned}$$

Thus, we have obtained a recurrence relation for a in isolation – that is, without b –; the “old” equations (2) and (3) can now be used to define b in terms of a , if so desired. That is to say, if we are able to derive an explicit – that is, non-recursive – definition for a then (2) and (3) provide an equally explicit definition for b in terms of a . Notice, however, the proviso “if so desired”: we may very well be interested in b too, but our original problem was about a , and we have only introduced b as an additional, auxiliary variable to be able to formulate equation (1).

* * *

Equations (4) through (6) can be used to calculate as many values a_i as we like, preferably in the order of increasing i . For example, the first 7 values are:

i	a_i
0	1
1	2
2	3
3	5
4	8
5	13
6	21

We see that the values a_i , as a function of i , increase rather quickly. This is not so strange: equation (6) is very similar to the following, only slightly different, equation.

$$a_{i+2} = a_{i+1} + a_{i+1} \text{ , for all } i \in \mathbb{N} \text{ ,}$$

which is equivalent to $a_{i+2} = 2 * a_{i+1}$. Now for this equation, together with (4) and (5), it is quite easy to *guess* that the solution might be $a_i = 2^i$, for all $i \in \mathbb{N}$, and, once we have guessed this, it requires only ordinary Mathematical Induction to prove that our guess is correct.

In our case, however, we have to deal with equation (6). Because of the smaller index, i instead of $i+1$, in one of the terms of the right-hand side, we suspect that the solution does not increase as quickly as 2^i , but, perhaps, it still increases exponentially? This idea deserves further investigation. Notice that we have tacitly decided to confine our attention to equation (6), and to ignore, at least for the time being, equations (4) and (5).

6.2.2 The characteristic equation

The last considerations give rise to the idea to investigate solutions of the form $\gamma * \alpha^i$ (for a_i), for some, non-zero, constants γ and α yet to be determined. To investigate this we substitute this expression for a_i in equation (6), such that we can try to calculate solutions for γ and α :

$$\gamma * \alpha^{i+2} = \gamma * \alpha^{i+1} + \gamma * \alpha^i \text{ , for all } i \in \mathbb{N} \text{ .}$$

Firstly, we observe that, unless $\gamma = 0$, which would make the equation useless, this equation does not really depend on γ ; for every $\gamma \neq 0$, this equation is equivalent to the following one:

$$\alpha^{i+2} = \alpha^{i+1} + \alpha^i \text{ , for all } i \in \mathbb{N} \text{ .}$$

Although this equation has to be met, for one-and-the-same α , for *all* natural i , this is not as bad as it may seem; for $\alpha \neq 0$, this equation, in turn, is equivalent to the following one, in which i does not occur anymore:

$$\alpha^2 = \alpha^1 + \alpha^0 .$$

Using that $\alpha^0 = 1$ and bringing all terms to one side of the equality we rewrite the equation into this form:

$$(7) \quad \alpha^2 - \alpha^1 - 1 = 0 .$$

This is called the *characteristic equation* of recurrence relation (6). What we have obtained now is the knowledge that if the solution to (6) is to be of the shape $\gamma * \alpha^i$ then α must satisfy (7).

Equation (7) is a quadratic one with two solutions α_0 and α_1 , say, given by:

$$(8) \quad \alpha_0 = (1 + \sqrt{5})/2 \text{ and: } \alpha_1 = (1 - \sqrt{5})/2 .$$

Apparently, we have two possibilities for α here, so $a_i = \alpha_0^i$ and $a_i = \alpha_1^i$ are both solutions to our original equation (6), but there is more. We recall equation (6):

$$(6) \quad a_{i+2} = a_{i+1} + a_i , \text{ for all } i \in \mathbb{N} .$$

As we will discuss more extensively later, this is a so-called *linear* and *homogeneous* recurrence relation, which has the property that any *linear combination* of solutions is a solution as well. In our case, this means that for all possible constants γ_0 and γ_1 , the definition:

$$(9) \quad a_i = \gamma_0 * \alpha_0^i + \gamma_1 * \alpha_1^i , \text{ for all } i \in \mathbb{N} ,$$

provides a solution to equation (6). In fact, it can be proved that *all* solutions have this shape, so now we have obtained *all* solutions to equation (6).

Recurrence relations (4) through (6), however, which also constitute a recursive definition for a , suggest that the solution should be *unique*. After all, we are able to construct a table containing the values a_i , for increasing i , as we did in the previous section. So, which one out of the infinitely many solutions of the shape given by (9) is the one we are looking for? This means: what should be the values of constants γ_0 and γ_1 such that we obtain the correct solution? To answer this question we have to consider equations (4) and (5) again, the ones we have temporarily ignored:

$$(4) \quad a_0 = 1$$

$$(5) \quad a_1 = 2$$

If we now instantiate (9), with $i := 0$ and $i := 1$, and using that $\alpha^0 = 1$ and $\alpha^1 = \alpha$, for any α , we obtain these two special cases:

$$a_0 = \gamma_0 + \gamma_1 \text{ and: } a_1 = \gamma_0 * \alpha_0 + \gamma_1 * \alpha_1 .$$

By substituting these values for a_0 and a_1 into equations (4) and (5) we obtain the following two new equations, in which γ_0 and γ_1 are the unknowns now:

$$\begin{aligned} \gamma_0 + \gamma_1 &= 1 \\ \gamma_0 * \alpha_0 + \gamma_1 * \alpha_1 &= 2 \end{aligned}$$

With γ_0 and γ_1 as the unknowns, these are just two linear equations that can be solved by standard algebraic means. In this case we obtain:

$$\gamma_0 = (2 - \alpha_1) / (\alpha_0 - \alpha_1) \quad \text{and:} \quad \gamma_1 = (\alpha_0 - 2) / (\alpha_0 - \alpha_1) \quad ,$$

where α_0 and α_1 are given, above, by definition (8). Thus we obtain, for all $i \in \mathbb{N}$:

$$(10) \quad a_i = ((3 + \sqrt{5}) / 2\sqrt{5}) * ((1 + \sqrt{5}) / 2)^i + ((\sqrt{5} - 3) / 2\sqrt{5}) * ((1 - \sqrt{5}) / 2)^i .$$

6.2.3 A strange, but beautiful, phenomenon

We recall that, originally, we have introduced function a to represent a counting problem; that is, we “defined” a_i to be “the number of admissible sequences”. So, a_i is a natural number, for every $i \in \mathbb{N}$. In the mean time we have obtained a solution of the form (10), in which both α_0 , α_1 , and the coefficients γ_0 and γ_1 are defined in terms of true – non-integer, even irrational – real numbers like $\sqrt{5}$. Apparently, however complicated formula (10) is, its value is a natural number nevertheless. Isn’t that strange?

This observation does have some practical consequences. From a mathematical point of view nothing is wrong with making an excursion into \mathbb{R} to solve a problem regarding integers only. From a computational point of view, however, this is awkward: in a discrete computer real numbers can only be represented approximately. Hence, if we wish to use formula (10) to actually calculate the values of a we have to see to it that the approximation errors are kept small enough not to disturb the answers. In order to avoid these problems, we may prefer to keep the calculations in the domain of the integers altogether. That is to say, from a computational point of view, the original recurrence relations (4) through (6) may be more attractive than the excursion into \mathbb{R} via formula (10).

If we are interested in all values a_i for all $i: 0 \leq i \leq n$, for some given natural n , the computation will unavoidably require an amount of time proportional to $\mathcal{O}(n)$, and the recurrence relations fit this very well. Moreover, if we are only interested in a_n , for some given, fixed (an possibly rather large) value of n , programming techniques exist to construct algorithms that allow a_n to be computed in $\mathcal{O}(\log(n))$ time, and entirely within the domain of the integers. We return to this in Subsection 6.2.7.

6.2.4 Linear recurrence relations

The recurrence relation studied in the previous section belongs to a class of recurrence relations known as *linear recurrence relations* with *constant coefficients*. They are called linear because function values, like a_{i+2} , are defined as linear combinations of other function values, like a_{i+1} and a_i in (6). Notice that equation (6) can be written, slightly more explicitly, as:

$$a_{i+2} = 1 * a_{i+1} + 1 * a_i \quad , \text{ for all } i \in \mathbb{N} \quad .$$

Any formula of the shape $c_0 * x_0 + c_1 * x_1 + c_2 * x_2 + \dots$ is called a linear combination of the x s, with the c s being the coefficients. In the case of recurrence

relations, we speak of constant coefficients if they *do not depend* on i . For example, in our example the coefficients, of a_{i+1} and a_i , are 1 and 1, respectively, which, indeed, do not depend on i .

The general shape of a linear recurrence relation with constant coefficients is the following one, in which the c_j are the coefficients, for all $j: 0 \leq j < k$, and in which k is a constant called the *order* of the recurrence relation:

$$(11) \quad a_{i+k} = c_{k-1} * a_{i+k-1} + \dots + c_1 * a_{i+1} + c_0 * a_i, \text{ for all } i \in \mathbb{N} .$$

So, in a k -th order recurrence relation value a_{i+k} is defined recursively in terms of its k direct predecessors, which are $a_{i+k-1}, \dots, a_{i+1}, a_i$, only.

Relation (11) defines a_{i+k} recursively in terms of $a_{i+k-1}, \dots, a_{i+1}, a_i$, as a linear combination, for all $i \in \mathbb{N}$, but it does not define the first k elements of function a ; that is, relation (11) gives no information on the values a_{k-1}, \dots, a_1, a_0 . These values, therefore, may be, and must be, defined separately. So, a complete k -th order recurrence relation consists of relation (11) together with k separate definitions for the values a_i , for all $i: 0 \leq i < k$. These definitions also are known as the *initial conditions* of the recurrence relation.

* * *

Relation (11) is *homogeneous*, which means that if α_i is a solution for a_i , then so is $\gamma * \alpha_i$, for any constant γ . Relation (11) is *linear*, which means that if α_i and β_i both are solutions for a_i , then so is $\alpha_i + \beta_i$. Combining these two observations we conclude that any linear combination of solutions to (11) is a solution as well. Note that this conclusion *only* pertains to equation (11) in isolation, so without regard for the initial conditions. If the initial conditions are taken into account the recurrence relation has only one, unique, solution.

As in the example before, we now investigate to what extent equation (11) admits solutions of the shape α^i for a_i . Notice that, because of the homogeneity, we do not need to incorporate a constant coefficient: as was the case with the solution to the example –Section 6.2.2–, this coefficient will drop out of the equation anyhow. Substitution of α^i for a_i in (11) transforms it into the following equation for α :

$$\alpha^{i+k} = c_{k-1} * \alpha^{i+k-1} + \dots + c_1 * \alpha^{i+1} + c_0 * \alpha^i, \text{ for all } i \in \mathbb{N} .$$

As we are looking for solutions with $\alpha \neq 0$, this equation can be further simplified into this, equivalent, form:

$$(12) \quad \alpha^k - c_{k-1} * \alpha^{k-1} - \dots - c_1 * \alpha - c_0 * 1 = 0 .$$

Now α^i is a solution for a_i in equation (11) if and only if α is a solution to equation (12), which, as before, is called the *characteristic equation* of the recurrence relation. Notice that this is an algebraic equation of the same order as the order of the recurrence relation.

* * *

The simplest case arises when the k -th order characteristic equation has k *different* real roots, α_j , say, for $j: 0 \leq j < k$. Then, any linear combination of the powers of these roots, that is, for all possible coefficients γ_j , with $0 \leq j < k$, function a defined by:

$$(13) \quad a_i = (\sum_{j: 0 \leq j < k} \gamma_j * \alpha_j^i) \quad , \text{ for all } i \in \mathbb{N} \quad ,$$

is a solution to (11). Moreover, not only does this yield just a solution to (11), it can even be proved that *all* solutions are thus characterized.

As stated before, when the k initial conditions, that is, the defining relations for a_i , for $0 \leq i < k$, are taken into account, the solution to the recurrence relation is unique. This means that if the values for a_i , for $0 \leq i < k$, have been given, and once the α_j , for $0 \leq j < k$ haven been solved from (11), definition (13), for all $i: 0 \leq i < k$, is not a definition anymore but a *restriction* on the possible values for γ_j . That is, the set of relations:

$$(\sum_{j: 0 \leq j < k} \gamma_j * \alpha_j^i) = a_i \quad , \text{ for all } i: 0 \leq i < k \quad ,$$

now constitutes a system of k linear equations with unknowns γ_j , with $0 \leq j < k$, from which these can be solved by means of standard linear-equation solving techniques.

* * *

A somewhat more complicated situation arises if the characteristic equation has *multiple* roots. A simple example of this phenomenon is the equation:

$$\alpha^2 - 2.8 * \alpha + 1.96 = 0 \quad ,$$

which can be rewritten as:

$$(\alpha - 1.4)^2 = 0 \quad .$$

This means that both roots α_0 and α_1 are equal to 1.4.

In a general k -th order algebraic equation a root may have a, so-called, multiplicity upto k . It can now be proved that, if a certain root α_p , say, has, multiplicity q , say, then α_p^i , $i * \alpha_p^i$, $i^2 * \alpha_p^i$, \dots , $i^{q-1} * \alpha_p^i$ are q independent solutions for a_i in equation (11). These q different solutions account for the multiplicity, also q , of root α_p . Thus, in total we still obtain k different basis solutions from which linear combinations can be formed to obtain, again, *all* solutions of the recurrence relation.

As a simple example, the characteristic equation:

$$\alpha^2 - 2.8 * \alpha + 1.96 = 0 \quad ,$$

has a single root, 1.4, with multiplicity 2. Hence, all solutions to the (homogeneous) recurrence relation of which this is the characteristic equation are of the form $\gamma_0 * 1.4^i + \gamma_1 * i * 1.4^i$.

* * *

The situation becomes even more complicated if the characteristic equation has less than k , possibly multiple, roots in \mathbb{R} : then the equation still has k roots, but some (or all) of them are complex numbers. A very simple example is the equation:

$$\alpha^2 - \alpha + 1 = 0 .$$

As the determinant, namely $1 - 4$, of this equation is negative, this equation has no roots in \mathbb{R} at all, but it has two complex numbers as its roots: $(1 + \sqrt{-3})/2$ and $(1 - \sqrt{-3})/2$, which are usually written as $-$ with $i = \sqrt{-1}$ -: $(1 + i * \sqrt{3})/2$ and $(1 - i * \sqrt{3})/2$.

These complex roots can be used, in the same way as described earlier, to obtain all solutions to the recurrence relation. Although, definitely, there is quite some mathematical beauty in this, such solutions are not very useful from a practical point of view. In Subsection 6.2.7 we present a more computational way to cope with such situations.

6.2.5 Keeping simple things simple

In the previous subsection we have presented a general technique to solve linear recurrence relations. Fortunately, for sufficiently simple recurrences we can get away with a more direct, albeit ad hoc, approach. In this subsection we present a few of these simple cases.

These cases have in common that we can guess the solutions quite easily, after which a straightforward proof by Mathematical Induction suffices to confirm the correctness of our guess. In each of the following cases we leave these proofs as exercises to the reader.

* * *

The simplest possible recurrence relation is: $a_{i+1} = a_i$, for all $i \in \mathbb{N}$. Its solutions are the *constant* functions, as all solutions will satisfy: $a_i = a_0$, for all $i \in \mathbb{N}$.

* * *

Next we consider the relation: $a_{i+1} = a_i + d$, for all $i \in \mathbb{N}$, where d is some given constant. Although, because of the presence of d , this recurrence relation is not homogeneous, its solution is still simple, as d is the difference between any two successive values of a ; therefore: $a_i = a_0 + i * d$, for all $i \in \mathbb{N}$.

* * *

Similarly, the following recurrence relation, which is a first order linear recurrence, can be solved as easily: $a_{i+1} = c * a_i$, for all $i \in \mathbb{N}$, where c is some given constant. Its solution is: $a_i = c^i * a_0$, for all $i \in \mathbb{N}$.

* * *

Let d be an integer function on \mathbb{N} . We consider the recurrence relation: $a_{i+1} = a_i + d_i$, for all $i \in \mathbb{N}$. Then we obtain as solution: $a_i = a_0 + (\sum j: 0 \leq j < i: d_j)$, for all $i \in \mathbb{N}$. Of course, whether we can derive a (sufficiently simple) expression for the summation $(\sum j: 0 \leq j < i: d_j)$ depends on how d has been defined.

A very simple instance of this pattern is the recurrence: $a_{i+1} = a_i + i$, for all $i \in \mathbb{N}$. Because $(\sum j: 0 \leq j < i: j) = i * (i-1) / 2$, we obtain: $a_i = a_0 + i * (i-1) / 2$, for all $i \in \mathbb{N}$. Similarly, because $(\sum j: 0 \leq j < i: 2^j) = 2^i - 1$, the recurrence relation: $a_{i+1} = a_i + 2^i$, for all $i \in \mathbb{N}$, is solved by: $a_i = a_0 + 2^i - 1$, for all $i \in \mathbb{N}$.

* * *

In the same vein the following recurrence relation can be solved equally easily: $a_{i+1} = a_i * d_i$, for all $i \in \mathbb{N}$. Then we obtain as solution: $a_i = a_0 * (\prod j: 0 \leq j < i: d_j)$, for all $i \in \mathbb{N}$. For instance: $a_{i+1} = a_i * (i+1)$, for all $i \in \mathbb{N}$. Its solution is, of course: $a_i = a_0 * (\prod j: 1 \leq j \leq i: j)$, for all $i \in \mathbb{N}$, which amounts to: $a_i = a_0 * i!$.

6.2.6 Slightly more complicated cases

The following recurrence relation expresses a_{i+k} as a linear combination of $a_{i+k-1}, \dots, a_{i+1}, a_i$ plus an additional constant, d :

$$(14) \quad a_{i+k} = c_{k-1} * a_{i+k-1} + \dots + c_1 * a_{i+1} + c_0 * a_i + d, \text{ for all } i \in \mathbb{N} .$$

This relation is not homogeneous anymore, but we can try to transform it into a homogeneous one in the following way. We introduce a new function b , say, which we couple to a by means of:

$$b_i = a_i - \delta, \text{ for all } i \in \mathbb{N} ,$$

for some, yet to be chosen, constant δ . This means that $a_i = b_i + \delta$ and if we substitute this into equation (14) we obtain, after some massaging in which all occurrences of δ are “collected” – here we introduce $C = (\sum j: 0 \leq j < k: c_j)$ as an abbreviation –

$$b_{i+k} = c_{k-1} * b_{i+k-1} + \dots + c_1 * b_{i+1} + c_0 * b_i + (C-1) * \delta + d, \\ \text{for all } i \in \mathbb{N} .$$

If we now succeed in choosing δ in such a way that the subexpression $(C-1) * \delta + d$ becomes equal to 0, then we have effectively eliminated d from the recurrence relation, because then the relation becomes:

$$b_{i+k} = c_{k-1} * b_{i+k-1} + \dots + c_1 * b_{i+1} + c_0 * b_i, \text{ for all } i \in \mathbb{N} ,$$

which is a homogeneous equation of the standard form. To achieve this we must, of course, define δ as follows:

$$\delta = \frac{d}{(1-C)} ,$$

which is well-defined only if $C \neq 1$: if $C = 1$ this, so-called, *translation trick* does not work.

* * *

The following recurrence relation expresses a_{i+k} as a linear combination of $a_{i+k-1}, \dots, a_{i+1}, a_i$ plus the value of an additional *disturbing* function, b :

$$a_{i+k} = c_{k-1} * a_{i+k-1} + \dots + c_1 * a_{i+1} + c_0 * a_i + b_i , \text{ for all } i \in \mathbb{N} .$$

This case differs from the previous one in that here the value added to the linear combination depends on i , whereas in the previous case it was constant.

To what extent this recurrence relation admits a systematic solution depends very much on the properties of function b . Sometimes, however, we can formulate a linear recurrence relation for b very similar to the one for a . This relation may define b_{i+k} as a linear combination of both $b_{i+k-1}, \dots, b_{i+1}, b_i$ and $a_{i+k-1}, \dots, a_{i+1}, a_i$. Thus we obtain two recurrence relations for a and b together that may be mutually dependent. This is called a *multiple* recurrence relation.

To make such a multiple recurrence complete we also need, of course, initial conditions for b_i , for all $i: 0 \leq i < k$, similarly to the initial conditions for a .

* * *

A special case of the previous example is a recurrence relation of the shape:

$$a_{i+1} = c * a_i + i , \text{ for all } i \in \mathbb{N} .$$

The factor c suggests solutions of the shape c^i ; therefore, to deal with the disturbing term “ $+i$ ” we introduce a new function b , say, which we couple to a by means of:

$$b_i = a_i / c^i , \text{ for all } i \in \mathbb{N} .$$

Notice that this transformation is similar to the translation trick discussed earlier, but now the relation between b and a is division by c^i instead of subtraction. Conversely we have:

$$a_i = b_i * c^i , \text{ for all } i \in \mathbb{N} ,$$

and if we substitute this into the recurrence relation for a we obtain:

$$b_{i+1} = b_i + i / c^{i+1} , \text{ for all } i \in \mathbb{N} .$$

The initial condition for b becomes, of course, $b_0 = a_0$. The solution to the new recurrence relation for b now is – according to one of the earlier cases in this subsection –:

$$b_i = b_0 + (\sum_{j: 0 \leq j < i: j/c^{j+1}}) , \text{ for all } i \in \mathbb{N} .$$

* * *

We recall the recurrence relation from the example discussed in Subsection 6.2.1:

$$\begin{aligned} (0) \quad a_0 &= 1 \\ (1) \quad a_{i+1} &= a_i + b_i , \text{ for all } i \in \mathbb{N} \\ (2) \quad b_0 &= 1 \\ (3) \quad b_{i+1} &= a_i , \text{ for all } i \in \mathbb{N} \end{aligned}$$

In order to be able to formulate it we had to introduce an additional variable, b , with b_i representing the number of admissible sequences, of length i , that satisfy an additional restriction – in our example: not ending with a “0” –. So, this is a multiple recurrence relation too. Because a_{i+1} and b_{i+1} are defined in terms of a_i and b_i only, we call this a *first-order* multiple recurrence relation.

In Subsection 6.2.1 we have solved this recurrence by viewing equations (2) and (3) as definitions for b in terms of a and, using this, subsequently eliminating b from equation (1). The result was the following, *second-order* recurrence relation for a in isolation:

$$\begin{aligned} (4) \quad a_0 &= 1 \\ (5) \quad a_1 &= 2 \\ (6) \quad a_{i+2} &= a_{i+1} + a_i , \text{ for all } i \in \mathbb{N} \end{aligned}$$

In the previous example, above, we also have seen how a multiple recurrence relation can come into existence: to eliminate a disturbing function by formulating recurrences for it.

Generally, every multiple linear recurrence relation with constant coefficients, involving $n+1$ variables, can be transformed into a recurrence relation involving n variables only, by expressing one of the variables in terms of the remaining ones and, subsequently, eliminating it. This transformation increases the order of the recurrence by one. By repeating this transformation as often as needed every recurrence relation, involving $n+1$ variables, can be transformed into a recurrence relation involving a single variable only, at the expense of increasing the order by n .

6.2.7 A more computational approach

In Subsection 6.2.4 we have introduced a general technique for arbitrary linear recurrence relations with constant coefficients. This technique gives rise to a so-called characteristic equation, which is an algebraic equation of the same order as the order of the recurrence relation.

From a practical point of view, however, this technique has several drawbacks. Firstly, the solutions of the characteristic equation generally are real numbers, if they exist, and else they even are complex numbers. Secondly, algebraic equations of order 4 and higher cannot be solved algebraically: for such equations the solutions can only be approximated numerically.

In this subsection, therefore, we present an approach that can be considered a basis for the development of efficient computer algorithms for the solutions of recurrence relations.

* * *

Firstly, as we have already seen – Subsection 6.2.5 – that a first-order linear recurrence relation has the shape: $a_{i+1} = c * a_i$, for all $i \in \mathbb{N}$, where c is some given constant. As we have seen, its solution is: $a_i = c^i * a_0$, for all $i \in \mathbb{N}$.

Actually, even *multiple* first-order linear recurrences can be solved this easily, if only we are prepared to make a little excursion into Linear Algebra. We demonstrate this by means of the same example used earlier, involving only two variables, but the technique is applicable to any number of variables:

$$\begin{aligned} (0) \quad a_0 &= 1 \\ (1) \quad a_{i+1} &= a_i + b_i, \text{ for all } i \in \mathbb{N} \\ (2) \quad b_0 &= 1 \\ (3) \quad b_{i+1} &= a_i, \text{ for all } i \in \mathbb{N} \end{aligned}$$

By rearranging the order of the equations we group the two initial conditions together, and we group the two homogeneous parts together, thus:

$$\begin{aligned} (0) \quad a_0 &= 1 \\ (2) \quad b_0 &= 1 \\ (1) \quad a_{i+1} &= a_i + b_i, \text{ for all } i \in \mathbb{N} \\ (3) \quad b_{i+1} &= a_i, \text{ for all } i \in \mathbb{N} \end{aligned}$$

The crucial observation now is that equations (1) and (3) define both a_{i+1} and b_{i+1} , respectively, as *linear combinations* of a_i and b_i : $a_i + b_i$ equals $1 * a_i + 1 * b_i$ whereas a_i equals $1 * a_i + 0 * b_i$. Hence, the function mapping the pair (a_i, b_i) to the pair of defining expressions for a_{i+1} and b_{i+1} is a *linear mapping*.

In the world of Linear Algebra linear mappings are studied from vector spaces to vector spaces. In this world linear mappings from an n -dimensional space to an m -dimensional space are represented by $n \times m$ matrices, and the application of a linear mapping to a vector is (represented by) the product of the corresponding matrix and the vector.

nice to know: Moreover, composition of two linear mappings amounts to matrix times matrix multiplication. Because n -dimensional vectors are $n \times 1$ matrices, this means that in this world both function application and function composition are represented by the very *same* binary operator, namely matrix multiplication.

□

We now consider a_i and b_i as the components of a (two-dimensional) vector, and we also consider a_{i+1} and b_{i+1} as the components of a vector, and now we can combine equations (1) and (3) into a single linear, and still first-order, recurrence relation:

$$(15) \quad \begin{pmatrix} a_{i+1} \\ b_{i+1} \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} \times \begin{pmatrix} a_i \\ b_i \end{pmatrix}, \text{ for all } i \in \mathbb{N},$$

$$(16) \quad \begin{pmatrix} a_0 \\ b_0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix},$$

where the matrix has been chosen such that the correct linear combinations of a_i and b_i are obtained for a_{i+1} and b_{i+1} , respectively: the matrix contains the coefficients of the required linear combinations.

To abstract from the vectors and matrices we introduce a new function \mathbf{v} , say, of type $\mathbb{N} \rightarrow \mathbb{Z}^2$, to represent the vectors formed by a and b :

$$\mathbf{v}_i = \begin{pmatrix} a_i \\ b_i \end{pmatrix}, \text{ for all } i \in \mathbb{N},$$

and we denote the matrix containing the coefficients of the recurrences for a and b by \mathbf{C} :

$$\mathbf{C} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}.$$

In terms of \mathbf{v} and \mathbf{C} equations (15) and (16) can now be rewritten as follows:

$$(17) \quad \mathbf{v}_{i+1} = \mathbf{C} \times \mathbf{v}_i, \text{ for all } i \in \mathbb{N}, \text{ and:}$$

$$(18) \quad \mathbf{v}_0 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

This is an ordinary first-order, linear recurrence relation, albeit that the type of the values \mathbf{v}_i now is a two-dimensional integer vector – that is: a pair of integers – instead of a single integer. The vector space, however, has all required algebraic properties for the solution, which is:

$$\mathbf{v}_i = \mathbf{C}^i \times \mathbf{v}_0, \text{ for all } i \in \mathbb{N}.$$

* * *

From Subsection 6.2.4 we recall the general form of a k -th order linear recurrence relation with constant coefficients; for the sake of clarity we also introduce names, α_i – not to be confused with the roots of the characteristic equation in Subsection 6.2.4 – , for the k values needed in the initial conditions for a_i , for all $i: 0 \leq i < k$:

$$(19) \quad \begin{aligned} a_0 &= \alpha_0 \\ a_1 &= \alpha_1 \\ &\vdots \\ a_{k-1} &= \alpha_{k-1} \\ a_{i+k} &= c_{k-1} * a_{i+k-1} + \dots + c_1 * a_{i+1} + c_0 * a_i, \text{ for all } i \in \mathbb{N} \end{aligned}$$

If the order, k , of this equation is larger than 1, then we can decrease the order by one, at the expense of the introduction of an additional variable, as follows. We introduce a new function b , say, which we define in terms of a by:

$$b_i = a_{i+1} \text{ , for all } i \in \mathbb{N} \text{ .}$$

Conversely, reading this from right to left, we obtain a new recurrence relation for a :

$$a_{i+1} = b_i \text{ , for all } i \in \mathbb{N} \text{ .}$$

In equation (19) we may now substitute b_i for a_{i+1} , and b_{i+1} for a_{i+2}, \dots , and substitute b_{i+k-1} for a_{i+k} . This transforms equation (19) into the following, equivalent equation:

$$(20) \quad b_{i+k-1} = c_{k-1} * b_{i+k-2} + \dots + c_1 * b_i + c_0 * a_i \text{ , for all } i \in \mathbb{N} \text{ ,}$$

which, together with the initial conditions: $b_{k-2} = \alpha_{k-1}, \dots, b_1 = \alpha_2, b_0 = \alpha_1$ forms a $(k-1)$ -th order linear recurrence relation for b . Notice that equation (20) still contains a_i : the new recurrence becomes a true mutual recurrence relation, in which a and b are mutually dependent. Thus we obtain:

$$\begin{aligned} a_0 &= \alpha_0 \\ b_0 &= \alpha_1 \\ b_1 &= \alpha_2 \\ &\vdots \\ b_{k-2} &= \alpha_{k-1} \\ (21) \quad a_{i+1} &= b_i \text{ , for all } i \in \mathbb{N} \\ (22) \quad b_{i+k-1} &= c_{k-1} * b_{i+k-2} + \dots + c_1 * b_i + c_0 * a_i \text{ , for all } i \in \mathbb{N} \end{aligned}$$

Thus, by introducing an additional variable, b , we have decreased the order of the recurrence relation by 1: equation (21) is a first-order recurrence, and equation (22) is $(k-1)$ -th order recurrence. By repeating this step as often as needed, the order of the recurrence relation can be reduced to 1, but notice that this transforms a k -th recurrence relation for a *single* variable into a *first-order* multiple recurrence for k variables: in each step an additional variable is introduced.

We can now represent the final result as a first-order recurrence for a k -dimensional vector variable, \mathbf{v} , say, of type $\mathbb{N} \rightarrow \mathbb{Z}^k$. Vector \mathbf{v}_i contains a_i and the $k-1$ additional variables b_i, \dots needed to reduce the k -th order recurrence relation to a first-order one:

$$\mathbf{v}_i = \begin{pmatrix} a_i \\ b_i \\ \vdots \end{pmatrix} \text{ , for all } i \in \mathbb{N} \text{ .}$$

The $k \times k$ matrix \mathbf{C} , say, of coefficients then takes the following shape:

$$\mathbf{C} = \begin{pmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ & & & \vdots & & \\ 0 & 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ c_{k-1} & c_{k-2} & c_{k-3} & \cdots & c_1 & c_0 \end{pmatrix} .$$

In terms of \mathbf{v} and \mathbf{C} the original homogeneous equation (20) now is represented as follows:

$$\mathbf{v}_{i+1} = \mathbf{C} \times \mathbf{v}_i , \text{ for all } i \in \mathbb{N} ,$$

whereas the original initial conditions for a_i , for all $i: 0 \leq i < k$, take the shape of a single initial condition for \mathbf{v}_0 :

$$\mathbf{v}_0 = \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_{k-1} \end{pmatrix} .$$

Finally, the solution of this recurrence relation simply is given by:

$$\mathbf{v}_i = \mathbf{C}^i \times \mathbf{v}_0 , \text{ for all } i \in \mathbb{N} ,$$

and this can be used as a starting point for algorithms for the efficient computation of a_n , for any given natural number n . Because \mathbf{C}^n can be calculated in $\mathcal{O}(\log(n))$ matrix multiplications and because multiplication of $k \times k$ integer matrices requires no more than $\mathcal{O}(k^3)$ integer multiplications, the computation of \mathbf{v}_n and, hence, of a_n can be performed in $\mathcal{O}(\log(n) * k^3)$ time.

6.3 Binomial Coefficients

6.3.1 Factorials

We already have seen that the product $n * (n-1) * (n-2) * \cdots * 2 * 1$ is denoted as $n!$ – “ n factorial” –, for any $n \in \mathbb{N}^+$. Even for $n=0$ this product is meaningful: then it is the *empty* product, consisting of 0 factors, for which the best possible definition is the identity element of multiplication, that is: 1. A recursive definition for $n!$ is:

$$\begin{aligned} 0! &= 1 \\ (n+1)! &= (n+1) * n! , \text{ for all } n \in \mathbb{N} \end{aligned}$$

Factorials occur in solutions to many counting problems. In Section 6.1 we already have argued that the number of ways to arrange n given, different objects into a sequence of length n equals $n!$. Here we repeat the argument, in a slightly more precise way. We do this recursively; for this purpose, we introduce a function a on \mathbb{N} with the idea that:

$a_i =$ “the number of sequences of length i , containing i different objects” ,
for all $i \in \mathbb{N}$.

The number of ways to arrange 0 different objects into a sequence of length 0 equals 1, because the only sequence of length 0 is the empty sequence, and it contains 0 objects. So, we conclude: $a_0 = 1$. Next, to arrange $i+1$ different objects into a sequence of length $i+1$ we can, firstly, select one object that will be the first element of the sequence, and, secondly, arrange the remaining i objects into a sequence of length i . The first object can be selected in $i+1$ different ways and, independently, the remaining i objects can be arranged in a_i ways into a sequence of length i . So, we obtain: $a_{i+1} = (i+1) * a_i$. Thus, we obtain as recurrence relation for a :

$$\begin{aligned} a_0 &= 1 \\ a_{i+1} &= (i+1) * a_i \text{ , for all } i \in \mathbb{N} \end{aligned}$$

The solution to this recurrence is $a_i = i!$, for all $i \in \mathbb{N}$.

Notice that, while we are speaking here of “different objects”, their actual nature is irrelevant: all that matters is that they are different. Actually, we even have used this tacitly in the above argument: after we have selected one object from a collection of $i+1$ different objects, the remaining collection, after removal of the object selected, is a collection of i different objects, independently of which object was selected to be removed.

* * *

A slightly more complicated problem we also have discussed in Section 6.1 was: what is the number of ways to select and arrange n different objects from a given collection of m different objects into a sequence of length n , for given $n, m : 0 \leq n \leq m$?

One way to approach this problem is to observe that, as we now know, the m given objects can be arranged into a sequence in $m!$ ways. The first n objects of such a sequence then constitute a sequence of length n that contains n different objects chosen from the given collection of m objects. The order of the last $m-n$ objects (in the sequence of length m), however, is completely irrelevant: these remaining objects can be ordered in $(m-n)!$ ways without affecting the first n objects in the sequence, so there are actually $(m-n)!$ sequences of length n , that all begin with the same n objects in the same order, and that only differ in the order of their last $m-n$ elements. So, to count the number of sequences of length n we have to *divide* the number of sequences of length m by $(m-n)!$. Thus we obtain as solution: $m! / (m-n)!$.

6.3.2 Binomial coefficients

The number of ways to select an (unordered) *collection* –instead of an (ordered) *sequence*– of n different objects from a given collection of m different objects, for $n, m : 0 \leq n \leq m$, can be determined by the following argument. The number of ways to arrange n different objects into a sequence is, as we know, $n!$. If we are only

interested in a collection of such objects, their order is irrelevant, and this means that all $n!$ such arrangements are equivalent. So, to obtain the number of possible collections we must, again, divide the answer to the previous problem by $n!$. Thus, we obtain as solution:

$$\frac{m!}{n! * (m-n)!} .$$

It so happens that this is an important quantity that has many interesting properties. These quantities are called *binomial coefficients*, and we introduce a notation for it, which is pronounced as “ m over n ”:

$$\binom{m}{n} = \frac{m!}{n! * (m-n)!} , \text{ for all } n, m : 0 \leq n \leq m .$$

For example, we consider all sequences of length m consisting of (binary) digits – “bits”, for short – 0 or 1. As every bit in such a sequence is either 0 or 1 – so, one out of two possibilities –, independently of the other bits in the sequence, the total number of bit sequences of length m equals 2^m . If we number the positions in such a sequence from and including 0 and upto and excluding m , a one-to-one correspondence – that is, a bijection – exists between the collection of all subsets of the interval $[0..m)$ and the collection of all bit sequences of length m : number i is in a given subset if and only if the corresponding sequence contains a bit 1 at position i , for all $i \in [0..m)$.

The collection of bit sequences of length m can be *partitioned* into $m+1$ disjoint classes, according to the number of bits 1 in the sequence. That is, we consider the bit sequences, of length m , containing exactly n bits 1, for $n, m : 0 \leq n \leq m$. In the subset \leftrightarrow bit-sequence correspondence every bit sequence containing exactly n bits 1 now corresponds to a subset with exactly n elements.

We already know that the number of subsets containing exactly n elements chosen from a given collection of m elements equals $\binom{m}{n}$. Hence, because of the one-to-one correspondence, the number of bit sequences of length m and containing exactly n bits 1 also equals $\binom{m}{n}$.

Because the bit sequences of length m and containing exactly n bits 1, for all n with $0 \leq n \leq m$, together are all bit sequences of length m , we obtain the following, quite interesting result:

$$(\sum_{n: 0 \leq n \leq m} \binom{m}{n}) = 2^m , \text{ for all } m \in \mathbb{N} .$$

Binomial coefficients satisfy an interesting recurrence relation, which is also known as *Pascal's triangle*. As a basis for the recurrence we have, for all $m \in \mathbb{N}$:

$$\binom{m}{0} = 1 \quad \text{and:} \quad \binom{m}{m} = 1 \quad ,$$

and for all n, m with $1 \leq n \leq m$ we have:

$$\binom{m+1}{n} = \binom{m}{n} + \binom{m}{n-1} .$$

By means of these relations the binomial coefficients can be arranged in an (infinite) table of triangular shape – hence the name “Pascal’s triangle” –, such that, for any $m \in \mathbb{N}$, row m in the table contains the $m+1$ binomial coefficients $\binom{m}{n}$ in the order of increasing n . For instance, here are the first 7 rows of this triangular table:

			1						
			1	1					
			1	2	1				
			1	3	3	1			
			1	4	6	4	1		
			1	5	10	10	5	1	
			1	6	15	20	15	6	1

6.3.3 The Shepherd’s Principle

Sometimes it is difficult to count the things we wish to count directly, and it may be easier to count the elements of a larger set that is so tightly related to our set of interest that we can obtain the count we want by means of a straightforward correction.

This is known as the *Shepherd’s Principle*, because of the following metaphor. To count the number of sheep in a flock of sheep may be difficult: when the sheep stick closely together one observes a single, cluttered mass of wool in which no individual sheep can be distinguished easily. But, knowing – and assuming! – that every sheep has 4 legs, we can count the legs in the flock and conclude that the number of sheep equals this number of legs divided by 4.

As a matter of fact, we have already applied this technique, in Subsection 6.3.1: to count the number of sequences of length n and containing n different objects from a given collection of m different objects – the sheep –, we have actually counted the number of sequences of length m – the legs – and divided this by the number of ways the irrelevant $m-n$ remaining objects can be ordered – the number of legs per sheep –.

* * *

To demonstrate the Shepherd's Principle we discuss a few simple examples. Firstly, the number of anagrams of the word "FLIPJE" just equals the numbers of ways in which the 6 different letters can be arranged into a sequence of length 6. This is simple because the 6 letters are different, and the answer just is $6!$.

Secondly, what is the number of anagrams of the word "SHEPHERD"? As this word contains 8 letters there would be, if all letters would be different, $8!$ ways to arrange them into an anagram. But the letters are not different: the word contains 2 letters "H" and 2 letters "E"; here we assume, of course, that the 2 letters "H" are indistinguishable and also that the 2 letters "E" are indistinguishable. Well, we make them different *temporarily*, by *tagging* them, for example, by means of subscripts: "SH₀E₀PH₁E₁RD". Now all letters are different, and the number of anagrams of this word – the legs – equals $8!$. The number of legs per sheep now is the number of ways in which the tagged letters can be arranged: E₀ and E₁ can be arranged in $2!$ ways and, similarly, H₀ and H₁ can be arranged in $2!$ ways. So, the number of legs per sheep equals $2! * 2!$, and, hence, the number of anagrams of "SHEPHERD" – the number of sheep – equals $8! / (2! * 2!)$.

Thirdly, what is the number of anagrams of the word "HAHAHAH"? Well, this 7-letter word contains 4 letters "H" and 3 letters "A"; so, on account of the Shepherd's Principle, the answer is $7! / (4! * 3!)$, which equals $\binom{7}{4}$. Is this a surprise? No, because the word is formed from only 2 two different letters, "H" and "A"; so, the number of anagrams of "HAHAHAH" is equal to the number of 7-letter words, formed from letters "H" and "A" only and containing exactly 4 letters "H".

As the exact nature of the objects is irrelevant – it really does not matter whether we use "H" and "A" or "1" and "0" –, this last question is the same as the question for the number of 7-bit sequences containing exactly 4 bits 1. Hence, the answers are the same as well.

6.3.4 Newton's binomial formula

In subsection 6.3.2 we have derived the following relation for binomial coefficients:

$$\left(\sum_{n: 0 \leq n \leq m} \binom{m}{n}\right) = 2^m, \text{ for all } m \in \mathbb{N}.$$

Actually, this is an instance of the following, more general relation, for all – integer, rational, real, ... – numbers x, y :

$$(23) \quad (x+y)^m = \left(\sum_{n: 0 \leq n \leq m} \binom{m}{n} * x^n * y^{m-n}\right), \text{ for all } m \in \mathbb{N}.$$

This is known as *Newton's binomial formula*, although it appears to be much older than Newton. Newton, however, has generalized the formula by not restricting m to the naturals but by allowing it even to be a complex number.

We will not elaborate this here; instead we confine ourselves to the observation that there is a connection between the recurrence relation for binomial coefficients we have seen earlier:

$$\binom{m+1}{n} = \binom{m}{n} + \binom{m}{n-1},$$

and relation (23), by means of the recurrence relation for exponentiation when applied to $(x+y)$:

$$(x+y)^{m+1} = (x+y) * (x+y)^m.$$

6.4 A few examples

We conclude this chapter with a few more examples, some of which we have already solved, albeit in a slightly different form.

example 1:

A vase contains m balls which have been numbered $0, 1, \dots, m-1$. From this vase we draw, at random, n balls, for some $n: 0 \leq n \leq m$. What is the number of possible results if the order in which the balls are drawn is relevant?

Notice that, although the balls themselves may not be distinguishable by their shapes, the fact that they have been numbered makes them distinguishable. So, what matters here is that we have a collection of m different objects. The problem at hand now is the same as a question we have already answered earlier: What is the number of sequences of length n and containing different objects chosen from a given collection of m different objects, for given $n, m: 1 \leq n \leq m$? As we have seen, the answer to this question is: $m! / (m-n)!$.

example 2:

A vase contains m balls which have been numbered $0, 1, \dots, m-1$. From this vase we draw, at random, n balls, for some $n: 0 \leq n \leq m$. What is the number of possible results if the order in which the balls are drawn is *not* relevant? We recall that this question is equivalent to the question: What is the number of subsets of size n of a given finite set of size m ? As we have seen, the answer is $\binom{m}{n}$.

We also have seen that an alternative way to answer this question is to apply the Shepherd's Principle: the n balls drawn from the vase can be ordered in $n!$ ways; as this order is irrelevant we have to divide the answer to the previous question by this very $n!$; as a result, we obtain the same answer, of course: $\binom{m}{n}$.

example 3:

A vase contains m balls which have been numbered $0, 1, \dots, m-1$, for some positive natural m . From this vase we draw, at random, n balls, for some $n \in \mathbb{N}$, but now every ball drawn is placed back into the vase immediately, that is, before any next balls are drawn. What is the number of possible results if the order in which the balls are drawn is relevant? (Notice that in this example we do not need the restriction

$n \leq m$.) The result of this game simply is a sequence, of length n , of numbers in the range $[0..m)$, and the numbers in this sequence are mutually independent. Hence, the answer is: m^n .

example 4:

A vase contains m balls which have been numbered $0, 1, \dots, m-1$, for some positive natural m . From this vase we draw, at random, n balls, for some $n \in \mathbb{N}$, but now every ball drawn is placed back into the vase immediately, that is, before any next balls are drawn. What is the number of possible results if the order in which the balls are drawn is *not* relevant? This example is new, and more difficult than the previous ones.

An effective technique to solve this and similar problems is to choose the “right” *representation* of the possible results of the experiment, namely in such a way that these results can be counted. In the current example, we might write down the numbers of the balls drawn, giving rise to sequences, of length n , and containing numbers in the range $[0..m)$. In this case, however, this representation is not *unique*, because the order of these numbers is irrelevant. For instance, the sequence “3011004” is a possible result, but so are the sequences “4011030” and “0001134”; because the order of the numbers in the sequences is irrelevant and because these three sequences contain the same numbers, albeit in different orders, these sequences actually represent only *one* result, which should be counted only once. So, we need a *unique* representation, because only then the number of possible results is equal to the number of possible representatives of these results.

We obtain a unique representation by writing down the numbers of the balls drawn in *ascending order*: then every result corresponds in a unique way to an ascending sequence, of length n , and containing numbers in the range $[0..m)$. For instance, from the three example sequences “3011004”, “4011030”, and “0001134”, the last one is ascending: it uniquely represents the, one-and-only, common result corresponding to these sequences.

So, the answer to our original question equals the number of *ascending sequences*, of length n , and containing numbers in the range $[0..m)$. To be able to count these smoothly, yet another representation happens to be convenient, based on the following observation. The sequence “0001134”, for instance, consists of 3 numbers 0, followed by 2 numbers 1, followed by 0 numbers 2, followed by 1 number 3, followed by 1 number 4, followed by 0 numbers 5, where we have assumed $m=6$ and $n=7$. So, by counting, for every number in the range $[0..m)$, how often it occurs in the ascending sequence we can also represent the result of the experiment. For the sequence “0001134”, for instance, the corresponding sequence of counts is “320110”: it contains m numbers, namely one count for every number in the range $[0..m)$; each count itself is a number in the range $[0..m]$ and the sum of all counts equals n , which was the length of the ascending sequence.

The restriction that the sum of the counts equals n is a bit awkward but can be made more manageable by means of the following *coding trick*: we represent each count by a sequence of as many zeroes as the count – in the *unary* representation,

so to speak – and we separate these sequences of zeroes by means of a digit 1. For instance, again with $m=6$, the counts of the original ascending sequence “0001134” are 3, 2, 0, 1, 1, and 0, which in the zeroes-encoding become “000”, “00”, “”, “0”, “0”, and “”. Concatenated and separated by digits 1 we obtain: “000100110101”. This is a bit sequence containing exactly n digits 0, because now the sum of the counts simply equals the total number of zeroes, and containing exactly $m-1$ digits 1, because we have m counts separated by $m-1$ digits 1. Consequently, the length of this bit sequence is $m-1+n$.

Generally, every ascending sequence of length n and containing numbers in the range $[0..m)$ can be represented *uniquely* by a bit sequence of length $m-1+n$ containing exactly n digits 0 and $m-1$ digits 1. So, also the results of our balls drawing experiment are uniquely represented by these bit sequences. Therefore, the number of all possible results equals the number of these bit sequences, which is:

$$\binom{m-1+n}{n}.$$

afterthought: It may appear that we have solved this problem in a rather ad hoc fashion. This is true as far as the particular representations we have used are concerned, but choosing a suitable representation for a problem at hand is a crucial ingredient of solutions to very many mathematical problems, particularly in the world of discrete structures. What actually constitutes a “suitable” representation is very problem dependent, which we will not further investigate here, but the method deserves to be remembered.

□

6.5 Exercises

1. Let V be a finite set with n elements.
 - (a) What is the number of functions in $V \rightarrow \{0, 1\}$?
 - (b) What is the number of functions in $V \rightarrow \{0, 1, 2\}$?
 - (c) What is the number of functions in $V \rightarrow W$, with $\#W = m$?
2. (a) What is the number of “words” consisting of 5 *different* letters?
 - (b) What is the number of *injective* functions in $\{0, 1, 2, 3, 4\} \rightarrow \{a, b, c, \dots, z\}$?
 - (c) What is the number of *injective* functions in $V \rightarrow W$, with $\#V = n$ and $\#W = m$?
3. Let w_i be the number of sequences, of length i , consisting of letters from $\{a, b, c\}$, in which each two successive letters are different, for all $i \in \mathbb{N}$.
 - (a) Formulate a recurrence relation for w .
 - (b) Derive an explicit definition for w .

4. Let v_i be the number of sequences, of length i , consisting of letters from $\{a, b, c\}$, in which no two letters a occur in direct succession, for all $i \in \mathbb{N}$. It is given that $v_{i+2} = 2 * v_{i+1} + 2 * v_i$, for all $i \in \mathbb{N}$.
- Determine v_0 and v_1 .
 - Values c_0, c_1, α_0 , and α_1 exist such that $v_i = c_0 * \alpha_0^i + c_1 * \alpha_1^i$, for all $i \in \mathbb{N}$. Determine c_0, c_1, α_0 , and α_1 .
5. Let a_i be the number of 0/1-sequences, of length i , not containing isolated zeroes, for all $i \in \mathbb{N}$. Such sequences are called “admissible”. For example, “1” and “0011” are admissible, but “0”, “0110”, and “0100” are not. The empty sequence – with length 0 – is considered admissible.
- Formulate a recurrence relation for a .
 - Solve this recurrence relation.
6. (a) What is the number of 2-letter combinations in which the 2 letters occur in alphabetical order? For example: “kx” is such a combination, whereas “xk” is not.
- (b) What is the number of 3-letter combinations in which the 3 letters occur in alphabetical order?
- (c) What is the number of n -letter combinations in which the n letters occur in alphabetical order, for every natural $n : n \leq 26$?
7. We consider (finite) sequences g , of length n , consisting of natural numbers, with the additional property that $g_i \leq i$, for all $i : 0 \leq i < n$. (Here g_i denotes the element at position i in the sequence, where positions are numbered from 0, as usual.)
- What is the number of this type of sequences, as a function of $n \in \mathbb{N}^+$?
 - Define a bijection between the set of these sequences and the set of all *permutations* of $[0..n)$.
8. This exercise is about strings of beads. The answers to the following questions depend on which strings of beads one considers “the same”: therefore, give this some thought to start with.
- What is the number of ways to thread n different beads onto a string?
 - We have $2 * n$ beads in n different colours; per colour we have exactly 2 beads, which are indistinguishable. What is the number of ways to thread these beads onto a string?
 - The same question, but now for $3 * n$ beads, with 3 indistinguishable beads per colour.
 - What is the number of ways to thread m red and n blue beads onto a string, if, again, beads of the same colour are indistinguishable?

9. The sequence of, so-called, *Fibonacci numbers* is the function F defined recursively by: $F_0 = 0$, $F_1 = 1$, and $F_{i+2} = F_{i+1} + F_i$, for $i \in \mathbb{N}$. The sequence S of, so-called, *partial sums* of F is defined by: $S_n = (\sum_{i: 0 \leq i < n} F_i)$, for all $n \in \mathbb{N}$. Prove that $S_n = F_{n+1} - 1$, for all $n \in \mathbb{N}$.
10. What is the number of anagrams of the word “STRUCTURES”?
11. Prove that $(\sum_{j: 0 \leq j < i} 2^j) = 2^i - 1$, for all $i \in \mathbb{N}$.
12. A given function a on \mathbb{N} satisfies: $a_0 = 0$ and $a_{i+1} = a_i + i$, for all $i \in \mathbb{N}$. Derive an explicit definition for a .
13. A given function a on \mathbb{N} satisfies: $a_0 = 0$ and $a_{i+1} = 2 * a_i + 2^i$, for all $i \in \mathbb{N}$. Derive an explicit definition for a .
14. A given function a on \mathbb{N} satisfies: $a_0 = 0$ and $a_{i+1} = 2 * a_i + 2^{i+1}$, for all $i \in \mathbb{N}$. Derive an explicit definition for a .
15. A given function a on \mathbb{N}^+ satisfies: $a_1 = 1$ and $a_{i+1} = (1 + 1/i) * a_i$, for all $i \in \mathbb{N}^+$. Derive an explicit definition for a .
16. At the beginning of some year John deposits 1000 Euro in a savings account, on which he receives 8% interest at the end of every year. At the beginning of each next year he withdraws 100 Euro. How many years can John maintain this behaviour, that is, after how many years the balance of his account is insufficient to withdraw yet another 100 Euro?
17. Solve the following recurrence relations for a function a ; in each of the cases we have $a_0 = 0$ and $a_1 = 1$, and for all $i \in \mathbb{N}$:
 - (a) $a_{i+2} = 6 * a_{i+1} - 8 * a_i$;
 - (b) $a_{i+2} = 6 * a_{i+1} - 9 * a_i$;
 - (c) $a_{i+2} = 6 * a_{i+1} - 10 * a_i$.
18. Solve the recurrence relation $a_{i+2} = 6 * a_{i+1} + 9 * a_i$, for all $i \in \mathbb{N}$, with $a_0 = 6$ and $a_1 = 9$.
19. The, so-called, Lucas sequence is the function L defined recursively by: $L_0 = 2$, $L_1 = 1$, and $L_{i+2} = L_{i+1} + L_i$, for all $i \in \mathbb{N}$. Determine the first 8 elements of this sequence, and derive an explicit definition for L .
20. Let V be a finite set with m elements. Let n satisfy: $0 \leq n \leq m$. What is the number of functions in $V \rightarrow \{0, 1\}$, with the additional property that $\#\{v \in V \mid f(v) = 1\} = n$?
21. A chess club with 20 members must elect a board consisting of a chairman, a secretary, and a treasurer.
 - (a) What is the number of ways to form a board (from members of the club) if the three positions must be occupied by different persons?

- (b) What is the number of ways if it is permitted that a person occupies several positions?
 - (c) What is the number of ways if the rule is that every person occupies at most 2 positions?
22. What is the number of sequences, of length 27, consisting of exactly 9 symbols 0, 9 symbols 1, and 9 symbols 2?
23. In a two-dimensional plane, a robot walks from gridpoint $(0,0)$ to gridpoint $(12,5)$, by successively taking a unit step, either in the positive x -direction or in the positive y -direction; it never takes a step backwards. What is the number of possible routes the robot can walk?
24. What is the number of 8-digit decimal numbers in which the digits occur in *descending* order only? For example, “77763331” is such a number, whereas “98764511” is not.

7 Number Theory

7.1 Introduction

The basic arithmetic operations we have learned in primary school are addition, subtraction, multiplication, and division on natural numbers. These operations are meaningful, not only for numbers, but also for more general objects, like functions and, in particular, polynomials. Some properties of the arithmetic operations remain valid in these more general structures, whereas other properties lose their validity.

In this chapter we study some of the properties of integer numbers, which are the numbers $\dots, -2, -1, 0, 1, 2, 3, \dots$. The set of these numbers is called \mathbb{Z} ; it has the set \mathbb{N} of the natural numbers $0, 1, 2, 3, \dots$ as a subset. So, \mathbb{N} includes 0. In the earlier days of mathematics 0 was not considered a natural number, but if we “define” the natural numbers as the numbers used for *counting* then 0 is a very natural number: at the moment I write this, for instance, I have 0 coins in my pocket. (As a child I have learned that “zero is nothing”, but this is not true, of course: although I have 0 coins in my pocket, it is not empty: it contains 0 coins but 1 handkerchief and 1 keyring holding 5 keys.) Also, at any given moment my wallet may contain 3 Euros and 0 U.S. dollars. Moreover, 0 has the, very important, algebraic property that it is the identity element of addition: $x+0 = x$, for all $x \in \mathbb{N}$ (and also for x in $\mathbb{Z}, \mathbb{Q}, \mathbb{R}$, of course). In addition, the natural numbers have the property that every natural number is equal to the number of its *predecessors*, where the predecessors of x are all natural numbers less than x : for every $x \in \mathbb{N}$, we have that x is equal to the number of elements of the set $\{y \in \mathbb{N} \mid y < x\}$, and this is also true if $x = 0$.

The set of *positive naturals* is denoted by \mathbb{N}^+ ; it equals \mathbb{N} but without 0, so we have $\mathbb{N}^+ = \mathbb{N} \setminus \{0\}$. When we discuss the *prime numbers* we will have need of the set of all natural numbers that are at least 2, so this is \mathbb{N} without 0 and 1. We will call such numbers “multiples” and denote its set as \mathbb{N}^{+2} . So, $\mathbb{N}^{+2} = \mathbb{N} \setminus \{0, 1\}$.

7.2 Divisibility

We start our subject with an exploration of *divisibility* and its, hopefully well-known, properties.

7.1 Definition. For $a, d \in \mathbb{Z}$ we say that “ a is *divisible by* d ” or, equivalently, that “ a is a *multiple of* d ” or, equivalently, “ d is a *divisor of* a ” if and only if:

$$(\exists q: q \in \mathbb{Z}: a = q * d) \ .$$

□

By this definition, every integer is a divisor of 0, even 0 itself, although $0/0$ is not defined! If, however, $d \neq 0$, then for every $a \in \mathbb{Z}$ the value q , if it exists, for which $a = q * d$ is *unique*, and we write it as a/d , called the “quotient of” a and d . So, note that a/d is well-defined *if and only if* both a is divisible by d and $d \neq 0$.

Other simple properties are that every integer is divisible by 1 and by itself, and, as a consequence, 1 is a divisor of every integer.

The relation “is a divisor of” is often denoted by the (infix) symbol $|$. That is, we write $d|a$ for the proposition “ d is a *divisor of* a ”. Then, as we have seen earlier, $(\mathbb{N}^+, |)$ is a poset: the relation $|$ is reflexive, anti-symmetric, and transitive.

7.2 Lemma. $(\forall a, d: a, d \in \mathbb{N}^+ : d|a \Rightarrow d \leq a)$.

A direct consequence of this is that the set of all (positive) divisors of a positive natural number is *finite*: the set of divisors of $a \in \mathbb{N}^+$ is a subset of the (finite) interval $[1..a]$. Because $1|a$ the set of divisors of a is non-empty, for every $a \in \mathbb{N}^+$.

□

Another important property of divisibility (by d) is that it is invariant under addition of multiples (of d). We call this a *translation property*.

7.3 Lemma. $(\forall a, d, x: a, d \in \mathbb{N}^+ \wedge x \in \mathbb{Z} : d|a \Leftrightarrow d|(a + x * d))$.

□

* * *

We have seen that not every number is divisible by every other number: a/d is not defined for all a, d , even if $d \neq 0$. Division can, however, be defined more generally, if only we allow the possibility of a, so-called, *remainder*. For the sake of this discussion we restrict ourselves to *positive* d , so $d \in \mathbb{N}^+$.

The equation, with $q \in \mathbb{Z}$ as the unknown, $a = q * d$ may not have a solution, but we can weaken the equation in such a way that it has a solution, and then the solution still happens to be unique.

7.4 Theorem. For all $a \in \mathbb{Z}$ and $d \in \mathbb{N}^+$ unique integers q, r exist satisfying:

$$a = q * d + r \wedge 0 \leq r < d .$$

Proof. We prove existence of the solution and its uniqueness separately.

Existence. We distinguish the cases $0 \leq a$ and $a < 0$. For the first case we prove, for all $a \in \mathbb{N}$, existence of a solution by Mathematical Induction on a . Firstly, if $a < d$, then $q=0$ and $r=a$ are a solution. Secondly, if $d \leq a$ then $a-d \in \mathbb{N}$ and, because $1 \leq d$, we have $a-d < a$. Now, we assume, by Induction Hypothesis, that q and r satisfy:

$$a - d = q * d + r \wedge 0 \leq r < d .$$

Then we also have:

$$a = (q+1) * d + r \wedge 0 \leq r < d ,$$

hence, $q+1$ and r are a solution for a and d .

The proof for the case $a < 0$ is very similar, but now by Mathematical Induction on $-a$. Firstly, if $-d \leq a$ then $q=-1$ and $r=a+d$ are a solution. Secondly, if $a < -d$ then we have $a+d < 0$ and $-(a+d) < -a$. So let, again by Induction Hypothesis, q and r satisfy:

$$a + d = q * d + r \wedge 0 \leq r < d .$$

Then we also have:

$$a = (q-1) * d + r \wedge 0 \leq r < d ,$$

hence, now $q-1$ and r are a solution for a and d .

Uniqueness. Assume that q_0 and r_0 satisfy: $a = q_0 * d + r_0 \wedge 0 \leq r_0 < d$, and, similarly, assume that q_1 and r_1 satisfy: $a = q_1 * d + r_1 \wedge 0 \leq r_1 < d$. To prove uniqueness of the solution, then, we must prove $q_0 = q_1$ and $r_0 = r_1$. We now derive:

$$\begin{aligned} & a = q_0 * d + r_0 \wedge a = q_1 * d + r_1 \\ \Rightarrow & \quad \{ \text{transitivity of } = \} \\ & q_0 * d + r_0 = q_1 * d + r_1 \\ \Leftrightarrow & \quad \{ \text{algebra} \} \\ & r_0 - r_1 = (q_1 - q_0) * d , \end{aligned}$$

from which we conclude that $r_0 - r_1$ is a multiple of d . From the restrictions on r_0 and r_1 , in the above equations, however, it follows that $-d < r_0 - r_1 < +d$, and the only multiple of d in this range is 0. So, we conclude that $r_0 - r_1 = 0$, which is equivalent to $r_0 = r_1$. But now we also have $(q_1 - q_0) * d = 0$, which, because $d \neq 0$, is equivalent to $q_0 = q_1$, as required.

□

7.5 Definition. The unique value q mentioned in the theorem is called the “quotient of” a and d , and is denoted as $a \operatorname{div} d$. The unique value r mentioned in the theorem is called the “remainder of” a and d , and is denoted as $a \operatorname{mod} d$. As a result we obtain the following relation for div and mod , which we consider their definition, albeit an implicit one:

$$a = (a \operatorname{div} d) * d + a \operatorname{mod} d \wedge 0 \leq a \operatorname{mod} d < d .$$

□

warning: Most programming languages have operators for quotient and remainder, even for negative values of d . The definitions of these operators not always are consistent with the definition given here. They do, however, always yield values q and r that satisfy $a = q * d + r$, but differences may arise in the additional restrictions imposed upon r . If both a and d are natural, however, so $0 \leq a$ and $1 \leq d$, then the operators for quotient and remainder yield the same values as $a \operatorname{div} d$ and $a \operatorname{mod} d$ as defined here. Be careful, though, in cases where either a or d may be negative. In particular, for negative a or d , the operations in most programming languages do *not* have the translation properties in Lemma 7.7!

□

Operators div and mod are a true generalisation of division, as they have the following properties.

7.6 Lemma. For all $a \in \mathbb{Z}$ and $d \in \mathbb{N}^+$ we have:

$$\begin{aligned} a \bmod d = 0 &\Leftrightarrow d|a, \text{ and} \\ a \bmod d = 0 &\Rightarrow a \operatorname{div} d = a/d. \end{aligned}$$

□

Operators div and \bmod have many other useful properties, such as the following, so-called, *translation properties*. For more properties we refer the reader to the exercises.

7.7 Lemma. For all $a \in \mathbb{Z}$ and $d \in \mathbb{N}^+$, and for all $x \in \mathbb{Z}$ we have:

$$\begin{aligned} (a+x*d) \operatorname{div} d &= a \operatorname{div} d + x, \text{ and} \\ (a+x*d) \bmod d &= a \bmod d \end{aligned}$$

□

7.3 Greatest common divisors

In this section we consider positive natural numbers only. Throughout this chapter we use names a, b, c, d for variables of type \mathbb{N}^+ and variables x, y, z to denote variables of type \mathbb{Z} .

As we have seen already in Lemma 7.2 the set of (positive) divisors of $a \in \mathbb{N}^+$ is non-empty, as it contains 1 and a , and it is finite. We denote this set as $\mathcal{D}(a)$.

7.8 Definition. For $a \in \mathbb{N}^+$ the set $\mathcal{D}(a)$ of (positive) divisors of a is defined by:

$$\mathcal{D}(a) = \{d \in \mathbb{N}^+ \mid d|a\}.$$

□

For all $a, b \in \mathbb{N}^+$ their respective sets $\mathcal{D}(a)$ and $\mathcal{D}(b)$ have a non-empty intersection, because both contain 1, and this intersection is finite as well. The elements of the set $\mathcal{D}(a) \cap \mathcal{D}(b)$ are called *common divisors* of a and b . As an abbreviation we also denote this intersection as $\mathcal{D}(a, b)$. So, by definition $\mathcal{D}(a, b)$ satisfies:

$$\mathcal{D}(a, b) = \{d \in \mathbb{N}^+ \mid d|a \wedge d|b\}.$$

Because $\mathcal{D}(a, b)$ is non-empty and finite it has a maximum. This maximum is called the *greatest common divisor* of a and b . This depends on a and b , of course, so it is a function, which we call gcd .

7.9 Definition. Function gcd , of type $\mathbb{N}^+ \times \mathbb{N}^+ \rightarrow \mathbb{N}^+$, is defined by, for all $a, b \in \mathbb{N}^+$:

$$gcd(a, b) = \max \mathcal{D}(a, b),$$

or, more explicitly, by:

$$gcd(a, b) = (\max d : d \in \mathbb{N}^+ \wedge d|a \wedge d|b : d).$$

□

The common divisors of a and a itself just are the divisors of a , that is, we have $\mathcal{D}(a, a) = \mathcal{D}(a)$; hence the greatest common divisor of a and a itself just is the greatest divisor of a , which is a . If $b < a$ then $a - b \in \mathbb{N}^+$, and on account of translation Lemma 7.3, we conclude that $\mathcal{D}(a, b) = \mathcal{D}(a - b, b)$. In words: if $b < a$ then a and b have the same common divisors as $a - b$ and b ; hence, their greatest common divisors are equal as well. Similarly, if $a < b$ then $\mathcal{D}(a, b) = \mathcal{D}(a, b - a)$ and the greatest common divisor of a and b is equal to the greatest common divisor of b and $b - a$. Thus we obtain the following lemma.

7.10 Lemma. For all $a, b \in \mathbb{N}^+$:

$$\begin{aligned} \gcd(a, a) &= a \\ \gcd(a, b) &= \gcd(a - b, b) \quad , \text{ if } b < a \\ \gcd(a, b) &= \gcd(a, b - a) \quad , \text{ if } a < b \end{aligned}$$

□

Greatest common divisors also have the following, quite surprising, property that the common divisors of a and b are the divisors of $\gcd(a, b)$.

7.11 Lemma. For all $a, b, c \in \mathbb{N}^+$ with $c = \gcd(a, b)$:

$$\mathcal{D}(a, b) = \mathcal{D}(c) .$$

Proof. By Mathematical Induction on the value $a + b$. Firstly, if $a = b$ then, as we have seen, $\mathcal{D}(a, b) = \mathcal{D}(a)$ and, by Lemma 7.10, we have $c = a$, so also $\mathcal{D}(c) = \mathcal{D}(a)$; hence, $\mathcal{D}(a, b) = \mathcal{D}(c)$. Secondly, if $b < a$ then, as we have seen, $\mathcal{D}(a, b) = \mathcal{D}(a - b, b)$, and, by Lemma 7.10, we have $c = \gcd(a - b, b)$; now we assume, by Induction Hypothesis – because $(a - b) + b < a + b$ – that $\mathcal{D}(a - b, b) = \mathcal{D}(c)$; then it also follows that $\mathcal{D}(a, b) = \mathcal{D}(c)$. Thirdly, the case $a < b$ is similar to the previous case, because the situation is symmetric in a and b .

□

A direct consequence of translation Lemma 7.3 is that every common divisor of a and b also is divisor of any linear combination of a and b .

7.12 Lemma. For all $a, b, d \in \mathbb{N}^+$ and for all $x, y \in \mathbb{Z}$:

$$d | a \wedge d | b \Rightarrow d | (x * a + y * b) .$$

□

In particular, $\gcd(a, b)$ is a common divisor of a and b ; hence, $\gcd(a, b)$ also is a divisor of every linear combination of a and b . There is more to this, however, as the following theorem shows.

7.13 Theorem. For all $a, b \in \mathbb{N}^+$, integers $x, y \in \mathbb{Z}$ exist satisfying:

$$\gcd(a, b) = x * a + y * b$$

Proof. A constructive proof is given in the next section, in the form of Euclid's extended algorithm, which shows how suitable numbers x and y can be calculated.

□

A consequence of this theorem is that $\gcd(a, b)$ is the *smallest* of all positive linear combinations of a and b .

7.14 Theorem. For all $a, b \in \mathbb{N}^+$ we have:

$$\gcd(a, b) = \min \{ x*a + y*b \mid x, y \in \mathbb{Z} \wedge 1 \leq x*a + y*b \} .$$

Proof. Let xm and ym be integers for which $xm*a + ym*b$ is positive and minimal. Let $c = \gcd(a, b)$ and let xc and yc be integers for which $c = xc*a + yc*b$; on account of Theorem 7.13 such numbers exist. Now we must prove: $c = xm*a + ym*b$, which we do by proving $c \leq xm*a + ym*b$ and $xm*a + ym*b \leq c$ separately:

$$\begin{aligned} & c \leq xm*a + ym*b \\ \Leftrightarrow & \quad \{ \text{Lemma 7.2, using that both } c \text{ and } xm*a + ym*b \text{ are positive} \} \\ & c \mid (xm*a + ym*b) \\ \Leftrightarrow & \quad \{ \text{Lemma 7.12} \} \\ & c \mid a \wedge c \mid b \\ \Leftrightarrow & \quad \{ c = \gcd(a, b) \} \\ & \text{true} , \end{aligned}$$

and:

$$\begin{aligned} & xm*a + ym*b \leq c \\ \Leftrightarrow & \quad \{ \text{definition of } xc \text{ and } yc \} \\ & xm*a + ym*b \leq xc*a + yc*b \\ \Leftrightarrow & \quad \{ \text{both sides of the inequality are positive, and the LHS is minimal} \} \\ & \text{true} \end{aligned}$$

□

* * *

Numbers of which the greatest common divisor equals 1 are called *relatively prime* or also *co-prime*. As we have seen –Theorem 7.13–, for all $a, b \in \mathbb{N}^+$ integers x, y exist such that

$$\gcd(a, b) = x*a + y*b .$$

If $\gcd(a, b) = 1$ this amounts to the existence of integers x and y satisfying:

$$x*a + y*b = 1 .$$

The following two lemmata are useful consequences of this property.

7.15 Lemma. For all $a, b, c \in \mathbb{N}^+$: $\gcd(a, b) = 1 \wedge a \mid (b * c) \Rightarrow a \mid c$.

Proof. Let $\gcd(a, b) = 1$ and let $a \mid (b * c)$; that is, assume that $x, y, z \in \mathbb{Z}$ satisfy:

$$(24) \quad x * a + y * b = 1$$

$$(25) \quad b * c = z * a$$

Now we derive:

$$\begin{aligned} & \text{true} \\ \Leftrightarrow & \quad \{ (24) \} \\ & x * a + y * b = 1 \\ \Rightarrow & \quad \{ \text{Leibniz} \} \\ & x * a * c + y * b * c = c \\ \Leftrightarrow & \quad \{ (25) \} \\ & x * a * c + y * z * a = c \\ \Leftrightarrow & \quad \{ \text{algebra} \} \\ & (x * c + y * z) * a = c \\ \Rightarrow & \quad \{ \exists\text{-introduction, with } q := x * c + y * z \} \\ & (\exists q : q \in \mathbb{Z} : c = q * a) \\ \Leftrightarrow & \quad \{ \text{Definition of } \mid \} \\ & a \mid c \end{aligned}$$

□

7.16 Lemma. For all $a, b, c \in \mathbb{N}^+$: $\gcd(a, b) = c \Rightarrow \gcd(a/c, b/c) = 1$.

□

7.4 Euclid's algorithm and its extension

The relations in Lemma 7.10 can be considered as a recursive definition of function \gcd ; that, thus, function \gcd is well-defined is, again, proved by Mathematical Induction on the value $a + b$. So, the following recursive definition actually constitutes an algorithm for the computation of the greatest common divisor of two positive naturals. This is known as “Euclid's algorithm”. For all $a, b \in \mathbb{N}^+$:

$$\begin{aligned} \gcd(a, b) = & \text{ if } a = b \rightarrow a \\ & \square a > b \rightarrow \gcd(a - b, b) \\ & \square a < b \rightarrow \gcd(a, b - a) \\ & \text{ fi} \end{aligned}$$

This version of the algorithm is not particularly *efficient*, but it is the simplest possible. If, for instance, a is very much larger than b the calculation of $\text{gcd}(a, b)$ gives rise to the repeated subtraction $a - b$, until a does not exceed b anymore. Therefore, a more efficient algorithm can be constructed by means of `div` and `mod` operations.

* * *

According to Theorem 7.13 we have that $\text{gcd}(a, b)$ is a *linear combination* of a and b ; this means that, for every $a, b \in \mathbb{N}^+$, integers x, y exist satisfying:

$$(26) \quad \text{gcd}(a, b) = x * a + y * b .$$

In what follows we call such integers “matching numbers” for $\text{gcd}(a, b)$. Matching numbers are not *unique*: if, for instance, x and y are matching numbers for $\text{gcd}(a, b)$ then so are $x+b$ and $y-a$.

Because Theorem 7.13 is about *existence* of integers, we can try to prove it constructively by showing how these numbers can be computed. It so happens that Euclid’s algorithm can be *extended* in such a way that, in addition to $\text{gcd}(a, b)$, integers x and y are calculated that satisfy (26) as well. As a result, provided we have proved the correctness of the extended algorithm, we not only have a proof of the theorem but we also obtain an algorithm to compute these numbers. (And, from the point of view of proving the theorem, efficiency is of no concern and the simplest possible algorithm yields the simplest possible proof.)

As was the case with function gcd we present Euclid’s extended algorithm in the form of a recursively defined function. For this purpose we simply call this function F here; it maps a pair of positive naturals to a *triple* consisting of a positive natural and two integers, namely the GCD of the pair together with matching numbers. We denote such a triple as $\langle c, x, y \rangle$, in which c , x , and y are the elements of the triple¹². This means that function F is required to satisfy the following *specification*.

specification: Function F has type $\mathbb{N}^+ \times \mathbb{N}^+ \rightarrow \mathbb{N}^+ \times \mathbb{Z} \times \mathbb{Z}$, and for all $a, b, c \in \mathbb{N}^+$ and for all $x, y \in \mathbb{Z}$, function F satisfies:

$$F(a, b) = \langle c, x, y \rangle \Rightarrow c = \text{gcd}(a, b) \wedge c = x * a + y * b$$

□

Notice that this specification does not specify F uniquely: because, as we have seen, matching numbers are not unique, several different functions F will satisfy this specification. This specification only states, firstly, that for every pair of positive naturals a and b its value $F(a, b)$ is a triple consisting of a positive natural and two integers, and, secondly, that for every such triple its first element is equal to $\text{gcd}(a, b)$ and its second and third elements are matching numbers for $\text{gcd}(a, b)$.

A simple recursive definition for F can now be constructed, based on the following considerations, using Mathematical Induction on $a + b$ again. Firstly, if $a = b$

¹²The datatype of, so-called, *tuples*, of which triples are a special case, is common in functional programming; it corresponds to “records” in PASCAL and “structs” in JAVA.

then $\gcd(a, b) = a$, and $a = 1 * a + 0 * b$: hence, in this case $x=1$ and $y=0$ is an acceptable solution for x and y . Because $a=b$ we also have $\gcd(a, b) = b$; therefore, $x=0$ and $y=1$ is an acceptable solution too: this illustrates once more that the numbers x and y are not unique.

Secondly, if $a > b$ then we have $\gcd(a, b) = \gcd(a-b, b)$. Now suppose, by Induction Hypothesis, that x, y are integers satisfying:

$$\gcd(a-b, b) = x * (a-b) + y * b .$$

The right-hand side of this equality can be rewritten to:

$$\gcd(a-b, b) = x * a + (y-x) * b ,$$

and because $\gcd(a, b) = \gcd(a-b, b)$ this is equivalent to:

$$\gcd(a, b) = x * a + (y-x) * b .$$

From this we conclude that if x and y are matching numbers for $\gcd(a-b, b)$ then x and $y-x$ are matching numbers for $\gcd(a, b)$.

Finally and similarly, for the case $a < b$ we can show that if x and y are matching numbers for $\gcd(a, b-a)$ then $x-y$ and y are matching numbers for $\gcd(a, b)$.

We now combine these results into the following recursive definition for F ; this we call *Euclid's extended algorithm*:

```

F(a, b) = if a = b → ⟨ a, 1, 0 ⟩
          [] a > b → ⟨ c, x, y-x ⟩
                    where ⟨ c, x, y ⟩ = F(a-b, b) end
          [] a < b → ⟨ c, x-y, y ⟩
                    where ⟨ c, x, y ⟩ = F(a, b-a) end
          fi

```

This recursive definition is an example of a, so-called, *functional program*, but it is not difficult to encode this as a recursive function in languages like PASCAL or JAVA. As was the case with Euclid's algorithm proper, this algorithm is not very efficient, but it can be transformed into a more efficient one by means of division and remainder operations.

7.5 Equations and their solutions

As we have seen, if $\gcd(a, b) = 1$, for $a, b \in \mathbb{N}^+$, then integers x, y exist satisfying:

$$(27) \quad x * a + y * b = 1 .$$

Conversely, if x and y satisfy (27) then, because 1 is the smallest of all positive naturals, we also have that 1 is the smallest of all positive linear combinations of a and b . Hence, we conclude, by Theorem 7.14, that if (27) then $\gcd(a, b) = 1$. Thus, we obtain:

$$\gcd(a, b) = 1 \Leftrightarrow (\exists x, y: x, y \in \mathbb{Z}: x * a + y * b = 1) \text{ .}$$

This shows that the equation $x * a + y * b = 1$, for given $a, b \in \mathbb{N}^+$ and with $x, y \in \mathbb{Z}$ the unknowns, has a solution *if and only if* $\gcd(a, b) = 1$. As we have seen, Euclid's extended algorithm can be used to calculate such a solution, but we also have seen that the solution is not unique. So, the question arises whether we can characterize *all* solutions of this equation. The answer is that we can, as the following theorem shows.

7.17 Theorem. For $a, b \in \mathbb{N}^+$ with $\gcd(a, b) = 1$, let $xc, yc \in \mathbb{Z}$ satisfy: $xc * a + yc * b = 1$. Then the set of all pairs $(xc + z * b, yc - z * a)$, for all $z \in \mathbb{Z}$, is the set of all solutions to the equation $x * a + y * b = 1$.

Proof. We prove this by proving, firstly, that every such pair is a solution, and, secondly, that every solution is such a pair. For $z \in \mathbb{Z}$ we derive:

$$\begin{aligned} & (xc + z * b) * a + (yc - z * a) * b \\ = & \quad \{ \text{algebra} \} \\ & xc * a + z * b * a + yc * b - z * a * b \\ = & \quad \{ \text{algebra} \} \\ & xc * a + yc * b \\ = & \quad \{ \text{definition of } xc \text{ and } yc \} \\ & 1 \text{ ,} \end{aligned}$$

which shows that the pair $(xc + z * b, yc - z * a)$ is a solution indeed. Conversely:

$$\begin{aligned} & x * a + y * b = 1 \wedge xc * a + yc * b = 1 \\ \Rightarrow & \quad \{ \text{algebra} \} \\ & (x - xc) * a + (y - yc) * b = 0 \\ \Leftrightarrow & \quad \{ \text{algebra} \} \\ & (x - xc) * a = (yc - y) * b \\ \Rightarrow & \quad \{ b \mid ((x - xc) * a) \text{ and } \gcd(a, b) = 1: \text{Lemma 7.15} \} \\ & b \mid (x - xc) \\ \Leftrightarrow & \quad \{ \text{Definition 7.1, of } \mid \} \\ & (\exists z: z \in \mathbb{Z}: x - xc = z * b) \\ \Leftrightarrow & \quad \{ \text{algebra} \} \\ & (\exists z: z \in \mathbb{Z}: x = xc + z * b) \text{ ,} \end{aligned}$$

which shows that if x and y solve the equation then $x = xc + z * b$ for some $z \in \mathbb{Z}$. Now it is easy to show that, for this *very same* z we also have $y = yc - z * a$. Thus we conclude that all solutions to the equation are of the shape $(xc + z * b, yc - z * a)$, for some $z \in \mathbb{Z}$. Notice that here (xc, yc) may be *any* solution to the equation, which can be obtained, for instance, by means of Euclid's extended algorithm.

□

* * *

We now consider the more general equation $x * a + y * b = d$, for given $a, b, d \in \mathbb{N}^+$ and with $x, y \in \mathbb{Z}$ the unknowns. Because $\gcd(a, b)$ is a divisor of $x * a + y * b$ we conclude that, if the equation has a solution then $\gcd(a, b)$ is a divisor of d too; if not, the equation has no solutions at all.

Conversely, if $\gcd(a, b)$ is a divisor of d then, as we will show, the equation does have solutions and again, we wish to characterize all solutions. So, we assume that $\gcd(a, b)$ is a divisor of d . Now, $\gcd(a, b)$ also is a divisor of a and b ; therefore, with $c = \gcd(a, b)$ we have that the equation $x * a + y * b = d$ is equivalent to:

$$x * (a/c) + y * (b/c) = (d/c) \quad ,$$

and, because of the divisibility, the coefficients a/c , b/c and d/c all are integers. So, by means of the substitutions $a := a/c$, $b := b/c$, and $d := d/c$, the equation can be transformed into this new equation:

$$(28) \quad x * a + y * b = d \quad ,$$

which has the same shape as the original equation but with the additional property, on account of Lemma 7.16, that $\gcd(a, b) = 1$. Moreover, this equation has exactly the same solutions as the original one. By means of Theorem 7.17 we have characterized all solutions to the equation:

$$(29) \quad x * a + y * b = 1 \quad .$$

Now if $x * a + y * b = 1$ then also $(x * d) * a + (y * d) * b = d$, so for such x, y the pair $(x * d, y * d)$ is a solution to equation (28). By means of the same reasoning as in (the second part of) the proof to Theorem 7.17 we can prove that all solutions to equation (28) are thus obtained.

Summarizing, let (x, y) may be any solution to equation (29), then the pairs $(x * d + z * b, y * d - z * a)$, for all $z \in \mathbb{Z}$, constitute all solutions to equation (28).

7.6 The prime numbers

In this section we study the set \mathbb{N}^{+2} of *multiples*, which are the natural numbers from 2 onwards. As we have seen, every integer, and, hence, also every multiple is divisible by 1 and by itself. A multiple with the property that it is *not* divisible by any other number is called a *prime (number)*.

7.18 Definition. A *prime* is a multiple that is divisible by 1 and itself only.

□

If we would not restrict ourselves to multiples but to positive naturals instead, 1 would be a prime too, according to this definition. There are sound, technical reasons, however, not to consider 1 as a prime, which is why we define the primes as a subset of the multiples. So, the smallest prime is 2.

7.19 Example. The primes less than 100 are: 2, 3, 5, 7, 11, 13, 17, 19, 23, 29, 31, 37, 41, 43, 47, 53, 59, 61, 67, 71, 73, 79, 83, 89, 97. Notice that 2 is even and that it is the *only* even prime.

□

The following lemma expresses that every multiple is divisible by at least one prime. If we would have allowed 1 as a prime number this lemma would have been void.

7.20 Lemma. For every $a \in \mathbb{N}^{+2}$ a prime p exists such that $p|a$.

Proof. By Mathematical Induction on a . Firstly, if a is a prime then a is a prime and $a|a$. Secondly, if a is not prime then a multiple $b \in \mathbb{N}^{+2}$ exists satisfying $b \neq a$, and $b|a$. From Lemma 7.2, using $b|a$, we conclude that $b \leq a$, but because $b \neq a$ this amounts to $b < a$. So, by Induction Hypothesis, we may assume that p is a prime such that $p|b$. Because $b|a$, by the transitivity of divisibility, we then conclude $p|a$, as required.

□

The following theorem is important; it has been proved already by Euclides.

7.21 Theorem. The set of all primes is infinite.

Proof. One way to prove that a set is infinite is to prove that every *finite subset* of it differs from the whole set; that is, for every finite subset the whole set contains an element not in that subset. So, let V be a finite subset of the primes. Now we define multiple a by:

$$a = (\prod_{p \in V} p) + 1 .$$

Because the product $(\prod_{p \in V} p)$ is divisible by every $p \in V$, the number a is *not* divisible by p , for every $p \in V$. On account of Lemma 7.20, however, a is divisible by at least one prime, which therefore, is not an element of V .

□

* * *

A very old algorithm to compute “all” primes is known as *Erathostenes’s sieve*. This involves infinite enumerations of infinite subsets of the multiples, which is unfeasible, of course, but for the purpose of computing any finite number of primes, finite prefixes of these infinite enumerations will do. To compute all, infinitely many, primes would, of course, take an infinite amount of time. Yet, we call it an algorithm to compute “all” primes because it can be used to compute as many primes as needed in a finite amount of time.

Informally, the algorithm is presented as follows. One starts with writing down all –that is: sufficiently many– multiples in increasing order:

2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23,

The first number of this sequence, 2, is the first prime number, and we now construct a new sequence from the first one by eliminating all multiples of 2 from it:

3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, \dots .

This second sequence is an enumeration, again in increasing order, of all multiples that are not divisible by the first prime, 2. The first number, 3, of this second sequence is the second prime, and, again, we construct a third sequence from this second one, this time by eliminating all multiples of 3:

5, 7, 11, 13, 17, 19, 23, \dots .

This sequence contains all multiples that are not divisible by either 2 or 3; its first element, 5, is the *next* prime number, which is the smallest prime that is larger than 2 and 3. And so on...

The general properties on account of which this algorithm is correct are easily formulated. After n steps, for some $n \in \mathbb{N}$, a sequence is obtained that contains, in increasing order, all multiples that are *not* divisible by the smallest n prime numbers. The first number of this sequence then, because the sequence is increasing, is its minimum, and it can be proved that this minimum is the *next* prime, that is, the smallest prime number exceeding the smallest n prime numbers. By eliminating all multiples of this next prime the next sequence is obtained, containing all multiples not divisible by the first $n+1$ primes.

* * *

We have seen that the set of primes is infinite, but we may still ask for the *density* of the primes; that is, for any given multiple n we may ask *how many* primes are less than n . Because the number of potential divisors of n increases with n , the likelihood – not in the mathematical meaning of the word – that an arbitrary number is prime may be expected to decrease with increasing numbers.

The, so-called, *Prime Number Theorem* states that the number of primes less than n is approximately $n / \ln(n)$. This formula really gives an approximation only; for example, the number of primes less than 10^9 equals 50 847 534, whereas $10^9 / \ln(10^9)$ is 48 254 942 (rounded). An elementary proof of the Prime Number Theorem has been constructed by the famous Hungarian mathematician Pál Erdős.

The greatest common divisor of a prime p and a positive natural a can have only one out of two possible values: either $p|a$ and then $\gcd(p, a) = p$, or $\neg(p|a)$ and then $\gcd(p, a) = 1$. An important consequence of this is the following lemma.

7.22 Lemma. For every prime p and for all $a, b \in \mathbb{N}^+$: $p|(a*b) \Rightarrow p|a \vee p|b$.

Proof. By distinguishing the cases $p|a$ and $\neg(p|a)$ and, for the latter case, using Lemma 7.15.

□

A corollary of this lemma is the more general property that if a prime p is a divisor of any finite product $(\prod a : a \in V : a)$, where V is some finite *bag* of positive naturals, then $p|a$ for some element $a \in V$. These observations lead to the following, important, theorem, known as the Unique Prime Factorization Theorem.

7.23 Theorem. For every positive natural $a \in \mathbb{N}^+$ a *unique* bag V of prime numbers exists such that: $a = (\prod p : p \in V : p)$.

Proof. We prove *existence* and *uniqueness* separately.

Existence. By Mathematical Induction on a . Firstly, if $a=1$ then a has no other divisors than 1, so a has no prime divisors, and the only bag of primes the product of which equals 1 is the empty bag, which we denote as $\{\}$ here: $1 = (\prod p : p \in \{\} : p)$ ¹³. Primes are multiples, so the product of any non-empty bag of primes differs from 1.

Secondly, if $a \geq 2$ then, by Lemma 7.20, a prime q , say, exists such that $q|a$; so, a/q is a positive natural and $a/q < a$. Now, by Induction Hypothesis, let V be a bag of primes such that $a/q = (\prod p : p \in V : p)$. Then we have that a is equal to $q * (\prod p : p \in V : p)$, which is equal to $(\prod p : p \in V + \{q\} : p)$

Uniqueness. We prove that any two bags U and V of primes satisfy:

$$(\prod p : p \in U : p) = (\prod p : p \in V : p) \Rightarrow U = V ,$$

by Mathematical Induction on $\#U + \#V$. Firstly, if both U and V are empty then the proposition is true because then $U=V$. Secondly, if one of U and V is empty and the other one is non-empty, the proposition is true because the product of the empty bag equals 1 whereas the product of a non-empty bag of primes is larger than 1; so, in this case the left-hand side of the implication is false. Thirdly, remains the case that both U and V are non-empty. We assume that their products are equal. Now let $q \in U$ then q is a divisor of $(\prod p : p \in U : p)$; hence, because the products are equal, q also is a divisor of $(\prod p : p \in V : p)$. Hence, by Lemma 7.22, we have $q|p$ for some $p \in V$, but, because such p is prime this amounts to $q=p$. So, we have $q \in V$ as well. Thus, we obtain that $(\prod p : p \in U - \{q\} : p) = (\prod p : p \in V - \{q\} : p)$, which, by Induction Hypothesis, implies that $U - \{q\} = V - \{q\}$; this is equivalent to $U = V$, as required.

□

* * *

By means of prime factorization we can look at Greatest Common Divisors (*gcd*) and Least Common Multiples (*lcm*) in a different way. Let $a, b \in \mathbb{N}^+$. If p is a prime and $p|a$ and $p|b$ then, of course, we also have $p|gcd(a, b)$. Now let m be the *largest* natural number satisfying both $p^m|a$ and $p^m|b$, then also $p^m|gcd(a, b)$, and m also is the largest natural for which this is true. Moreover, let k be the largest natural for which $p^k|a$ and let l be the largest natural for which $p^l|b$ then we have that $m = k \min l$. As a matter of fact, we now have that $gcd(a, b)$ is the product

¹³Recall that the product of an empty bag of numbers, by definition, equals 1 because 1 is the identity element of multiplication.

of all such numbers p^m . More precisely, the bag of primes of which $\gcd(a, b)$ is the product is obtained by joining the bags of primes factorizing a and b by taking the *minimums* of the occurrences of all primes.

The Least Common Multiple of a and b , denoted as $\text{lcm}(a, b)$, is defined as the smallest positive natural number that is a multiple of both a and b . So, a and b both are divisors of $\text{lcm}(a, b)$, and $\text{lcm}(a, b)$ is the smallest of all positive naturals with this property. So, with k and l for the (maximal) numbers of occurrences of prime p in a and b , respectively, we have that p occurs $k \max l$ times in $\text{lcm}(a, b)$.

Thus, we obtain the following result. If prime p occurs (exactly) k times in a and if p occurs l times in b , then p occurs $k \min l$ times in $\gcd(a, b)$ and $k \max l$ times in $\text{lcm}(a, b)$. Now we also have that $(k \min l) + (k \max l) = k + l$, and p occurs $k + l$ times in $a * b$. Because this is true for every prime p we conclude that the bag of primes factorizing $\gcd(a, b) * \text{lcm}(a, b)$ is equal to the bag of primes factorizing $a * b$; hence:

$$\gcd(a, b) * \text{lcm}(a, b) = a * b .$$

7.7 Modular Arithmetic

7.7.1 Congruence relations

Almost everybody probably knows that the product of two *even* integers is even, and that the product of two *odd* integers is odd. Also, the sum of two even integers is even too, and even the sum of two odd numbers is even. The point is that, apparently, whether the result of an operation, like addition or multiplication, is even or odd *only* depends on whether the arguments of the operation are even or odd.

The proposition that integer “ x is even” is equivalent to “ x is divisible by 2”, which in turn is equivalent to $x \bmod 2 = 0$; similarly, the proposition “ x is odd” is equivalent to $x \bmod 2 = 1$. That the property “being even” of the sum of two integers only depends on the “being even” of these two numbers know means that $(x + y) \bmod 2$ only depends on $x \bmod 2$ and $y \bmod 2$ (and not on $x \text{div} 2$ or $y \text{div} 2$). In formula this is rendered as:

$$(30) \quad (x + y) \bmod 2 = (x \bmod 2 + y \bmod 2) \bmod 2 \quad , \text{ for all } x, y \in \mathbb{Z} .$$

Properties like these are not specific for 2 as a divisor: similar properties hold for all positive divisors. From the chapter on relations we recall that, for every function of type $B \rightarrow V$, the relation, on its domain B , “having the same function value” is an equivalence relation. For any fixed $d \in \mathbb{N}^+$, the function $(\bmod d)$ that maps every $x \in \mathbb{Z}$ to $x \bmod d$ has type $\mathbb{Z} \rightarrow [0..d)$. This function induces an equivalence relation, on \mathbb{Z} , of “having the same remainder when divided by d ”. This relation partitions \mathbb{Z} into d different (and, as always, disjoint) equivalence classes, namely one for every value of the function $(\bmod d)$: the equivalence class corresponding to $a \in [0..d)$ is the set

$$\{ x \in \mathbb{Z} \mid x \bmod d = a \} ,$$

which can also be formulated as:

$$\{ q * d + a \mid q \in \mathbb{Z} \} .$$

In particular, of course, $a \in \{ x \in \mathbb{Z} \mid x \bmod d = a \}$, because $a \bmod d = a$, for every $a \in [0..d)$.

Now properties similar to (30) also hold in this case; that is, we now have:

$$(31) \quad (x + y) \bmod d = (x \bmod d + y \bmod d) \bmod d \quad , \text{ for all } x, y \in \mathbb{Z} .$$

A similar property holds for subtraction and multiplication:

$$(32) \quad (x * y) \bmod d = (x \bmod d * y \bmod d) \bmod d \quad , \text{ for all } x, y \in \mathbb{Z} .$$

A consequence of propositions like (31) and (32) is that equivalence is *preserved under* arithmetic operations like addition and multiplication. Using (31), for instance, we can now derive, for all $x, y, z \in \mathbb{Z}$:

$$(33) \quad x \bmod d = y \bmod d \Rightarrow (x + z) \bmod d = (y + z) \bmod d .$$

An equivalence relation that is preserved under a given set of operations is called a *congruence relation*. In our case, the relation “having the same remainder when divided by d ” is congruent with the operations addition, subtraction, and multiplication. Conversely, we also say that the operations addition, subtraction, and multiplication are *compatible with* the relation.

Notation: According to mathematical convention, the fact that x and y are congruent modulo d is often denoted as:

$$x = y \pmod{d} .$$

This notation is somewhat awkward, though, because it is not very clear what the *scope* is of the suffix “ \pmod{d} ”. Apparently, its scope extends over the complete equality textually preceding it; that is, if we would use some sort of brackets to delineate the scope of “ \pmod{d} ” more explicitly, we should write something like:

$$[[x = y \pmod{d}]] .$$

Because the relation is a congruence relation and because it is not the same as sheer equality, although it resembles it, it seems better to denote the relation as an infix symbol resembling but different from “ $=$ ”. For example, the symbol “ $=_{\bmod d}$ ” would be appropriate, as the subscript explicitly indicates the nature of the congruence. In this text we will abbreviate this to “ $=_d$ ”; so, by definition we now have, for all $d \in \mathbb{N}^+$ and $x, y \in \mathbb{Z}$:

$$x =_d y \Leftrightarrow x \bmod d = y \bmod d .$$

For example, congruence property (33) can now be rendered as, for all $x, y, z \in \mathbb{Z}$:

$$x =_d y \Rightarrow x + z =_d y + z .$$

□

Other algebraic properties are compatible with congruence modulo d too. The numbers 0 and 1, for instance, are the identity elements of addition and multiplication, respectively, and this remains so under congruence. In addition, the property that 0 is a zero-element of multiplication –that is: $0 * x = 0$ – is retained. Finally, that multiplication distributes over addition remains true as well.

For any given $d \in \mathbb{N}^+$ we can now define binary operations \oplus and \otimes , say, by, for all $x, y \in \mathbb{Z}$:

$$x \oplus y = (x + y) \bmod d \quad , \text{ and:}$$

$$x \otimes y = (x * y) \bmod d \quad , \text{ and:}$$

(If we would be very strict we should make the dependence on d explicit by writing \oplus_d and \otimes_d .) Now \oplus and \otimes have type $[0..d) \times [0..d) \rightarrow [0..d)$ and with these operators various algebraic structures can be formed, which we mention here without further elaboration or proofs.

7.24 Properties. For all $d \in \mathbb{N}^+$:

- (a) $([0..d), \oplus, 0)$ is a group.
- (b) $([0..d), \otimes, 1)$ is a monoid but not a group.
- (c) $([1..d), \otimes, 1)$ is a group if and only if d is prime.
- (d) $([0..d), \oplus, \otimes, 0, 1)$ is a, so-called, *ring*.
- (e) $([0..d), \oplus, \otimes, 0, 1)$ is a, so-called, *field* if, but not only if, d is prime.

In traditional mathematical parlance these structures are usually denoted differently. The ring $([0..d), \oplus, \otimes, 0, 1)$, for instance, then is denoted as $(\mathbb{Z}/d\mathbb{Z}, +, *, 0, 1)$, where $\mathbb{Z}/d\mathbb{Z}$ denotes the set $[0..d)$ –actually, the set of the equivalence classes– and where $+$ and $*$ must be read as addition and multiplication modulo d , that is, as our operations \oplus and \otimes .

□

* * *

Another important property is that two integers are congruent modulo d if and only if their difference is divisible by d ; that is, for all $x, y \in \mathbb{Z}$ we have:

$$x =_d y \Leftrightarrow (x - y) \bmod d = 0 \quad .$$

7.7.2 An application: the nine and eleven tests

A technique that was commonly applied to verify manual calculations is the, so-called, *nine test*. This is based on the property that, in our decimal number representation, the remainder of a number when divided by 9 can be easily calculated: it equals the remainder of the sum of the number's digits modulo 9. As the sum of the digits of a number usually is much smaller than the number itself, the problem has been reduced. This process is repeated until a number is obtained that is less than 10: this last number, then, is the remainder of the original number modulo 9, except when it is 9 in which case the remainder is 0, of course.

For example, the sum of the digits of the number 123456789 equals 45, and the sum of the digits of 45 is 9. Hence, the remainder of 123456789 modulo 9 is 0.

Now to verify a calculation, for instance the addition or multiplication of two large numbers, one calculates the remainders modulo 9 of both numbers and of the result of the calculation, and one performs the same operation, modulo 9, to these remainders. If the results match we can be pretty confident that our calculation was correct, although we do not have certainty, of course. But, if the results do not match we certainly have made an error!

* * *

The property that the remainder modulo 9 of a number equals the remainder modulo 9 of the sum of the digits of that number's decimal representation is based on the observation that $10 \bmod 9 = 1$. Now we have that a number like, for instance, 1437 is equal to $143 * 10 + 7$. Therefore, we have:

$$\begin{aligned}
 & 1437 \bmod 9 \\
 = & \quad \{ \text{above property} \} \\
 & (143 * 10 + 7) \bmod 9 \\
 = & \quad \{ 10 = 9 + 1 \} \\
 & (143 * 9 + 143 * 1 + 7) \bmod 9 \\
 = & \quad \{ \text{mod over +; multiples of 9 may be discarded and } 7 \bmod 9 = 7 \} \\
 & (143 \bmod 9 + 7) \bmod 9 .
 \end{aligned}$$

This calculation shows that $1437 \bmod 9$ is equal to $143 \bmod 9$ plus 7, modulo 9; if now, by Induction Hypothesis, $143 \bmod 9$ is equal to the sum, modulo 9, of the digits of 143 then $1437 \bmod 9$ also is equal to the sum of its digits, modulo 9.

* * *

In general, the decimal representation of natural numbers can be defined in a recursive way, as follows. A sequence of n decimal digits " $d_{n-1} \cdots d_2 d_1 d_0$ " represents the natural number d_0 if $n = 1$: in this case the sequence just is a single digit, " d_0 ". If $n \geq 2$, the number represented by the sequence is equal to the number represented by the sequence of $n - 1$ digits " $d_{n-1} \cdots d_2 d_1$ " times 10 plus d_0 .

By means of this recursive definition it can be proved, by Mathematical Induction, of course, that the number represented by a sequence of decimal digits and the sum of these digits are congruent modulo 9. By means of the recursive definition it also is possible to prove that the number represented by the sequence of digits “ $d_{n-1} \cdots d_2 d_1 d_0$ ” is equal to:

$$(\sum_i: 0 \leq i < n: d_i * 10^i) \text{ ,}$$

but in most cases the recursive definition is more manageable than this expression.

* * *

In a very similar, albeit slightly more complicated way we observe that 10 is congruent to -1 modulo 11, and that $100 \bmod 11$ is equal to 1. This is the basis of the *eleven test*: the remainder of a natural number modulo 11 is equal to the remainder modulo 11 of the sum of the digits of that number’s decimal representation, but here the digits are added with *alternating signs*, starting at the least-significant digit with a positive sign. The number 123456789, for example, is congruent modulo 11 to: $+9 - 8 + 7 - 6 + 5 - 4 + 3 - 2 + 1$, which equals 5.

As was the case with the nine test, the eleven test can be used to “verify” the results of calculations.

7.7.3 Linear congruence equations

We consider, for given $a \in \mathbb{N}^+$ and $b \in [0..d)$, the equation $a * x =_d b$, so x is the unknown here. Can we determine all integers x satisfying this, preferably in a systematic way?

Well, the relation “ $=_d$ ” means that we are actually considering this equation:

$$(34) \quad (a * x) \bmod d = b \bmod d \text{ ,}$$

which, because $b \bmod d = b$, is equivalent to:

$$(a * x) \bmod d = b \text{ ,}$$

and also to:

$$d \mid (a * x - b) \text{ .}$$

This means that an integer y exists satisfying:

$$a * x - b = d * y \text{ .}$$

Demonstrating the existence of such y amounts to determining a suitable value of it, so in fact we are dealing with an equation with *two* unknowns, x and y . This equation can be rewritten equivalently as:

$$a * x - d * y = b \text{ .}$$

Because the range of y is \mathbb{Z} we can apply a substitution –comparable to a dummy transformation– $y := -y$, so as to transform the equation into this form:

$$(35) \quad a * x + d * y = b \ .$$

This equation, however, we know how to solve: this equation is, except for the names of the variables, the same as equation (28) in Section 7.5! Moreover, a method to solve it has been presented there too. Recall that this equation has solutions if and only if $\gcd(a, d)$ is a divisor of b .

In our case, of course, we are not really interested in the solutions for y : our original equation (34) only has x as the unknown. This means that, in the set of all solutions for x and y obtained for (35) we can simply discard the values for y .

* * *

The method in Section 7.5 shows that if x is a solution to (34) then so is $x + z * d$, for all $z \in \mathbb{Z}$. But this means that x may be replaced by $x \bmod d$ without losing generality; that is, if x is a solution then all solutions are $x \bmod d + z * d$, for all $z \in \mathbb{Z}$. This also shows that, modulo d , the solution is *unique* because all numbers $x \bmod d + z * d$ are congruent modulo d .

But now a second method to solve equation (34) emerges. If all solutions will be numbers of the shape $r + z * d$, for all $z \in \mathbb{Z}$, for some yet to be determined $r \in [0..d)$, then we might as well start with substituting r for x in the equation, and solve r from this under the additional restriction $r \in [0..d)$! Particularly if d is not very large, the number of potential values for r is small enough to make enumeration and verification of all of them feasible.

7.7.4 An example

As an example, we will determine all $x \in \mathbb{Z}$ satisfying: $11 * x =_{17} 13$. This proposition is equivalent to the proposition that $11 * x - 13$ is divisible by 17, which, in turn, means that an integer y exists satisfying:

$$11 * x - 13 = 17 * y \ .$$

By means of the substitution $y := -y$ and some simple rewriting this equation is transformed into the following one:

$$(36) \quad 11 * x + 17 * y = 13 \ .$$

In Section 7.5 we have seen that this equation has a solution if and only if $\gcd(11, 17)$ is a divisor of 13. This is the case, because $\gcd(11, 17) = 1$. Now we first solve this equation:

$$11 * x + 17 * y = 1 \ .$$

This equation can be solved by means of Euclid's extended algorithm, because the 1 in its right-hand side equals $\gcd(11, 17)$. This yields $x = 14$ and $y = -9$ as a solution. Multiplication by 13 then yields $x = 182$ and $y = -117$ as a solution to equation (36), so $x = 182$ is a solution to our original equation $11 * x = 13 \pmod{17}$. In addition, as we have seen in Section 7.5, adding (or subtracting) multiples of 17

–that is: the coefficient of y in (36)– to a solution yields new solutions again. So, all solutions to $11 * x = 13 \pmod{17}$ are the numbers $182 + z * 17$, for all $z \in \mathbb{Z}$. Because multiples of 17 may be added or subtracted *ad libitum*, the number 182 in this expression may be reduced modulo 17, if so desired; that is, it may be replaced by $182 \pmod{17}$. Thus, all solutions to our equation also may be expressed as, somewhat simpler: $12 + z * 17$, for all $z \in \mathbb{Z}$.

* * *

To illustrate the alternative method discussed in the previous subsection, we now solve the same problem by means of this method. This means that we are looking for an $r \in [0..17)$ satisfying $(11 * r - 13) \pmod{17} = 0$. We obtain it by tabulating the 17 possible values of r together with the corresponding values of the expression $(11 * r - 13) \pmod{17}$, until we encounter 0:

0	4
1	15
2	9
3	3
4	14
5	8
6	2
7	13
8	7
9	1
10	12
11	6
12	0

Thus we find that $r = 12$ solves the equation and, hence, that all solutions are the numbers $12 + z * 17$, for all $z \in \mathbb{Z}$. Notice that the process of constructing this table can be aborted as soon as the value 0 is encountered. (Exercise: why?)

7.7.5 Multiple linear congruences: an example

As another example we consider the system of equations $x =_{21} 5$ and $x =_{35} 12$, where x is the unknown. So, one-and-the-same x must satisfy both equations.

The proposition $x =_{21} 5$ means that $x - 5$ is divisible by 21. Because $21 = 3 * 7$ this is equivalent to the proposition that $x - 5$ is divisible both by 3 and by 7; So, the equation $x =_{21} 5$ can be reformulated equivalently as the conjunction of $x =_3 5$ and $x =_7 5$. In very much the same way we observe that $35 = 5 * 7$; therefore the equation $x =_{35} 12$ can be reformulated equivalently as the conjunction of $x =_5 12$ and $x =_7 12$.

Hence, the whole system of equations $x =_{21} 5$ and $x =_{35} 12$ is equivalent to the following system:

$$x =_3 5 \wedge x =_7 5 \wedge x =_5 12 \wedge x =_7 12 .$$

Because $12-5$ is a multiple of 7 , the two equations $x =_7 5$ and $x =_7 12$ are equivalent: they have the same solutions and, therefore, one of them may be omitted. (And: if $12-5$ would *not* be a multiple of 7 the two propositions would be *contradictory* and the system of equations would have no solutions at all.)

Thus, we obtain this, somewhat simplified, system of equations:

$$x =_3 5 \wedge x =_7 5 \wedge x =_5 12 .$$

We now solve this system by solving each of its three parts, one by one, but in such a way that, for each next part, we express the solution in terms of the solutions for preceding parts.

The proposition $x =_3 5$ means that $x - 5$ is a multiple of 3 , so its solutions are the integers $5 + 3 * u$, for all $u \in \mathbb{Z}$. The solutions of the second equation, $x =_7 5$, are a subset of the solutions of the first one. To obtain them, we substitute $7 * v + r$ for u , with $v \in \mathbb{Z}$ and $0 \leq r < 7$. Hence, the solutions of the first equation, the numbers $5 + 3 * u$, for all $u \in \mathbb{Z}$, are rewritten as the numbers $5 + 3 * (7 * v + r)$, for all $v \in \mathbb{Z}$ and $r \in [0..7)$. In order to identify which of these numbers are solutions to the second equation as well, we now calculate as follows:

$$\begin{aligned} & x =_7 5 \\ \Leftrightarrow & \quad \{ \text{definition of } =_7 \} \\ & (x - 5) \bmod 7 = 0 \\ \Leftrightarrow & \quad \{ \text{substitute } x := 5 + 3 * (7 * v + r) \} \\ & (5 + 3 * (7 * v + r) - 5) \bmod 7 = 0 \\ \Leftrightarrow & \quad \{ \text{algebra} \} \\ & (21 * v + 3 * r) \bmod 7 = 0 \\ \Leftrightarrow & \quad \{ \text{property of mod: } 21 * v \text{ is divisible by } 7 \} \\ & (3 * r) \bmod 7 = 0 . \end{aligned}$$

Thus, we conclude that the second equation poses restrictions on r only (and not on v). In view of the limited range of values of r , after all we have $r \in [0..7)$, we conclude that the only solution to $(3 * r) \bmod 7 = 0$ is $r = 0$. Actually, the reason to rewrite u as $7 * v + r$ was to introduce a factor 7 , in the hope that v would disappear from the equation, as it did indeed.

Thus, we obtain that the solutions of the system $x =_3 5$ and $x =_7 5$ are the integers $5 + 21 * v$, for all $v \in \mathbb{Z}$. (In retrospect, we could have discovered this in a more direct way: these numbers are the solutions to the equation $x =_{21} 5$, and the factorization of 21 into 3 and 7 has not been really necessary.)

Finally, we take the last equation, $x =_5 12$, into account, by applying the same technique once more. That is, we substitute $5 * w + s$ for v , in the above expression for the solutions thus far, where $w \in \mathbb{Z}$ and $s \in [0..5)$, and after this substitution we calculate again, as follows:

$$\begin{aligned}
& x =_5 12 \\
\Leftrightarrow & \quad \{ \text{definition of } =_5 \} \\
& (x - 12) \bmod 5 = 0 \\
\Leftrightarrow & \quad \{ \text{substitute } x := 5 + 21 * (5 * w + s) \} \\
& (5 + 21 * (5 * w + s) - 12) \bmod 5 = 0 \\
\Leftrightarrow & \quad \{ \text{algebra} \} \\
& (105 * w + 21 * s - 7) \bmod 5 = 0 \\
\Leftrightarrow & \quad \{ \text{property of mod: } 105 * w \text{ is divisible by } 5 \} \\
& (21 * s - 7) \bmod 5 = 0 \\
\Leftrightarrow & \quad \{ \text{properties of mod} \} \\
& (21 * s) \bmod 5 = 2 \ .
\end{aligned}$$

In view of the limited range of values of s we conclude that the only solution to $(21 * s) \bmod 5 = 2$ is $s = 2$.

Thus, we obtain the solution to the whole system of equations is: the set of all integers $5 + 21 * (5 * w + 2)$, for all $w \in \mathbb{Z}$. If so desired, this expression can be simplified to $47 + 105 * w$, for all $w \in \mathbb{Z}$.

7.7.6 Two linear congruences: the general case

We now consider the general situation of the following two congruences:

$$(37) \quad x =_p a \ \wedge \ x =_q b \ ,$$

for given positive naturals a, b and p, q . In addition we assume that $\gcd(p, q) = 1$. If $\gcd(p, q) \neq 1$ we can factor out the greatest common divisor by means of the same trick we employed in the previous subsection, where we transformed a system of two equations with moduli 21 and 35 into a system of three equations with moduli 3, 5, and 7.

Because, by the definition of $=_p$ and $=_q$, the equation is equivalent to:

$$x \bmod p = a \bmod p \ \wedge \ x \bmod q = b \bmod q \ ,$$

we may as well assume that $a \in [0..p)$ and $b \in [0..q)$, because $a \bmod p$ and $b \bmod q$ are all that matter. With these assumptions we may simplify $a \bmod p$ and $b \bmod q$ to a and b respectively. Thus, our equation becomes:

$$x \bmod p = a \ \wedge \ x \bmod q = b \ .$$

Now all solutions to the first conjunct are easily characterized as the integers $a + u * p$, for all $u \in \mathbb{Z}$, with a being the unique solution in the interval $[0..p)$. As with the example in the previous subsection, we now take the second equation into account by substituting $r + v * q$ for u in this formula, where $v \in \mathbb{Z}$ and $r \in [0..q)$, and we obtain $a + (r + v * q) * p$. As in the last example, we now calculate:

$$\begin{aligned}
& x \bmod q = b \\
\Leftrightarrow & \quad \{ \text{property of mod} \} \\
& (x - b) \bmod q = 0 \\
\Leftrightarrow & \quad \{ \text{substitute } x := a + (r + v * q) * p \} \\
& (a + (r + v * q) * p - b) \bmod q = 0 \\
\Leftrightarrow & \quad \{ \text{algebra} \} \\
& ((v * p) * q + r * p + a - b) \bmod q = 0 \\
\Leftrightarrow & \quad \{ \text{properties of mod} \} \\
& (r * p) \bmod q = (b - a) \bmod q .
\end{aligned}$$

This is a kind of equation we have seen earlier. Because $\gcd(p, q) = 1$ a *unique* solution for r exists in the range $[0..q)$. Calling this solution r_0 , we obtain as formula for all solutions for x in our original equation (37):

$$(38) \quad a + r_0 * p + v * (p * q) \quad , \text{ for all } v \in \mathbb{Z} .$$

Because $0 \leq a < p$ and $0 \leq r_0 < q$ we also have $0 \leq a + r_0 * p < p * q$, and this is the only solution in the interval $[0..p * q)$, because of the shape of the additional term $v * (p * q)$. So, modulo $p * q$, the solution is unique.

This also means that our original equation (37) is equivalent to the following, single congruence equation, with unknown x and in which r_0 is the solution for r in $(r * p) \bmod q = (b - a) \bmod q$, as defined above:

$$(39) \quad x =_{p * q} a + r_0 * p .$$

We will need this result to prove the Chinese Remainder Theorem, in the next subsection.

* * *

Here is a different way to solve equation (37), which does not destroy the symmetry between the two conjuncts of the equation. These two conjuncts of the equation are:

$$(37) \quad x =_p a \quad \text{and:} \quad x =_q b ,$$

and taken in isolation they can be solved easily. The solutions to $x =_p a$ are all integers $a + u * p$, where $u \in \mathbb{Z}$, and the solutions to $x =_q b$ are all integers $b + v * q$, where $v \in \mathbb{Z}$. The set of solutions to both equations is the intersection of these two sets, so we obtain a solution provided that $u \in \mathbb{Z}$ and $v \in \mathbb{Z}$ satisfy:

$$(40) \quad a + u * p = b + v * q ,$$

which can be rewritten to:

$$u * p - v * q = b - a .$$

This equation, now with unknowns u and v , however, we know how to solve: it is the type of equation we discussed in Section 7.5. There we have seen that a necessary (and sufficient) condition for this equation to have solutions at all is that $\gcd(p, q)$ is a divisor of $b - a$. This condition is satisfied, because of our assumption $\gcd(p, q) = 1$. By means of the technique from Section 7.5 we can calculate values u_0 and v_0 , say, satisfying:

$$a + u_0 * p = b + v_0 * q ,$$

and all solutions for u and v in (40) are the pairs $u_0 + z * q$ and $v_0 + z * p$, for all $z \in \mathbb{Z}$. Because both sides of equation (40) represent the solutions for x in our original equation (37), we obtain, by straightforward substitution, as formula for all solutions to (37):

$$a + (u_0 + z * q) * p , \text{ for all } z \in \mathbb{Z} ,$$

which can be rewritten to:

$$a + u_0 * p + z * (p * q) , \text{ for all } z \in \mathbb{Z} .$$

This formula is, except for the names of the variables, exactly the same as formula (38) obtained earlier. And, of course, the solutions also are characterized by this formula's symmetric counterpart:

$$b + v_0 * q + z * (p * q) , \text{ for all } z \in \mathbb{Z} .$$

7.7.7 The Chinese Remainder Theorem

A very old theorem in number theory is known as the “Chinese Remainder Theorem”. It states that *any* number of congruence equations, of the kind discussed in the previous subsections, has a unique solution.

7.25 Theorem. For all $k \in \mathbb{N}$ the system of equations $x = a_i \pmod{q_i}$, for $i: 0 \leq i \leq k$, has a solution that is *unique* modulo $(\prod i: 0 \leq i \leq k: q_i)$, provided the values q are pair-wise relatively prime, that is: provided $\gcd(q_i, q_j) = 1$, for all $i, j: 0 \leq i < j \leq k$.

Proof. By Mathematical Induction on k . Firstly, if $k = 0$ the system consists of a single equation, $x = a_0 \pmod{q_0}$, only, which we know how to solve. The solution, modulo q_0 , is $a_0 \pmod{q_0}$ and this is unique modulo q_0 .

Secondly, we consider a system of $k + 1$ such equations, for some $k \in \mathbb{N}$. The two equations with the largest indices then are:

$$x = a_k \pmod{q_k} \text{ and: } x = a_{k+1} \pmod{q_{k+1}} ,$$

where $\gcd(q_k, q_{k+1}) = 1$. As we have seen in Subsection 7.7.6 the conjunction of these two equations is equivalent to the single equation – see formula (39) –:

$$x = a_k + r_0 * q_k \pmod{q_k * q_{k+1}} ,$$

where r_0 is the solution for r in $(r * q_k) \bmod q_{k+1} = (a_{k+1} - a_k) \bmod q_{k+1}$. Thus, the system of $k+1$ equations can be transformed into an equivalent system of k equations, by replacing equations k and $k+1$ by the single equation modulo $q_k * q_{k+1}$.

Moreover, we have that $\gcd(q_i, q_k * q_{k+1}) = 1$, for all $i: 0 \leq i < k$, on account of Lemma 7.26, below. Therefore, by Induction Hypothesis, the new system of k equations has a unique solution modulo $(\prod i: 0 \leq i \leq k: q_i)$, and so does the original system of $k+1$ equations. Notice that the product $(\prod i: 0 \leq i \leq k: q_i)$ is not affected by the transformation.

□

In the proof of this theorem we have used the following lemma, the proof of which we leave to the exercises.

7.26 Lemma. For all $a, b, c, \in \mathbb{N}^+$ we have:

$$\gcd(a, c) = 1 \wedge \gcd(b, c) = 1 \Rightarrow \gcd(a * b, c) = 1 .$$

□

7.8 Fermat's little theorem

For a fixed $p \in \mathbb{N}^{+2}$ we define, on \mathbb{Z} , a binary operator \otimes by $x \otimes y = (x * y) \bmod p$, for all $x, y \in \mathbb{Z}$. Then the structure $([1..p], \otimes, 1)$ is a *monoid* and this structure is a group if and only if p is prime.

7.27 Lemma. For all $p \in \mathbb{N}^{+2}$ we have:

$$“([1..p], \otimes, 1) \text{ is a group}” \Leftrightarrow “p \text{ is prime}” .$$

□

Using this lemma we can prove the following theorem which is known as “Fermat's little theorem”.

7.28 Theorem. For every prime number p and for every $a \in \mathbb{N}^+$ we have:

$$\neg(p|a) \Rightarrow a^{p-1} \bmod p = 1 .$$

Proof. Let p be a prime and let $a \in \mathbb{N}^+$. Assume $\neg(p|a)$; this is equivalent to $a \bmod p \neq 0$, so we have: $1 \leq a \bmod p < p$; that is, $a \bmod p \in [1..p)$. For the sake of brevity, we define $b = a \bmod p$. Now we have that $a^{p-1} \bmod p$ is equal to b^{p-1} , where $b \in [1..p)$ and b^{p-1} is to be interpreted in terms of \otimes -operations, instead of $*$. Hence, we also have $b^{p-1} \in [1..p)$.

By Lemma 7.27 we have that $([1..p), \otimes, 1)$ is a group, because p is prime. The set $\{b^i \mid 0 \leq i\}$ together with \otimes and 1 is the subgroup of $([1..p), \otimes, 1)$ generated by b . Because the whole group is finite, of size $p-1$, so is this subgroup. Therefore, a number $n \in \mathbb{N}^+$ exists such that $b^n = 1$ and $b^i \neq 1$, for all $i: 1 \leq i < n$. Then we have: $\{b^i \mid 0 \leq i\} = \{b^i \mid 0 \leq i < n\}$, and n is the size of this set.

On account of Lagrange's Theorem we conclude $n \mid (p-1)$. So, let $z \in \mathbb{N}$ satisfy $p-1 = z * n$. Now we derive:

$$\begin{aligned}
& a^{p-1} \bmod p \\
= & \quad \{ \text{as observed above} \} \\
& b^{p-1} \\
= & \quad \{ \text{definition of } z \} \\
& b^{z*n} \\
= & \quad \{ \text{property of exponentiation} \} \\
& (b^n)^z \\
= & \quad \{ \text{definition of } n \} \\
& 1^z \\
= & \quad \{ 1 \text{ is the identity of } \otimes \} \\
& 1
\end{aligned}$$

□

7.9 Cryptography: the RSA algorithm

We conclude this chapter with a practical application of the theory, namely in the area of *cryptography*, which is the art of transmitting messages in a secure way, such that these messages cannot be read by anyone else than the intended receiver. For this purpose the messages are *encrypted* in such a way that they become unintelligible, except for the intended receiver who is the only one able to *decypher* the messages.

The algorithms for encryption and decryption themselves usually are not kept secret, but the parameters used in the process are. In cryptography such parameters usually are called *keys*.

Here we discuss the, so-called, RSA-algorithm, named after its inventors: Rivest, Shamir, and Adleman. This is an example of a so-called *public key* system. This means that the keys needed for encryption and decryption are chosen by the receiver of the messages, and the receiver makes the encryption key publicly known: everybody who wishes to send a message to this particular receiver now can use this public encryption key. The security of this arrangement is based on the assumption that it is (virtually) impossible to infer the (secret) decryption key from the (public) encryption key.

In older encryption schemes the encryption and decryption keys used to be chosen by the sender of the messages, and the sender now was faced with the problem how to communicate the decryption key to the intended receiver(s) in a secure way. In a public key system this difficulty is avoided.

The security of the RSA-algorithm rest on the assumption that it is very hard to factorize very large numbers into their prime factors. Here “very large numbers” means: numbers the decimal representation of which comprises several hundreds of digits.

* * *

The intended receiver of secure messages picks two different and very large prime numbers p and q , say, and calculates their product, which we call n . So, $n = p * q$. For any positive natural number a that is smaller than both p and q we have that neither p nor q divides a . Therefore, by Fermat's little theorem, we have:

$$a^{p-1} \bmod p = 1 \quad \wedge \quad a^{q-1} \bmod q = 1 \quad .$$

From this it follows that we also have, because 1^{p-1} and 1^{q-1} are equal to 1:

$$a^{(p-1)*(q-1)} \bmod p = 1 \quad \wedge \quad a^{(p-1)*(q-1)} \bmod q = 1 \quad .$$

By means of the equation solving technique from Section 7.7.6 we can show that we now also have, using that $\gcd(p, q) = 1$ as p and q are different primes, and recalling that $n = p * q$:

$$(41) \quad a^{(p-1)*(q-1)} \bmod n = 1 \quad .$$

Message encryption is now performed as follows. The product of the two primes p and q , that is, the number n , is the public key used for encryption. A (digital) message to be transmitted can be represented, one way or another, as a finite sequence of bits, and every such sequence can be interpreted as a natural number, according to the rules of the binary number system. It is, of course, not difficult to ensure that this number m , say, is positive. So, the message to be transmitted is represented as a positive natural m , of which we also assume that $m < p$ and $m < q$. (A message that is too large to meet these latter requirements can, of course, be partitioned into several smaller messages, to be encrypted and transmitted separately.) The sender now computes $m^{127} \bmod n$, and this is the encrypted message actually sent to the intended receiver. (The number 127 in this expression seems rather arbitrary; actually, it is, and in practical applications larger exponents are used.)

The receiver now receives a number, of which he knows that it is a number of the shape $m^{127} \bmod n$; he also knows that $m < p$ and $m < q$. To decypher this message he must solve m from this. For this purpose, the receiver first solves another equation, namely: $(127 * x) \bmod (p-1) * (q-1) = 1$. As this equation does not depend on m , it has to be solved only once: the same solution x can be used to decrypt all messages that have been encrypted by means of n . Solving this equation means that the receiver, in fact, calculates integers x and y satisfying:

$$127 * x = y * (p-1) * (q-1) + 1 \quad .$$

To decrypt the message, that is, to calculate m from $m^{127} \bmod n$ the receiver now performs the following calculation, in which all operations are operations modulo n :

$$\begin{aligned} & (m^{127})^x \\ =_n & \quad \{ \text{algebra, and the relation between } x \text{ and } y \} \\ & m^{y*(p-1)*(q-1)+1} \\ =_n & \quad \{ \text{algebra} \} \end{aligned}$$

$$\begin{aligned}
& (m^{(p-1)*(q-1)})^y * m \\
=_{n} & \quad \{ \text{property (41), using } m < p \text{ and } m < q \} \\
& 1^y * m \\
=_{n} & \quad \{ \text{algebra} \} \\
& m .
\end{aligned}$$

Hence, by raising the value received, $m^{127} \bmod n$, to the power x and calculating the remainder modulo n , the receiver obtains $m \bmod n$. Because $m < p$ and $m < q$ we certainly have $m < n$ as well, so $m \bmod n = m$. Hence, we have obtained our final result:

$$(m^{127} \bmod n)^x \bmod n = m .$$

The security of this encryption/decryption scheme rests on the assumption that, given a product, n , of two prime numbers p and q , it is unfeasible – in a reasonable amount of time, that is – to compute p and q from n alone. This assumption has proved to be fair enough in practice, provided the primes p and q are sufficiently large.

Question: In order that this algorithm be executable, an additional restriction must be imposed upon the primes p and q . What restriction?

□

7.10 Exercises

- Let $a, d \in \mathbb{Z}$ and $d \neq 0$. Assuming that a is divisible by d , prove that the value q satisfying $a = q * d$ is unique.
- What are the divisors of 1? Prove the correctness of your answer.
 - Prove $(\forall a, d: a \in \mathbb{Z} \wedge d \in \mathbb{N}^{+2} : d | a \Rightarrow \neg(d | (a+1)))$.
 - Prove $(\forall a, b, d: a, b \in \mathbb{Z} \wedge d \neq 0 : d | a \vee d | b \Rightarrow d | (a*b))$.
 - Give a simple counter-example illustrating that Lemma 7.22 does *not* hold if p is *not* prime, for every $p \in \mathbb{N}^{+2}$.
- Prove the following properties of div and mod , using Definition 7.5; it is given that $d \in \mathbb{N}^+$ and that $a, b, x \in \mathbb{Z}$:
 - $0 \leq a < d \Leftrightarrow (a \bmod d) = a$
 - $0 \leq a < d \Leftrightarrow (a \text{ div } d) = 0$
 - $0 \leq a \Leftrightarrow 0 \leq a \text{ div } d$
 - $(a+d) \bmod d = a \bmod d$
 - $(a+d) \text{ div } d = (a \text{ div } d) + 1$
 - $(a+x*d) \bmod d = a \bmod d$

- (g) $(a+x*d) \operatorname{div} d = (a \operatorname{div} d) + x$
- (h) $(a \operatorname{mod} d) \operatorname{mod} d = a \operatorname{mod} d$
- (i) $(a \operatorname{mod} d) \operatorname{div} d = 0$
- (j) $(a+b) \operatorname{mod} d = (a \operatorname{mod} d + b \operatorname{mod} d) \operatorname{mod} d$
- (k) $(a+b) \operatorname{div} d = (a \operatorname{div} d) + (b \operatorname{div} d) + (a \operatorname{mod} d + b \operatorname{mod} d) \operatorname{div} d$
- (l) Give (simple) counter examples illustrating that $(a+b) \operatorname{mod} d$ is not necessarily equal to $a \operatorname{mod} d + b \operatorname{mod} d$, and that $(a+b) \operatorname{div} d$ is not necessarily equal to $a \operatorname{div} d + b \operatorname{div} d$.
- (m) $(a * b) \operatorname{mod} d = ((a \operatorname{mod} d) * (b \operatorname{mod} d)) \operatorname{mod} d$
- (n) $a \operatorname{mod} d = 0 \Leftrightarrow d|a$
- (o) $a \operatorname{mod} d = 0 \Leftrightarrow a \operatorname{div} d = a/d$
- (p) $1 \leq a \Leftrightarrow a \operatorname{div} d < a$, provided that $2 \leq d$
- (q) $a \operatorname{mod} d = b \operatorname{mod} d \Leftrightarrow (a-b) \operatorname{mod} d = 0$
- (r) Determine $(-1) \operatorname{div} d$ en $(-1) \operatorname{mod} d$
4. Given are $c, d \in \mathbb{N}^+$. Prove that for all $a \in \mathbb{Z}$:
- $$(a * d) \operatorname{mod} (c * d) = (a \operatorname{mod} c) * d \text{ and:}$$
- $$(a * d) \operatorname{div} (c * d) = a \operatorname{div} c .$$
5. (a) Determine the gcd of the numbers 112 and 280.
- (b) Determine numbers x and y satisfying: $x * 112 + y * 280 = \operatorname{gcd}(112, 280)$.
6. Prove Lemma 7.16.
7. (a) Construct an efficient implementation, in your favorite programming language, of Euclid's extended algorithm and prove its correctness.
- (b) Extend your program in such a way that it computes the least common multiple of the two numbers as well.
8. Determine all integer solutions x, y of: $21 * x + 15 * y = 33$.
9. Prove Lemma 7.26.
10. Determine, by hand, all prime numbers between 100 and 200.
11. Prove that $x * (x+1) * (x+2)$ is divisible by 6, for all $x \in \mathbb{Z}$.
12. Prove that $(x^2 - 1) \operatorname{mod} 8 \in \{0, 3, 7\}$, for all $x \in \mathbb{Z}$.
13. Resolve $8! - 3 * 7!$ into prime factors.
14. Solve $x \in \mathbb{Z}$ from: $12 * x =_{18} 30$.
15. Solve $x \in \mathbb{Z}$ from: $12 * x =_{18} 13$.

16. Solve $x \in \mathbb{Z}/34\mathbb{Z}$ from: $11 * x = 13$.
17. Determine the *lcm* of the numbers 1500000021 and 3000000045.
18. We consider the number whose decimal representation consists of 38 digits 1. We call this number X .
 - (a) Give a, formally correct, mathematical expression for X .
 - (b) What is the remainder of the division of x by 9?
 - (c) What is the remainder of the division of x by 11?
 - (d) What is the remainder of the division of x by 99?
19. Determine all values $x \in \mathbb{Z}$ satisfying both: $x = 2 \pmod{11}$ and: $x = 3 \pmod{23}$.
20. Solve $x \in \mathbb{Z}$ from the system: $x =_5 1$, $x =_{11} 2$ and: $x =_{23} 3$.
21. Resolve $\binom{17}{5}$ into prime factors.
22. Prove Lemma 7.27.
23. Determine the *gcd* and *lcm* of $\binom{17}{5}$ and $\binom{18}{4}$.
24. Prove that $(n+1) * (n+2) * \cdots * (n+k)$ is divisible by $k!$, for all $n, k \in \mathbb{N}$.
25. Prove that a natural number is divisible by 4 if and only if, in the representation of that number in base 17, the sum of the digits is divisible by 4.
26. We know that $37 * 43 = 1591$. Determine $e \in \mathbb{N}$ in such a way that, for all $m \in \mathbb{Z}$, we have: $m^{127 * e} =_{1591} m$.
27. Represent the number 1000 in base m , for every $m \in \{2, 3, 4, 5, 6, 7, 8, 9, 10\}$.
28. Calculate $(\sum i : 0 \leq i \leq n : 5^i * \binom{n}{i})$.
29. What is the period of the recurring decimal fraction that represents $1/1001001$?
30. Determine all prime numbers p with the property that $33 * p + 1$ is a square.
31. A given number requires 10 (ternary) digits for its ternary representation. How many digits are needed to represent this number in the decimal system?
32. Solve $x \in \mathbb{Z}$ from: $x =_{21} 5$ and: $x =_{28} 13$.