

Friedrich Eisenbrand

# Discrete Optimization

These are notes of my course Discrete Optimization. They are constantly updated.

May 5, 2025



# Contents

<b>1</b>	<b>Preface</b> .....	5
1.1	Optimization problems .....	5
1.2	Certificates .....	8
<b>2</b>	<b>Linear programming</b> .....	13
2.1	Softdrink production .....	13
2.2	Proving optimality .....	14
2.3	Linear Programs .....	16
2.4	Fitting a line .....	18
2.5	Linear Programming solvers and modeling languages .....	19
2.6	Linear programming for longer OLED-lifetime .....	20
<b>3</b>	<b>Polyhedra and convex sets</b> .....	25
3.1	Extreme points and vertices .....	26
3.2	Linear, affine, conic and convex hulls .....	28
3.3	Radon's lemma and Carathéodory's theorem .....	32
3.4	Separation theorem and Farkas' lemma .....	33
3.5	Decomposition theorem for polyhedra .....	34
<b>4</b>	<b>The simplex method</b> .....	43
4.1	Adjacent vertices .....	43
4.2	Bases, feasible bases and vertices .....	44
4.3	Moving to an improving vertex .....	46
4.4	Termination in the degenerate case .....	49
4.5	Finding an initial basic feasible solution .....	50
4.6	Removing degeneracy by perturbation .....	51
<b>5</b>	<b>Duality</b> .....	55
5.1	Zero sum games .....	57
5.2	A proof of the duality theorem via Farkas' lemma .....	60

<b>6</b>	<b>Algorithms and running time analysis</b>	63
6.1	Analysis of Gaussian elimination	68
6.2	Fast matrix multiplication	71
6.3	One iteration of the simplex algorithm	75
<b>7</b>	<b>Integer Programming</b>	77
7.1	Integral Polyhedra	79
7.2	Applications of total unimodularity	81
7.2.1	Bipartite matching	81
7.2.2	Bipartite vertex cover	82
7.2.3	Flows	83
7.2.4	Doubly stochastic matrices	83
7.3	The matching polytope	84
<b>8</b>	<b>Paths, cycles and flows in graphs</b>	89
8.1	Graphs	89
8.2	Representing graphs and computing the distance of two nodes	90
8.2.1	Breadth-first search	90
8.3	Shortest Paths	92
8.4	Maximum $s - t$ -flows	96
8.5	Minimum cost network flows, MCNFP	100
8.6	Computing a minimum cost-to-profit ratio cycle	110
8.6.1	Parametric search	112
<b>9</b>	<b>The ellipsoid method</b>	115
9.1	The method	118
9.2	Deciding feasibility	119
9.3	The separation problem	121
9.4	The ellipsoid method for optimization	123
9.5	Numerical issues	123
<b>10</b>	<b>Primal-Dual algorithm</b>	125
10.1	Graphs and Matchings	125
10.2	Matching problems	126
10.2.1	The maximum cardinality matching problem	127
10.2.2	The maximum weight matching problem	128
	<b>References</b>	137

# Chapter 1

## Preface

### 1.1 Optimization problems

An *optimization problem* is a pair  $(\mathcal{F}, f)$ , where  $\mathcal{F}$  is a set, referred to as the set of *feasible solutions*, and  $f : \mathcal{F} \rightarrow \mathbb{R}$  is a function, that is called the *objective* or *cost function*. The problem is to find an element  $x \in \mathcal{F}$  such that

$$f(x) \geq f(y) \text{ holds for all } y \in \mathcal{F}.$$

Such an  $x \in \mathcal{F}$  is then an *optimal solution* of the optimization problem. Since the task is to find a global maximum solution, we also write

$$\max\{c(x) : x \in \mathcal{F}\}.$$

If the objective is to minimize  $f(x)$  or, in other words, to find an  $x \in \mathcal{F}$  with

$$f(x) \leq f(y) \text{ for all } y \in \mathcal{F},$$

then this can be understood as the optimization problem  $(\mathcal{F}, -f)$ , since

$$-f(x) \geq -f(y) \text{ for all } y \in \mathcal{F} \text{ if and only if } f(x) \leq f(y) \text{ for all } y \in \mathcal{F}.$$

*Example 1.1.* From linear algebra we know the problem of maximizing a *quadratic form*  $x^T A x$  over the sphere  $S^{(d-1)} = \{x \in \mathbb{R}^d : \|x\| = 1\} \subseteq \mathbb{R}^d$ , where  $A \in \mathbb{R}^{d \times d}$  is a symmetric matrix. Thus, in this case,  $\mathcal{F} = S^{(d-1)}$  and  $f(x) = x^T A x$  and we write

$$\max\left\{x^T A x : x \in S^{(d-1)}\right\}.$$

If  $v \in \mathbb{R}^d \setminus \{0\}$  is an eigenvector of  $A$  corresponding to the maximal eigenvalue  $\lambda_{\max} \in \mathbb{R}$  of  $A$ , then  $v / \|v\|$  is an optimal solution of the optimization problem.

Following the discussion above, we see that an  $x^* \in S^{(d-1)}$  is an optimal solution of

$$\min \left\{ x^T A x : x \in S^{(d-1)} \right\}$$

if and only if  $x^*$  is an optimal solution of

$$\max \left\{ x^T (-A) x : x \in S^{(d-1)} \right\}.$$

We have seen that such an optimal solution is an eigenvector of  $A$  corresponding to the minimal eigenvalue  $\lambda_{\min}$  of  $A$ .

*Example 1.2.* A *directed graph* is a tuple  $G = (V, A)$ , where  $V$  is a finite set of elements, called the *vertices* of  $G$  and  $A \subseteq (V \times V)$  is the set of *arcs* of  $G$ . We denote an arc by its two defining nodes  $(u, v) \in A$ . The graph is *weighted* if additionally equipped with a function  $w : A \rightarrow \mathbb{R}$  that maps the set of arcs to the reals. A *simple path* is a finite sequence of *distinct* vertices  $P = v_1, v_2, \dots, v_k$  such that  $(v_i, v_{i+1}) \in A$  for each  $i \in \{1, \dots, k-1\}$ . The path is from  $v_1$  to  $v_k$  and the *length* of the path  $P$  is defined as

$$\ell(P) = \sum_{i=1}^{k-1} w(v_i, v_{i+1}).$$

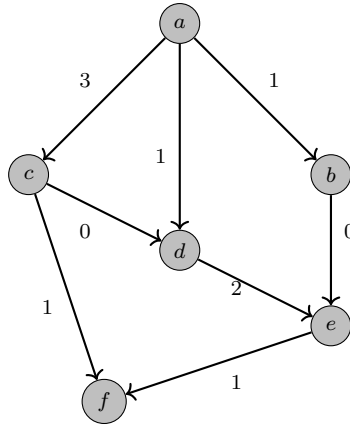


Fig. 1.1: Example of a weighted directed graph with 6 nodes and 8 arcs. A shortest path from  $a$  to  $f$  is the node sequence  $a, b, e, f$ . Its length is equal to 2.

The *shortest path problem* is now as follows. Given a directed graph  $G = (V, A)$  with non-negative weights  $w : A \rightarrow \mathbb{R}_+$  and two designated vertices  $s$  and  $t \in V$ , determine a path from  $s$  to  $t$  of minimal length. In our setting this is then the tuple  $(\mathcal{P}, f)$ , where  $\mathcal{P}$  is the set of all paths from  $s$  to  $t$  in  $G$  and  $f(P)$  is equal to  $-1$  times  $\ell(P)$ .

## Linear programming

Central to this course is the *linear programming* or *linear optimization* problem. This is an optimization problem  $(\mathcal{F}, f)$  where the set  $\mathcal{F}$  is described by *linear inequalities*

$$\mathcal{F} = \{x \in \mathbb{R}^n : Ax \leq b\},$$

where  $A \in \mathbb{R}^{m \times n}$  is a matrix and  $b \in \mathbb{R}^m$  is a vector. The objective function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is linear and of the form  $f(x) = c^T x$  for a given vector  $c \in \mathbb{R}^n$ .

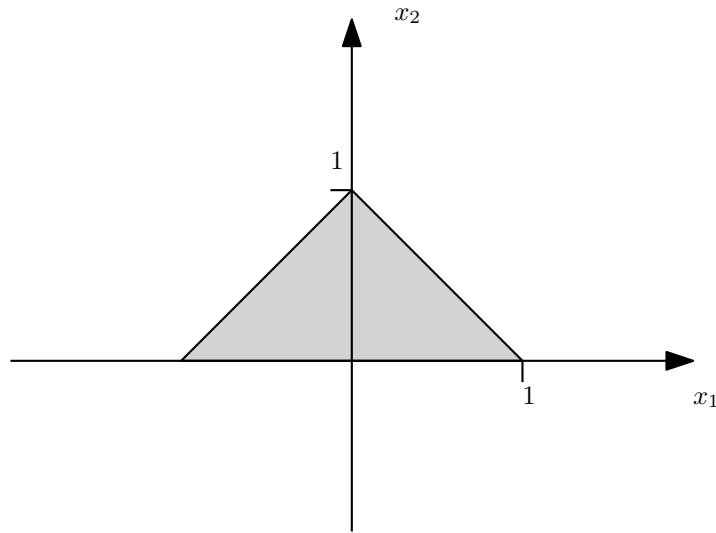


Fig. 1.2: The set  $\mathcal{F}$  in example 1.3

*Example 1.3.* Let  $A \in \mathbb{R}^{3 \times 2}$  be the matrix

$$A = \begin{pmatrix} -1 & 1 \\ 1 & 1 \\ 0 & -1 \end{pmatrix}$$

and  $b \in \mathbb{R}^3$  be the vector

$$b = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$$

The set  $\mathcal{F} \subseteq \mathbb{R}^2$  is the set of all points  $(x_1, x_2)^T \in \mathbb{R}^2$  in the plane that satisfy all three inequalities

$$-x_1 + x_2 \leq 1, \quad x_1 + x_2 \leq 1, \quad \text{and} \quad x_2 \geq 0.$$

This set  $\mathcal{F}$  is the triangle depicted in Figure 1.2. If we set the objective function as  $f(x) = x_2$ , then  $x = (0, 1)^T$  is an optimal solution.

The linear optimization problem is one of the most important types of mathematical optimization problems. It has numerous applications in science and engineering and, in particular, modern fields like *machine learning*.

In this course, we will learn the theory of linear optimization and develop *efficient algorithms* to solve linear programming problems. This means that we do not content ourselves with an algorithm that is correct, but we want this algorithm to be capable to solve large scale problems within a short time.

## 1.2 Certificates

A central topic in linear algebra is the theory around linear equations

$$Ax = b \quad (1.1)$$

where  $A \in \mathbb{R}^{m \times n}$  is a given matrix and  $b \in \mathbb{R}^m$  is a given vector. Every student of mathematics learns to appreciate *Gaussian elimination* which transforms the system (1.1) into an equivalent system

$$A'x = b'$$

that is in row-echelon form. This means that  $A' \in \mathbb{R}^{m \times n}$  is such that for each  $i \in \{1, \dots, m-1\}$

$$\min\{j : a'_{ij} \neq 0\} < \min\{j : a'_{i+1j} \neq 0\},$$

see Figure 1.3.

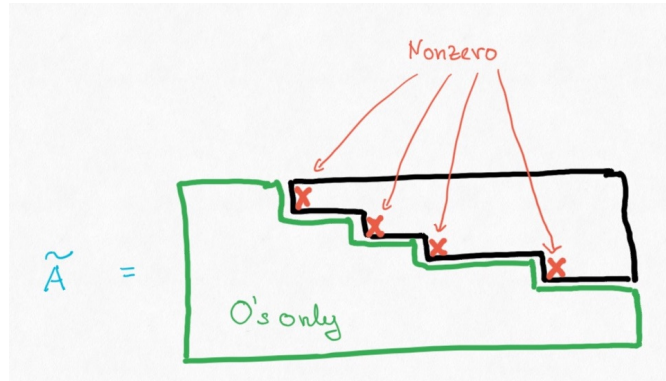


Fig. 1.3: A schematic image of a matrix in row-echelon form



The system (1.1) has a solution if and only if  $b'$  has only zero's in those components that correspond to the zero-rows of  $A'$  and a solution is readily computed. In fact, Gaussian elimination also provides an invertible matrix  $Q \in \mathbb{R}^{m \times m}$  such that  $Q \cdot A = A'$  and  $Q \cdot b = b'$ . If the  $i$ -th component of  $b'$  is nonzero and the  $i$ -th row of  $A'$  consists of zeros only, then with  $q$  being the column vector corresponding to the  $i$ -th row of  $Q$  one has

$$q^T A = 0 \text{ and } q^T b \neq 0.$$

Thus there is a convenient *certificate* of the fact that (1.1) is not solvable.

**Theorem 1.1.** *A linear system (1.1) is not solvable if and only if there exists a vector  $q \in \mathbb{R}^m$  such that*

$$q^T A = 0 \text{ and } q^T b \neq 0.$$

In this course, we will now also deal with systems of *linear inequalities*

$$Ax \leq b, \tag{1.2}$$

where  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ . A *solution* to such a system is a vector  $x^* \in \mathbb{R}^n$  such that

$$\text{for each } i = 1, \dots, m \quad : \quad a_{i1}x_1^* + \dots + a_{in}x_n^* \leq b_i. \tag{1.3}$$

If a solution exists, the system (1.2) is called *feasible* otherwise it is called *infeasible*. We will ask analogous questions for systems of linear inequalities as we did for linear equations. Can we efficiently find a solution of (1.2)? Is there a simple certificate to convince somebody of the infeasibility of (1.2)?

To shed a bit of light on this second question, let us consider a vector  $\lambda \in \mathbb{R}_{\geq 0}^m$ . If  $x^*$  is a solution, then  $x^*$  also satisfies the inequality

$$\lambda^T A x \leq \lambda^T b. \tag{1.4}$$

If there exists a  $\lambda \in \mathbb{R}_{\geq 0}^m$  such that

$$\lambda^T A = 0^T \text{ and } \lambda^T b = -1, \tag{1.5}$$

then this  $\lambda$  *certifies* that (1.2) is infeasible. In the case of infeasibility, can such a  $\lambda \geq 0$  always be found? Among the many results presented in this course, we will show that the answer is “yes”.

**Theorem 1.2 (Farkas' lemma).** *A system of linear inequalities (1.2) is infeasible if and only if there exists a  $\lambda \in \mathbb{R}_{\geq 0}^m$  such that*

$$\lambda^T A = 0 \text{ and } \lambda^T b = -1.$$

More generally, we will learn about certificates of optimality of linear optimization problems. Lets consider again example 2.1. How can we certify that  $x^* = (0, 1)^T$  is an optimal solution of the linear program.

Each point  $x^*$  in  $\mathcal{F}$  must satisfy the first two linear inequalities

$$-x_1 + x_2 \leq 1 \text{ and } x_1 + x_2 \leq 1.$$

But then such an  $x^*$  satisfies also the *sum* of these two inequalities

$$(-x_1 + x_2) + (x_1 + x_2) \leq 1 + 1$$

which simplifies to

$$x_2 \leq 1.$$

But  $x_2$  is the objective function. We have just shown that the objective function value of any point is at most 1. Since this objective function value at  $x^* = (0, 1)^T$  is *equal* to one, we have certified optimality of  $x^*$ .

The general theory that makes certification of optimality possible for linear optimization is the theory of *duality*. It is a rich theory with many other applications in discrete mathematics. We will touch upon many of them.

## Exercises

- 1) Provide a certificate as in Theorem 1.1 of the unsolvability of the linear equation

$$\begin{pmatrix} 2 & 1 & 0 \\ 5 & 4 & 1 \\ 7 & 5 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 4 \end{pmatrix}$$

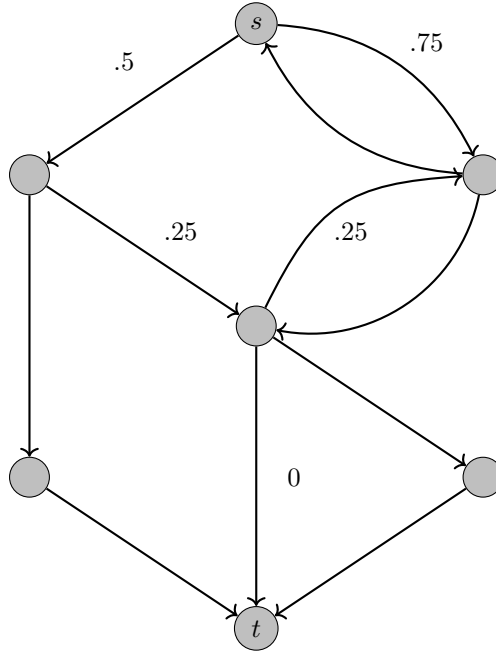
- 2) Show the “if” direction of the Farkas’ lemma (Theorem 1.2).
- 3) The directed graph  $G = (V, A)$  is *complete*, if  $A = \{(u, v) : u \neq v \in V\}$ . Let  $s \neq t \in V$  be two designated vertices. How many directed paths connect  $s$  and  $t$  in  $G$ ? Find a formula with parameter  $n = |V|$ .
- 4) Let  $G = (V, A)$  be a directed graph and  $s, t \in V$  be two designated vertices. For a vertex  $v \in V$  we let

$$\delta^+(v) = \{(u, v) : u \in V, (u, v) \in A\} \text{ and } \delta^-(v) = \{(v, u) : u \in V, (v, u) \in A\}$$

the *arcs entering* and *leaving*  $u$  respectively. Consider the following inequalities

$$\begin{aligned}
\sum_{a \in \delta^+(v)} x_a - \sum_{a \in \delta^-(v)} x_a &= 0 & v \in V \setminus \{s, t\} \\
\sum_{a \in \delta^+(s)} x_a - \sum_{a \in \delta^-(s)} x_a &= -1 \\
\sum_{a \in \delta^+(t)} x_a - \sum_{a \in \delta^-(t)} x_a &= 1 \\
x_a &\geq 0 & a \in A.
\end{aligned} \tag{1.6}$$

- a) Consider the following digraph with  $s$  and  $t$  and a partial assignment of arc variables. Can this partial assignment be completed to a feasible solution satisfying the inequalities (1.6)? If yes, complete the assignment.



- b) Show the following for a digraph  $G = (V, A)$  with  $s, t \in V$ : If there is a path connecting  $s$  and  $t$  in  $G$ , then the system of inequalities (1.6) has a feasible solution
- c) (\*) Show the following for a digraph  $G = (V, A)$  with  $s, t \in V$ : If the system of inequalities (1.6) has a feasible solution, then there is a path connecting  $s$  and  $t$  in  $G$ .



## Chapter 2

# Linear programming

We start by giving some examples of linear programs and how they are used in practice.

### 2.1 Softdrink production

Imagine that you own a company that produces the two softdrinks, *Spring* and *Nebsi*. These softdrinks are a mixture of *water*, an ingredient *A* and ingredient *B*. The recipes for Spring and Nebsi are different. Also, the profit for the two drinks is not the same. Those are as follows.

The use of ingredients A and B and the profit per 100l are as follows.

	A	B	Profit
Spring	3l	8l	100 CHF
Nebsi	6l	4l	125 CHF

While the supply of water is unlimited, your company has only 30l of ingredient *A* and 44l of ingredient *B*.

At the end of the production day, the local wholesaler picks up the drinks in two barrels. The capacity of the barrel for Spring is 500l while the barrel for Nebsi has a capacity of 400l.

As the manager of your small company, your goal is to come up with a *production plan* that maximizes your profit. A production plan is a two-dimensional vector  $(x_1, x_2) \in \mathbb{R}^2$  which means that you will produce  $x_1 \cdot 100l$  of Spring and  $x_2 \cdot 100l$  of Nebsi. A production plan is *feasible* if the produced drinks fit into the respective barrels and not more of A and B is used than what is on stock. Clearly,  $(5, 4)$  is not a feasible production plan, as this would require 39l of ingredient A which exceeds the capacity.

A feasible production plan that maximizes your profit can be found with the help of a *linear program*, a central object of study in this course.

$$\begin{aligned}
& \max. && 100 \cdot x_1 + 125 \cdot x_2 \\
& \text{s.t.} && 3 \cdot x_1 + 6 \cdot x_2 \leq 30 \\
& && 8 \cdot x_1 + 4 \cdot x_2 \leq 44 \\
& && x_1 \leq 5 \\
& && x_2 \leq 4 \\
& && x_1 \geq 0 \\
& && x_2 \geq 0
\end{aligned} \tag{2.1}$$

One has to maximize a linear *objective function*, in this case  $f(x_1, x_2) = 100 \cdot x_1 + 125 \cdot x_2$  where  $(x_1, x_2)$  satisfies linear inequalities. The linear inequalities  $x_1 \geq 0$  and  $x_2 \geq 0$  reflect the fact that only positive amounts can be produced, while  $x_1 \leq 5$  reflects the barrel capacity for Spring. The linear inequality  $3 \cdot x_1 + 6 \cdot x_2 \leq 30$  reflects the amount of 30l of ingredient A that is on stock.

We can now make a drawing of all feasible production plans.

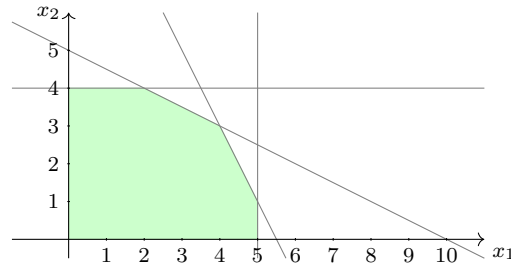


Fig. 2.1: The feasible production plans are the green area.

The set of points  $(x_1, x_2)$  that have objective function  $\beta$  is the line

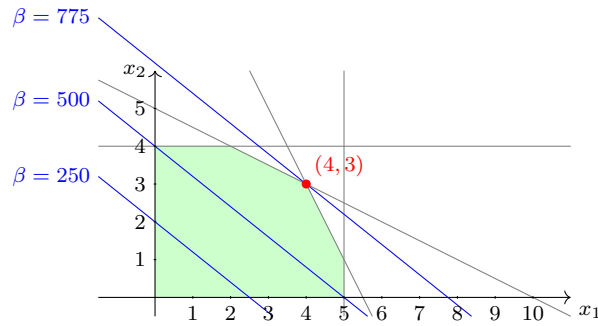
$$\{(x_1, x_2) \in \mathbb{R}^2 : 100 \cdot x_1 + 125 \cdot x_2 = \beta\}$$

and our task is now to find the largest value for  $\beta$  such that the corresponding line still intersects the set of feasible production plans.

Figure 2.2 reveals that  $(4, 3)$  is an optimal production plan and that the maximum profit that the manager can achieve is 775.

## 2.2 Proving optimality

How can the manager be convinced that  $(4, 3)$  is an optimal production plan? Maybe he has made a mistake in his drawing or with his calculations and  $(4, 3)$  is not optimal. What we will see now is a very important principle of linear programming. There is a simple way to prove optimality of solutions that we will explore later on.

Fig. 2.2: The optimal production plan is  $(4, 3)$ .

Inspecting the drawing, one can see that there are two inequalities that  $(4, 3)$  satisfies with equality, namely the inequalities

$$3 \cdot x_1 + 6 \cdot x_2 \leq 30 \quad (2.2)$$

$$8 \cdot x_1 + 4 \cdot x_2 \leq 44. \quad (2.3)$$

Clearly all feasible production plans satisfy these inequalities and inspecting Figure 2.2 it seems clear that  $(4, 3)$  is an optimal solution of the optimization problem (2.1) where each linear inequality but the inequalities (2.2) and (2.3) have been removed.

What now follows is a very important technique that we will apply later on again in greater generality. Since each feasible production plan satisfies the inequalities (2.2) and (2.3) it satisfies also these inequalities, after they have been multiplied by  $50/3$  and  $25/4$  respectively. In fact the inequalities (2.2) and (2.3) are equivalent to the following two inequalities

$$50 \cdot x_1 + 100 \cdot x_2 \leq 500 \quad (2.4)$$

$$50 \cdot x_1 + 25 \cdot x_2 \leq 275. \quad (2.5)$$

By adding up the inequalities (2.4) and (2.5) we obtain the inequality

$$100 \cdot x_1 + 125 \cdot x_2 \leq 775 \quad (2.6)$$

which in turn is also satisfied by each feasible production plan. The left-hand-side of inequality (2.6) is the objective function and 775 is the value of the objective function evaluated at  $(4, 3)$ . Thus each feasible production plan yields a profit of at most 775 which is the profit yielded by  $(4, 3)$ . This shows that  $(4, 3)$  is optimal.

## 2.3 Linear Programs

We use the following notation. For a matrix  $A \in \mathbb{R}^{m \times n}$ ,  $i \in \{1, \dots, m\}$  and  $j \in \{1, \dots, n\}$  we denote the  $i$ -th row of  $A$  by  $a_i^T$  and the  $j$ -th column of  $A$  by  $a^j$ . With  $a_{ij}$  we denote the element of  $A$  which is in the  $i$ -th row and  $j$ -th column of  $A$ . For a vector  $v \in \mathbb{R}^m$  and  $i \in \{1, \dots, m\}$  we denote the  $i$ -th element of  $v$  by  $v_i$ .

**Definition 2.1.** Let  $A \in \mathbb{R}^{m \times n}$  be a matrix,  $b \in \mathbb{R}^m$  and  $c \in \mathbb{R}^n$  be vectors and  $I_{\geq}, I_{\leq}, I_{=} \subseteq \{1, \dots, m\}$  and  $J_{\geq}, J_{\leq} \subseteq \{1, \dots, n\}$  be index sets. A *linear program (LP)* consists of

i) a linear *objective function*

$$\begin{aligned} & \max c^T x \\ \text{or } & \min c^T x \end{aligned}$$

ii) linear *constraints*

$$\begin{aligned} a_i^T x & \geq b_i, i \in I_{\geq} \\ a_j^T x & \leq b_j, j \in I_{\leq} \\ a_k^T x & = b_k, k \in I_{=} \end{aligned}$$

iii) and *bounds on the variables*

$$\begin{aligned} x_j & \geq 0, j \in J_{\geq} \\ x_j & \leq 0, j \in J_{\leq}. \end{aligned}$$

Notice that we can re-write the objective function  $\min c^T x$  as  $\max -c^T x$ . Similarly, the constraints  $a_i^T x \geq b_i, i \in I_{\geq}$  are equivalent to the constraints  $-a_i^T x \leq -b_i, i \in I_{\geq}$ . Also the constraints  $a_k^T x = b_k, k \in I_{=}$  can be replaced by the constraints  $a_k^T x \leq b_k, -a_k^T x \leq -b_k, k \in I_{=}$ . A lower bound  $x_j \geq 0$  can be written as  $-e_j^T x \leq 0$ , where  $e_j$  is the  $j$ -th unit vector which has zeroes in every component, except for the  $j$ -th component, which is 1. Similarly an upper bound  $x_j \leq 0$  can be written as  $e_j^T x \leq 0$ .

All-together, a linear program as in Definition 2.1 can always be written as

$$\max\{c^T x: \tilde{A}x \leq \tilde{b}, x \in \mathbb{R}^n\}$$

with a suitable matrix  $\tilde{A} \in \mathbb{R}^{m \times n}$  and a suitable vector  $\tilde{b} \in \mathbb{R}^m$ . This representation has a name.

**Definition 2.2.** A linear program is in *inequality standard form*, if it is of the form

$$\max\{c^T x: Ax \leq b, x \in \mathbb{R}^n\}$$

for some matrix  $A \in \mathbb{R}^{m \times n}$  and some vector  $b \in \mathbb{R}^m$ .

*Example 2.1.* Let us convert the following linear program to an equivalent linear program in inequality standard form. The objective is to minimize



$$2x_1 + 5x_2 - 8x_3$$

such that  $x_1, x_2, x_3$  satisfy the constraints

$$x_1 + x_3 \geq 6, \quad (2.7)$$

$$-x_1 + 3x_2 - 5x_3 = -4, \quad (2.8)$$

as well as the lower bounds

$$x_1 \geq 0, x_2 \geq 0, x_3 \leq 0. \quad (2.9)$$

The constraint (2.15) is equivalent to

$$-x_1 - x_3 \leq -6$$

and (2.16) is equivalent to the conjunction of the two constraints

$$\begin{aligned} -x_1 + 3x_2 - 5x_3 &\leq -4 \\ x_1 - 3x_2 + 5x_3 &\leq 4. \end{aligned}$$

The lower bounds  $x_1 \geq 0, x_2 \geq 0$  in (2.17) can be re-written as  $-x_1 \leq 0$  and  $-x_2 \leq 0$ . Finally, by multiplying the objective function by  $-1$  we can transfer to a maximization problem and obtain the linear program

$$\max(-2, -5, 8)x$$

subject to  $x \in \mathbb{R}^3$  satisfying

$$A = \begin{pmatrix} -1 & 0 & -1 \\ -1 & 3 & -5 \\ 1 & -3 & 5 \\ -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{pmatrix} x \leq \begin{pmatrix} -6 \\ -4 \\ 4 \\ 0 \\ 0 \\ 0 \end{pmatrix}.$$

**Definition 2.3.** A point  $x^* \in \mathbb{R}^n$  is called *feasible*, if  $x^*$  satisfies all constraints and bounds on the variables. If there are feasible solutions of a linear program, then the linear program is called *feasible* itself. A linear program is *bounded* if there exists a constant  $M \in \mathbb{R}$  such for all feasible  $x^* \in \mathbb{R}^n$   $c^T x^* \leq M$ , if the linear program is a maximization problem and  $c^T x^* \geq M$ , if the linear program is a minimization problem. A feasible solution  $x^*$  is an optimal solution if  $c^T x^* \geq c^T y^*$  for all feasible  $y^*$  if the linear program is a maximization problem and  $c^T x^* \leq c^T y^*$  if the linear program is a minimization problem.

We will see later that a feasible and bounded linear program has an optimal solution.

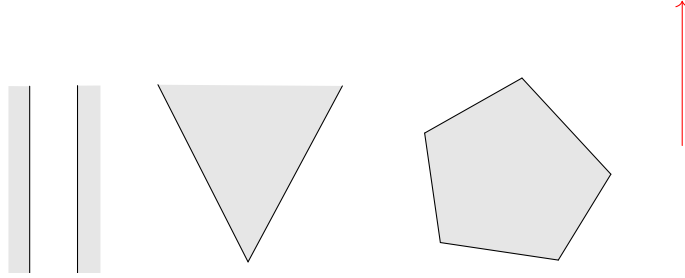


Fig. 2.3: With the objective function being to find the highest point, we have from left-to-right an infeasible linear program, an unbounded linear program and a bounded linear program.

## 2.4 Fitting a line

The following is an example which is well known in statistics. Suppose that you measure points  $(x_i, y_i) \in \mathbb{R}^2$   $i = 1, \dots, n$  and you are interested in a linear function  $y = a \cdot x + b$  that reflects the sample. One way to do that is by minimizing the expression

$$\sum_{i=1}^n (ax_i + b - y_i)^2, \quad (2.10)$$

where  $a, b \in \mathbb{R}$  are the parameters of the line that we are looking for. The number  $(ax_i + b - y_i)^2$  is the square of the vertical distance of the point  $(x_i, y_i)$  from the line  $y = ax + b$ .

Instead of using the method of least-squares, we could also minimize the following function, see also [13, Chapter 2.4],

$$\sum_{i=1}^n |ax_i + b - y_i|. \quad (2.11)$$

This objective has the advantage to be slightly more robust towards outliers. How can we model this as a linear program. The trick is to use an extra variable  $h_i$  which models the absolute value of  $ax_i + b - y_i$ .

$$\begin{aligned} \min \quad & \sum_{i=1}^n h_i \\ h_i \quad & \geq \quad ax_i + b - y_i, \quad i = 1, \dots, n \\ h_i \quad & \geq \quad -(ax_i + b - y_i), \quad i = 1, \dots, n \end{aligned} \quad (2.12)$$

The variables of this linear program are  $h_i$ ,  $i = 1, \dots, n$ ,  $a$  and  $b$ . For a fixed  $a \in \mathbb{R}$  and  $b \in \mathbb{R}$  the optimal  $h_i$ 's will be  $h_i = |ax_i + b - y_i|$  since the objective minimizes the sum of the  $h_i$ 's. If one of the  $h_i$ 's was strictly larger than  $|ax_i + b - y_i|$ , then the objective could be improved by making it smaller.

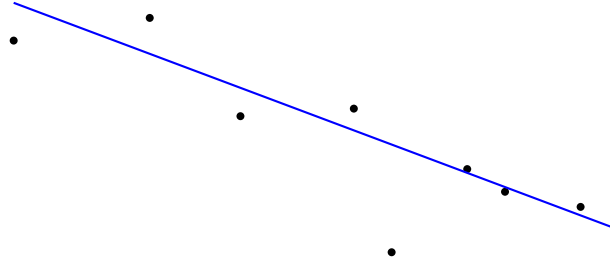


Fig. 2.4: A line that minimizes the sum of the vertical distances.

## 2.5 Linear Programming solvers and modeling languages

We will demonstrate now how to use a modeling language for linear programming and a linear programming solver to find a fitting line, as described in Section 2.4 for the points

$$(1, 3), (2.8, 3.3), (4, 2), (5.5, 2.1), (6, 0.2), (7, 1.3), (7.5, 1), (8.5, 0.8)$$

There are two popular formats for linear programming problems which are widely used by linear programming solvers, the *lp-format* and the *mps-format*. Both are not easy to read. To facilitate the modeling of a linear program, so-called modeling languages are used. We demonstrate the use of the popular open source modeling software called *zimpl* [8]. Below you see a way to model our fitting line linear program with *zimpl*:

```

set I := {1 to 8};
param X[I] := <1> 1, <2> 2.8, <3> 4 , <4> 5.5,
              <5> 6, <6> 7 , <7> 7.5, <8> 8.5 ;
param Y[I] := <1> 3 , <2> 3.3, <3> 2, <4> 2.1,
              <5> 0.2, <6> 1.3, <7> 1, <8> 0.8 ;
var h[I] >= -infinity <= infinity;
var a >= -infinity <= infinity ;
var b >= -infinity <= infinity ;

minimize cost: sum <i> in I: h[i];

subto c1: forall <i> in I:      h[i] >= ( a * X[i] + b -Y[i]);
subto c2: forall <i> in I:      h[i] >= - ( a * X[i] + b -Y[i]);

```

Zimpl creates a linear program which is readable by linear programming solvers like QSOPT or SoPLEX.

## 2.6 Linear programming for longer OLED-lifetime

*Organic Light Emitting Diodes* (OLEDs) are considered as the display technology of the future and more and more commercial products are equipped with such displays as shown in Fig. 2.5. However, the cheapest OLED technology suffers from short lifetimes. We will show in this section how linear programming can be used to increase the lifetime of such displays.

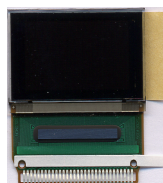


Fig. 2.5: Sample of a commercial OLED device with integrated driver chip

A (passive matrix) OLED display has a matrix structure consisting of  $n$  rows and  $m$  columns. At any crossover between a row and a column there is a vertical diode which works as a pixel. The image itself is given as an integral non-negative  $n \times m$  matrix  $(r_{ij}) \in [0, \dots, \varrho]^{n \times m}$  representing its RGB values. Consider the contacts for the rows and columns as switches. For the time the switch of row  $i$  and column  $j$  is closed, an electrical current flows through the diode of pixel  $(i, j)$  and it shines. Hence, we can control the intensity of a pixel by the two quantities *electrical current* and *time*. The value  $r_{ij}$  determines the amount of time within the time frame in which the switches  $i$  and  $j$  have to be simultaneously closed. At a sufficient high frame rate e.g. 50 Hz, the perception by the eye is the average value of the light emitted by the pixel and one sees the image.

The traditional addressing scheme is row-by-row. This means that the switch for the first row is closed for a certain time while the switches for the columns are closed for the necessary amount of time dictated by the entries  $r_{1j}$ ,  $j = 1, \dots, m$ . Consequently the first row can be displayed in time  $\max\{r_{1j} : j = 1, \dots, m\}$ . Then the second row is displayed and so on. With this addressing scheme, the pixels are idle most of the time and then have to shine with very high intensity. This puts the diodes under stress and is a major cause of the short lifetime of the displays.

How can this lifetime problem be dealt with? The main idea is to save time, or equivalently to lower the maximum intensity, by displaying several rows at once.

Consider the schematic image on the left of Fig. 2.6. Let us compute the amount of time which is necessary to display the image with this addressing scheme. The maximum value of the entries in the first row is 238. This is the amount of time which is necessary to display the first row. After that the second row is displayed in time 237. In total the time which is required to display the image is  $238 + 237 + 234 + 232 + 229 = 1170$  time units.

$$\begin{array}{|c|c|c|} \hline 109 & 238 & 28 \\ \hline 112 & 237 & 28 \\ \hline 150 & 234 & 25 \\ \hline 189 & 232 & 22 \\ \hline 227 & 229 & 19 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 82 & 25 \\ \hline 0 & 82 & 25 \\ \hline 0 & 41 & 22 \\ \hline 0 & 41 & 22 \\ \hline 0 & 0 & 0 \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 112 & 155 & 3 \\ \hline 112 & 155 & 3 \\ \hline 189 & 191 & 0 \\ \hline 189 & 191 & 0 \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline 109 & 156 & 3 \\ \hline 0 & 0 & 0 \\ \hline 38 & 38 & 0 \\ \hline 0 & 0 & 0 \\ \hline 38 & 38 & 19 \\ \hline \end{array}$$

Fig. 2.6: An example decomposition

Now consider the decomposition of the image as the sum of the three images on the right of Fig. 2.6. In the first image, each odd row is equal to its even successor. This means that we can close the switches for rows 1 and 2 simultaneously, and these two equal rows are displayed in 82 time units. Rows 3 and 4 can also be displayed simultaneously which shows that the first image on the right can be displayed in  $82 + 41$  time units. The second image on the right can be displayed in  $155 + 191$  time units while the third image has to be displayed traditionally. In total all three images, and thus the original image on the left via this decomposition, can be displayed in  $82 + 41 + 155 + 191 + 156 + 38 + 38 = 701$  time units. This means that we could reduce the necessary time via this decomposition by roughly 40%. We could equally display the image in the original 1170 time units but reduce the peak intensity, or equally the maximum electrical current through a diode by roughly 40%.

We now show how to model the time-optimal decomposition of an image as a linear program. To decompose  $R$  we need to find matrices  $F^{(1)} = (f_{ij}^{(1)})$  and  $F^{(2)} = (f_{ij}^{(2)})$  where  $F^{(1)}$  represents the singleline part and  $F^{(2)}$  the two doubleline parts. More precisely, the  $i$ -th row of matrix  $F^{(2)}$  represents the doubleline covering rows  $i$  and  $i + 1$ . Since the overlay (addition) of the subframes must be equal to the original image to get a valid decomposition of  $R$ , the matrices  $F^{(1)}$  and  $F^{(2)}$  must fulfill the constraint  $f_{ij}^{(1)} + f_{i-1,j}^{(2)} + f_{ij}^{(2)} = r_{ij}$  for  $i = 1, \dots, n$  and  $j = 1, \dots, m$ , where we now and in the following use the convention to simply omit terms with indices running out of bounds. Since we cannot produce “negative” light we require also non-negativity of the variables  $f_{ij}^{(\alpha)} \geq 0$ . The goal is to find an integral decomposition that

minimizes

$$\sum_{i=1}^n \max\{f_{ij}^{(1)} : 1 \leq j \leq m\} + \sum_{i=1}^{n-1} \max\{f_{ij}^{(2)} : 1 \leq j \leq m\} .$$

This problem can be formulated as a linear program by replacing the objective by  $\sum_{i=1}^n u_i^{(1)} + \sum_{i=1}^{n-1} u_i^{(2)}$  and by adding the constraints  $f_{ij}^{(\alpha)} \leq u_i^{(\alpha)}$ . This yields

$$\begin{aligned} \min \quad & \sum_{i=1}^n u_i^{(1)} + \sum_{i=1}^{n-1} u_i^{(2)} \\ \text{s.t.} \quad & f_{ij}^{(1)} + f_{i-1,j}^{(2)} + f_{ij}^{(2)} = r_{ij} \quad \text{for all } i, j \end{aligned} \quad (2.13)$$

$$f_{ij}^{(\alpha)} \leq u_i^{(\alpha)} \quad \text{for all } i, j, \alpha \quad (2.14)$$

$$f_{ij}^{(\alpha)} \in \mathbb{R}_{\geq 0} \quad \text{for all } i, j, \alpha$$

Note that the objective does not contain the  $f$ -variables. By decomposing images like this, the average lifetime of an OLED display can be increased by roughly 100%, see [4].

## Exercises

- 1) A company produces and sells two different products. Our goal is to determine the number of units of each product they should produce during one month, assuming that there is an unlimited demand for the products, but there are some constraints on production capacity and budget. There are 20000 hours of machine time in the month. Producing one unit takes 3 hours of machine time for the first product and 4 hours for the second product. Material and other costs for producing one unit of the first product amount to 3CHF, while producing one unit of the second product costs 2CHF. The products are sold for 6CHF and 5CHF per unit, respectively. The available budget for production is 4000CHF initially. 25% of the income from selling the first product can be used immediately as additional budget for production, and so can 28% of the income from selling the second product.
  - a. Formulate a linear program to maximize the profit subject to the described constraints.
  - b. Solve the linear program graphically by drawing its set of feasible solutions and determining an optimal solution from the drawing.
  - c. Suppose the company could modernize their production line to get an additional 2000 machine hours for the cost of 400CHF. Would this investment pay off?

- 2) Reformulate the following linear program in inequality standard form. The objective is to minimize

$$-3x_1 + 3x_2 + 5x_3$$

such that  $x_1, x_2, x_3$  satisfy the constraints

$$x_2 + 3x_3 \geq 4, \quad (2.15)$$

$$2x_1 + 2x_2 - 4x_3 = -4, \quad (2.16)$$

as well as the lower bounds

$$x_1 \geq 0, x_2 \geq 0, x_3 \geq 0. \quad (2.17)$$

- 3) A factory produces two different products. To create one unit of product 1, it needs one unit of raw material *A* and one unit of raw material *B*. To create one unit of product 2, it needs one units of raw material *B* and two units of raw material *C*. Raw material *B* needs preprocessing before it can be used, which takes one minute per unit. At most 20 hours of time is available per day for the preprocessing. Raw materials of capacity at most 1200 can be delivered to the factory per day. One unit of raw material *A*, *B* and *C* has size 4, 3 and 2 respectively.

At most 130 units of the first and 100 units of the second product can be sold per day. The first product sells for 6 CHF per unit and the second one for 9 CHF per unit.

Formulate the problem of maximizing turnover as a linear program in two variables and solve it.

- 4) Prove the following statement or give a counterexample: The set of optimal solutions of a linear program is always finite.  
 5) Let (2.18) be a linear program in inequality standard form, i.e.

$$\max\{c^T x \mid Ax \leq b, x \in \mathbb{R}^n\} \quad (2.18)$$

where  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ , and  $c \in \mathbb{R}^n$ .

Prove that there is an equivalent linear program (2.19) of the form

$$\max\{\tilde{c}^T x \mid \tilde{A}x = \tilde{b}, x \geq 0, x \in \mathbb{R}^{\tilde{n}}\} \quad (2.19)$$

where  $\tilde{A} \in \mathbb{R}^{\tilde{m} \times \tilde{n}}$ ,  $\tilde{b} \in \mathbb{R}^{\tilde{m}}$ , and  $\tilde{c} \in \mathbb{R}^{\tilde{n}}$  are such that every feasible point of (2.18) corresponds to a feasible point of (2.19) with the same objective function value and vice versa.

Linear programs of the form in (2.19) are said to be in *equality standard form*.

- 6) Model the linear program (2.13) to decompose the EPFL logo with Zimpl. An incomplete model containing the encoding of the grayscale values of the logo can be found here<sup>1</sup>. Use an LP solver library of your choice to compute an optimal solution.

---

<sup>1</sup> [http://disopt.epfl.ch/webdav/site/disopt/users/190205/public/logo\\_dec.zmpl](http://disopt.epfl.ch/webdav/site/disopt/users/190205/public/logo_dec.zmpl)



## Chapter 3

# Polyhedra and convex sets

**Definition 3.1.** A polyhedron  $P \subseteq \mathbb{R}^n$  is a set of the form  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  for some  $A \in \mathbb{R}^{m \times n}$  and some  $b \in \mathbb{R}^m$ .

We are interested in polyhedra, since the set of feasible solutions of a linear program  $\max\{c^T x : Ax \leq b\}$  is a polyhedron.

*Example 3.1.* Consider again the soft-drink production problem from chapter 2.1. The corresponding set of feasible solutions is the polyhedron  $P = \{x \in \mathbb{R}^2 : Ax \leq b\}$  with

$$A = \begin{pmatrix} 3 & 6 \\ 8 & 4 \\ 1 & 0 \\ 0 & 1 \\ -1 & 0 \\ 0 & -1 \end{pmatrix}, \quad \text{and} \quad b = \begin{pmatrix} 30 \\ 44 \\ 5 \\ 4 \\ 0 \\ 0 \end{pmatrix}.$$

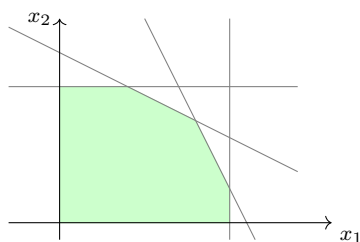


Fig. 3.1: The polyhedron of feasible solutions of the linear program (2.1).

**Definition 3.2.** A set  $K \subseteq \mathbb{R}^n$  is *convex* if for each  $u, v \in K$  and  $\lambda \in [0, 1]$  the point  $\lambda u + (1 - \lambda)v$  is also contained in  $K$ .

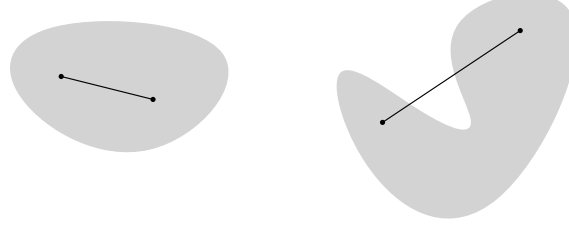


Fig. 3.2: The set on the left is convex, the set on the right is non-convex.

A *halfspace* is a set of solutions of one inequality  $a^T x \leq \beta$  where  $a \in \mathbb{R}^n$  and  $\beta \in \mathbb{R}$ , i.e., a set of the form

$$\{x \in \mathbb{R}^n : a^T x \leq \beta\}.$$

A *hyperplane* is a set of the form

$$\{x \in \mathbb{R}^n : a^T x = \beta\}.$$

It is easy to see that a halfspace is convex. Convexity is also maintained if convex sets are intersected.

**Proposition 3.1.** *Let  $I$  be an index set and  $C_i \subseteq \mathbb{R}^n$  be convex sets for each  $i \in I$ , then  $\cap_{i \in I} C_i$  is a convex set.*

Consequently, the set of feasible solutions of a linear program  $\{x \in \mathbb{R}^n : Ax \leq b\}$  is a convex set. This is our motivation to study properties of convex sets.

### 3.1 Extreme points and vertices

**Definition 3.3.** An inequality  $a^T x \leq \beta$  is *valid* for a set  $K \subseteq \mathbb{R}^n$  if each  $x^* \in K$  satisfies  $a^T x^* \leq \beta$ . If in addition  $(a^T x = \beta) \cap K \neq \emptyset$ , then  $a^T x \leq \beta$  is a *supporting inequality* and  $a^T x = \beta$  is a *supporting hyperplane*.

**Definition 3.4.** Let  $K \subseteq \mathbb{R}^n$  be a convex set. A point  $x^* \in K$  is an *extreme point* or *vertex* of  $K$  if there exists a valid inequality  $a^T x \leq \beta$  of  $K$  such that

$$\{x^*\} = K \cap \{x \in \mathbb{R}^n : a^T x = \beta\}.$$

In other words, if  $x^*$  is the only point of  $K$  that satisfies the valid inequality with equality.

We can now characterize the extreme points of polyhedra. In fact, there are only finitely many of them.

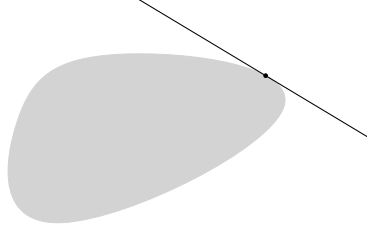


Fig. 3.3: An extreme point of a convex set.

**Theorem 3.1.** Let  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  be a polyhedron. A feasible point  $x^*$  is an extreme point of  $P$  if and only if there is a sub-system  $A'x \leq b'$  of  $Ax \leq b$  such that

- i)  $x^*$  satisfies all inequalities of  $A'x \leq b'$  with equality.
- ii)  $A'$  has  $n$  rows and  $A'$  is non-singular.

*Proof.* Let  $A'x \leq b'$  be such a sub-system and consider the valid inequality  $\mathbf{1}^T A'x \leq \mathbf{1}^T b'$ . Clearly  $x^*$  satisfies this inequality with equality. Any  $y^* \in P$  that satisfies this inequality with equality must satisfy  $A'x = b'$ . Since  $A'$  is non-singular,  $x^*$  is the unique solution of  $A'x = b'$  which means that  $x^*$  is the unique point of  $P$  that satisfies  $\mathbf{1}^T A'x \leq \mathbf{1}^T b'$  with equality.

Assume now that there does not exist a sub-system  $A'x \leq b'$  of  $Ax \leq b$  with properties i) and ii). Denote the sub-system of inequalities that are satisfied by  $x^*$  with equality by  $\tilde{A}x \leq \tilde{b}$ . Then  $\text{rank}(\tilde{A}) < n$  and there exists a  $d \neq 0 \in \mathbb{R}^n$  with  $\tilde{A}d = 0$ . Consequently there exists an  $\varepsilon > 0$  such that  $x^* \pm \varepsilon \cdot d \in P$ .

Clearly, any inequality that is satisfied by  $x^*$  with equality and that is satisfied by  $x^* \pm \varepsilon d$  is satisfied by  $x^* \pm \varepsilon d$  with equality as well. This implies that  $x^*$  is not an extreme point.

The relevance of vertices for linear programming is reflected in the following theorem.

**Theorem 3.2.** If a linear program  $\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\}$  is feasible and bounded and if  $\text{rank}(A) = n$ , then the linear program has an optimal solution that is an extreme point.

*Proof.* We use the following notation. If  $x^*$  is a feasible solution then  $A_{x^*}x \leq b_{x^*}$  is the subsystem of  $Ax \leq b$  that is satisfied by  $x^*$  with equality. The rank of  $x^*$ ,  $\text{rank}(x^*)$  is the rank of  $A_{x^*}$ . The following claim implies the assertion.

If  $x^*$  is feasible and  $\text{rank}(x^*) < n$ , then there exists a  $y^*$  with  $c^T y^* \geq c^T x^*$  and  $\text{rank}(y^*) > \text{rank}(x^*)$ .

To prove this, let  $d \neq 0 \in \mathbb{R}^n$  be a vector with  $A_{x^*}d = 0$ . We can assume  $c^T d \geq 0$  by switching to  $-d$  otherwise.

If  $c^T d > 0$ , then consider the points  $x^* + \lambda d$  with  $\lambda \geq 0$  and let  $\lambda_{max}$  be maximal with the corresponding point feasible. Clearly  $y^* = x^* + \lambda_{max} d$  satisfies the condition of the claim.

Suppose now that  $c^T d = 0$ . Then  $Ad \neq 0$  since  $\text{rank}(A) = n$ . Let  $\lambda_{max}$  be the maximum of the set  $\{\lambda \geq 0 : A(x^* \pm \lambda d) \leq b\}$ . Then  $y^* = x^* + \lambda_{max} d$  or  $y^* = x^* - \lambda_{max} d$  satisfies the condition of the claim.

**Corollary 3.1.** *A linear program  $\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\}$  which is feasible and bounded has an optimal solution.*

*Proof.* The linear program

$$\max\{c^T x^+ - c^T x^- : x^+, x^- \in \mathbb{R}^n, x^+ \geq 0, x^- \geq 0, A(x^+ - x^-) \leq b\}$$

is feasible and bounded and the rank of the constraint matrix is  $2n$ . Thus, by Theorem 3.2 possesses an optimal solution.  $x^+, x^-$ . This corresponds to an optimal solution  $x^+ - x^-$  of the linear program  $\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\}$ .

### 3.2 Linear, affine, conic and convex hulls

We now describe convex sets that are generated by a set  $X \subseteq \mathbb{R}^n$  of vectors of  $\mathbb{R}^n$ . The *linear hull*, *affine hull*, *conic hull* and *convex hull* of  $X$  are defined as follows.

$$\begin{aligned} \text{lin.hull}(X) = \{ & \lambda_1 x_1 + \cdots + \lambda_t x_t \mid t \geq 0, \\ & x_1, \dots, x_t \in X, \lambda_1, \dots, \lambda_t \in \mathbb{R} \} \end{aligned} \quad (3.1)$$

$$\text{affine.hull}(X) = \{ \lambda_1 x_1 + \cdots + \lambda_t x_t \mid t \geq 1, \quad (3.2)$$

$$\begin{aligned} & x_1, \dots, x_t \in X, \sum_{i=1}^t \lambda_i = 1, \lambda_1, \dots, \lambda_t \in \mathbb{R} \} \\ \text{cone}(X) = \{ & \lambda_1 x_1 + \cdots + \lambda_t x_t \mid t \geq 0, \\ & x_1, \dots, x_t \in X, \lambda_1, \dots, \lambda_t \in \mathbb{R}_{\geq 0} \} \end{aligned} \quad (3.3)$$

$$\begin{aligned} \text{conv}(X) = \{ & \lambda_1 x_1 + \cdots + \lambda_t x_t \mid t \geq 1, \\ & x_1, \dots, x_t \in X, \sum_{i=1}^t \lambda_i = 1, \lambda_1, \dots, \lambda_t \in \mathbb{R}_{\geq 0} \} \end{aligned} \quad (3.4)$$

**Proposition 3.2.** *Let  $X \subseteq \mathbb{R}^n$  and  $x_0 \in X$ . One has*

$$\text{affine.hull}(X) = x_0 + \text{lin.hull}(X - x_0),$$

where for  $u \in \mathbb{R}^n$  and  $V \subseteq \mathbb{R}^n$ ,  $u + V$  denotes the set  $u + V = \{u + v \mid v \in V\}$ .

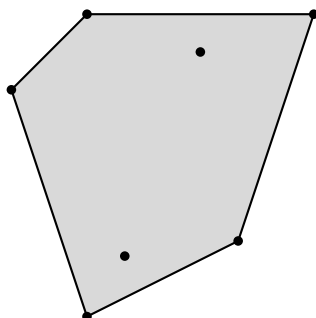


Fig. 3.4: The convex hull of 7 points in  $\mathbb{R}^2$ .

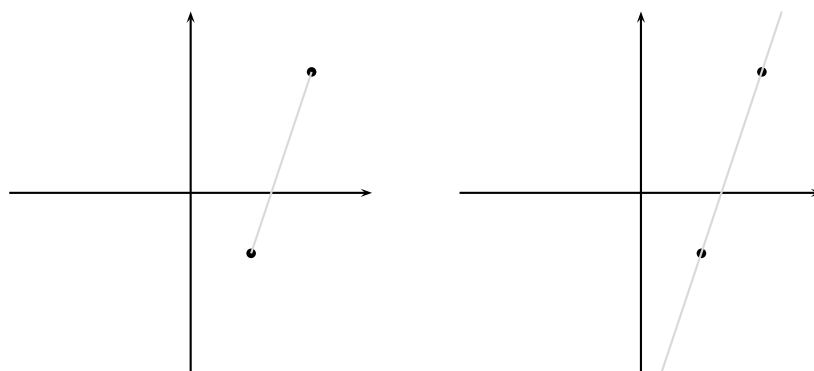


Fig. 3.5: Two points with their convex hull on the left and their affine hull on the right.

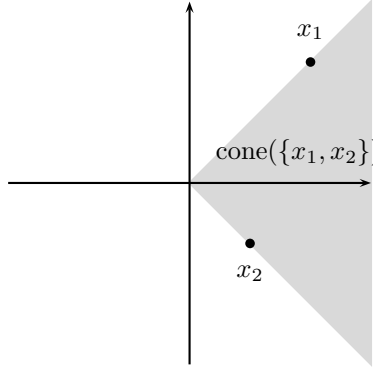


Fig. 3.6: Two points with their conic hull

*Proof.* We first show that each  $x \in \text{affine.hull}(X)$  is also an element of the set  $x_0 + \text{lin.hull}(X - x_0)$  and then we show that each point  $x \in x_0 + \text{lin.hull}(X - x_0)$  is also an element of  $\text{affine.hull}(X)$ .

Let  $x \in \text{affine.hull}(X)$ , i.e., there exists a natural number  $t \geq 1$  and  $\lambda_1, \dots, \lambda_t \in \mathbb{R}$ , with  $x = \lambda_1 x_1 + \dots + \lambda_t x_t$  and  $\sum_{i=1}^t \lambda_i = 1$ . Now

$$\begin{aligned} x &= x_0 - x_0 + \lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_t x_t \\ &= x_0 - \lambda_1 x_0 - \dots - \lambda_t x_0 + \lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_t x_t \\ &= x_0 + \lambda_1 (x_1 - x_0) + \dots + \lambda_t (x_t - x_0), \end{aligned}$$

which shows that  $x \in x_0 + \text{lin.hull}(X - x_0)$ .

Suppose now that  $x \in x_0 + \text{lin.hull}(X - x_0)$ . Then there exist  $\lambda_1, \dots, \lambda_t \in \mathbb{R}$  with  $x = x_0 + \lambda_1 (x_1 - x_0) + \dots + \lambda_t (x_t - x_0)$ . With  $\lambda_0 = 1 - \sum_{i=1}^t \lambda_i$  one has  $\sum_{i=0}^t \lambda_i = 1$  and

$$\begin{aligned} x &= x_0 + \lambda_1 (x_1 - x_0) + \dots + \lambda_t (x_t - x_0) \\ &= \lambda_0 x_0 + \dots + \lambda_t x_t \end{aligned}$$

and thus that  $x \in \text{affine.hull}(X)$ . □

**Theorem 3.3.** *Let  $X \subseteq \mathbb{R}^n$  be a set of points. The convex hull,  $\text{conv}(X)$ , of  $X$  is convex.*

*Proof.* Let  $u$  and  $v$  be points in  $\text{conv}(X)$ . This means that there exists a natural number  $t \geq 1$ , real numbers  $\alpha_i, \beta_i \geq 0$ , and points  $x_i \in X$ ,  $i = 1, \dots, t$  with  $\sum_{i=1}^t \alpha_i = \sum_{i=1}^t \beta_i = 1$  with  $u = \sum_{i=1}^t \alpha_i x_i$  and  $v = \sum_{i=1}^t \beta_i x_i$ . For  $\lambda \in [0, 1]$  one has  $\lambda \alpha_i + (1 - \lambda) \beta_i \geq 0$  for  $i = 1, \dots, t$  and  $\sum_{i=1}^t (\lambda \alpha_i + (1 - \lambda) \beta_i) = 1$ . This shows that

$$\lambda u + (1 - \lambda)v = \sum (\lambda_i \alpha_i + (1 - \lambda_i) \beta_i) x_i \in \text{conv}(X),$$

and therefore that  $\text{conv}(X)$  is convex.  $\square$

**Theorem 3.4.** *Let  $X \subseteq \mathbb{R}^n$  be a set of points. Each convex set  $K$  containing  $X$  also contains  $\text{conv}(X)$ .*

*Proof.* Let  $K$  be a convex set containing  $X$ , and let  $x_1, \dots, x_t \in X$  and  $\lambda_i \in \mathbb{R}$  with  $\lambda_i \geq 0$ ,  $i = 1, \dots, t$  and  $\sum_{i=1}^t \lambda_i = 1$ . We need to show that  $u = \sum_{i=1}^t \lambda_i x_i$  is contained in  $K$ . This is true for  $t \leq 2$  by the definition of convex sets.

We argue by induction. Suppose that  $t \geq 3$ . If one of the  $\lambda_i$  is equal to 0, then one can represent  $u$  as a convex combination of  $t - 1$  points in  $X$  and, by induction,  $u \in K$ . If  $t \geq 3$ , each  $\lambda_i > 0$  and  $\sum_{i=1}^t \lambda_i = 1$ , then one has  $0 < \lambda_i < 1$  for  $i = 1, \dots, t$  and thus we can write

$$u = \lambda_1 x_1 + (1 - \lambda_1) \sum_{i=2}^t \frac{\lambda_i}{1 - \lambda_1} x_i.$$

One has  $\lambda_i / (1 - \lambda_1) > 0$  and

$$\sum_{i=2}^t \frac{\lambda_i}{1 - \lambda_1} = 1,$$

which means that the point  $\sum_{i=2}^t \frac{\lambda_i}{1 - \lambda_1} x_i$  is in  $K$  by induction. Again, by the definition of convex sets, we conclude that  $u$  lies in  $K$ .  $\square$

Theorem 3.4 implies that  $\text{conv}(X)$  is the intersection of all convex sets containing  $X$ , i.e.,

$$\text{conv}(X) = \bigcap_{\substack{K \supseteq X \\ K \text{ convex}}} K.$$

**Definition 3.5.** A set  $C \subseteq \mathbb{R}^n$  is a *cone*, if it is convex and for each  $c \in C$  and each  $\lambda \in \mathbb{R}_{\geq 0}$  one has  $\lambda \cdot c \in C$ .

Similarly to Theorem 3.3 and Theorem 3.4 one proves the following.

**Theorem 3.5.** *For any  $X \subseteq \mathbb{R}^n$ , the set  $\text{cone}(X)$  is a cone.*

**Theorem 3.6.** *Let  $X \subseteq \mathbb{R}^n$  be a set of points. Each cone containing  $X$  also contains  $\text{cone}(X)$ .*

These theorems imply that  $\text{cone}(X)$  is the intersection of all cones containing  $X$ , i.e.,

$$\text{cone}(X) = \bigcap_{\substack{C \supseteq X \\ C \text{ is a cone}}} C.$$

### 3.3 Radon's lemma and Carathéodory's theorem

**Theorem 3.7 (Radon's lemma).** *Let  $A \subseteq \mathbb{R}^n$  be a set of  $n + 2$  points. There exist disjoint subsets  $A_1, A_2 \subseteq A$  with*

$$\text{conv}(A_1) \cap \text{conv}(A_2) \neq \emptyset.$$

*Proof.* Let  $A = \{a_1, \dots, a_{n+2}\}$ . We embed these points into  $\mathbb{R}^{n+1}$  by appending a 1 in the  $n + 1$ -st component, i.e., we construct

$$A' = \left\{ \begin{pmatrix} a_1 \\ 1 \end{pmatrix}, \dots, \begin{pmatrix} a_{n+2} \\ 1 \end{pmatrix} \right\} \subseteq \mathbb{R}^{n+1}.$$

The set  $A'$  consists of  $n + 2$  vectors in  $\mathbb{R}^{n+1}$ . Those vectors are linearly dependent. Let

$$0 = \sum_{i=1}^{n+2} \lambda_i \begin{pmatrix} a_i \\ 1 \end{pmatrix} \quad (3.5)$$

be a nontrivial linear representation of 0, i.e., not all  $\lambda_i$  are 0. Furthermore, let  $P = \{i: \lambda_i \geq 0, i = 1, \dots, n + 2\}$  and  $N = \{i: \lambda_i < 0, i = 1, \dots, n + 2\}$ . We claim that

$$\text{conv}(\{a_i: i \in P\}) \cap \text{conv}(\{a_i: i \in N\}) \neq \emptyset.$$

It follows from (3.5) and the fact that the  $n + 1$ -st component of the vectors is 1 that  $\sum_{i \in P} \lambda_i = -\sum_{i \in N} \lambda_i = s > 0$ . It follows also from (3.5) that

$$\sum_{i \in P} \lambda_i a_i = \sum_{i \in N} -\lambda_i a_i.$$

The point  $u = \sum_{i \in P} (\lambda_i/s) \cdot a_i = \sum_{i \in N} (-\lambda_i/s) a_i$  is contained in  $\text{conv}(\{a_i: i \in P\}) \cap \text{conv}(\{a_i: i \in N\})$ , implying the claim.  $\square$

**Theorem 3.8 (Carathéodory's theorem).** *Let  $X \subseteq \mathbb{R}^n$ , then for each  $x \in \text{cone}(X)$  there exists a set  $\tilde{X} \subseteq X$  of cardinality at most  $n$  such that  $x \in \text{cone}(\tilde{X})$ . The vectors in  $\tilde{X}$  are linearly independent.*

*Proof.* Let  $x \in \text{cone}(X)$ , then there exist  $t \in \mathbb{N}_+$ ,  $x_i \in X$  and  $\lambda_i \geq 0$ ,  $i = 1, \dots, t$ , with  $x = \sum_{i=1}^t \lambda_i x_i$ . Suppose that  $t \in \mathbb{N}_+$  is minimal such that  $x$  can be represented as above. We claim that  $t \leq n$ . If  $t \geq n + 1$ , then the  $x_i$  are linearly dependent. This means that there are  $\mu_i \in \mathbb{R}$ , not all equal to 0 with

$$\sum_{i=1}^t \mu_i x_i = 0. \quad (3.6)$$

By multiplying each  $\mu_i$  in (3.6) with  $-1$  if necessary, we can assume that at least one of the  $\mu_i$  is strictly larger than 0. One has for each  $\varepsilon \in \mathbb{R}$



$$x = \sum_{i=1}^t (\lambda_i - \varepsilon \cdot \mu_i) x_i. \quad (3.7)$$

What is the largest  $\varepsilon^* > 0$  that we can pick for  $\varepsilon$  such that (3.7) is still a conic combination? We need to have

$$\lambda_i - \varepsilon \cdot \mu_i \geq 0, \text{ for each } i \in \{1, \dots, t\}. \quad (3.8)$$

Let  $J$  be the set of indices  $J = \{j : j \in \{1, \dots, t\}, \mu_j > 0\}$ . We observed that we can assume  $J \neq \emptyset$ . We have (3.8) as long as

$$\varepsilon \leq \lambda_j / \mu_j \text{ for each } j \in J. \quad (3.9)$$

This means that  $\varepsilon^* = \min\{\lambda_j / \mu_j : j \in J\}$ . Let  $j^* \in J$  be an index where this minimum is attained. Since  $\lambda_i - \varepsilon^* \cdot \mu_i \geq 0$  for all  $i = 1, \dots, t$  and since  $\lambda_{j^*} - \varepsilon^* \cdot \mu_{j^*} = 0$ , we have  $x \in \text{cone}(\{x_1, \dots, x_t\} \setminus \{x_{j^*}\})$ , which is a contradiction to the minimality of  $t$ .  $\square$

**Corollary 3.2 (Carathéodory's theorem for convex hulls).** *Let  $X \subseteq \mathbb{R}^n$ , then for each  $x \in \text{conv}(X)$  there exists a set  $\tilde{X} \subseteq X$  of cardinality at most  $n + 1$  such that  $x \in \text{conv}(\tilde{X})$ .*

### 3.4 Separation theorem and Farkas' lemma

We recall a basic fact from analysis, see, e.g. [12, Theorem 4.4.1].

**Theorem 3.9.** *Let  $X \subseteq \mathbb{R}^n$  be compact and  $f : X \rightarrow \mathbb{R}$  be continuous. Then  $f$  is bounded and there exist points  $x_1, x_2 \in X$  with  $f(x_1) = \sup\{f(x) : x \in X\}$  and  $f(x_2) = \inf\{f(x) : x \in X\}$ .*

**Theorem 3.10.** *Let  $K \subseteq \mathbb{R}^n$  be a closed convex set and  $x^* \in \mathbb{R}^n \setminus K$ , then there exists an inequality  $a^T x \leq \beta$  such that  $a^T y < \beta$  holds for all  $y \in K$  and  $a^T x^* > \beta$ .*

*Proof.* Since the mapping  $f(x) = \|x^* - x\|$  is continuous and since for any  $k \in K$ ,  $K \cap \{x \in K : \|x^* - x\| \leq \|x^* - k\|\}$  is compact, there exists a point  $k^* \in K$  with minimal distance to  $x^*$ . Consider the midpoint  $m = 1/2(k^* + x^*)$  on the line-segment  $\overline{k^* x^*}$  and the hyperplane  $a^T x = \beta$  with  $\beta = a^T m$  and  $a = (x^* - k^*)$ .

Clearly  $a^T x^* > \beta$  and  $a^T k^* < \beta$ . Suppose now that there exists a point  $k' \in K$  with  $a^T k' \geq \beta$ . Since the line segment spanned by  $k^*$  and  $k'$  is contained in  $K$ , we can assume that  $a^T k' = \beta$  holds. This means that  $k' = k^* + a/2 + x'$ , where  $x'$  is orthogonal to  $a$ . Again, by convexity  $k^* + \lambda(a/2 + x') \in K$  for  $0 \leq \lambda \leq 1$ . But the square of the distance of such a point to  $x^*$  is

$$\begin{aligned}\|x^* - k^* - \lambda(a/2 + x')\|^2 &= \|(1 - \lambda/2)a - \lambda x'\|^2 \\ &= (1 - \lambda/2)^2 \|a\|^2 + \lambda^2 \|x'\|^2\end{aligned}$$

As a function of  $\lambda$ , this is decreasing at  $\lambda = 0$ . Thus there exists a point on the line segment between  $k^*$  and  $k'$  which is closer to  $x^*$  than  $k^*$  and this is a contradiction.  $\square$

**Theorem 3.11 (Farkas' lemma).** *Let  $A \in \mathbb{R}^{m \times n}$  be a matrix and  $b \in \mathbb{R}^m$  be a vector. The system  $Ax = b$ ,  $x \geq 0$  has a solution if and only if for all  $\lambda \in \mathbb{R}^m$  with  $\lambda^T A \geq 0$  one has  $\lambda^T b \geq 0$ .*

*Proof.* Suppose that  $x^* \in \mathbb{R}_{\geq 0}^n$  satisfies  $Ax^* = b$  and let  $\lambda \in \mathbb{R}^m$  with  $\lambda^T A \geq 0$ . Then  $\lambda^T b = \lambda^T Ax^* \geq 0$ , since  $\lambda^T A \geq 0$  and  $x^* \geq 0$ .

Now suppose that  $Ax = b$ ,  $x \geq 0$  does not have a solution. Then, with  $X \subseteq \mathbb{R}^m$  being the set of column vectors of  $A$ ,  $b$  is not in  $\text{cone}(X)$ . The set  $\text{cone}(X)$  is convex and closed, see exercise 8. Theorem 3.10 implies that there is an inequality  $\lambda^T x \geq \beta$  such that  $\lambda^T y > \beta$  for each  $y \in \text{cone}(X)$  and  $\lambda^T b < \beta$ . Since  $0 \in \text{cone}(X)$  it follows that  $0 > \beta$  and thus that  $\lambda^T b < 0$ .

### 3.5 Decomposition theorem for polyhedra

In the following we use the notation  $P(A, b) = \{x \in \mathbb{R}^n : Ax \leq b\}$  for the polyhedron that is defined by  $Ax \leq b$ . We prove the Minkowski-Weyl theorem in this section that shows that polyhedra can be decomposed into the Minkowski sum of a polytope and a cone.

**Definition 3.6.** An inequality  $a^T x \leq \beta$  is called an *implicit equality* of  $Ax \leq b$  if each  $x^* \in P(A, b)$  satisfies  $a^T x^* = \beta$ . We denote the subsystem consisting of implicit equalities of  $Ax \leq b$  by  $A^=x \leq b^=$  and the subsystem consisting of the other inequalities by  $A^<x \leq b^<$ . A constraint is *redundant* if its removal from  $Ax \leq b$  does not change the set of feasible solution of  $Ax \leq b$ .

In the following, a vector  $x$  satisfies  $Ax < b$  if and only if  $a_i^T x < b_i$  for all  $1 \leq i \leq m$ , where  $a_1, \dots, a_m$  are the rows of  $A$ .

**Lemma 3.1.** *Let  $P(A, b)$  be a non-empty polyhedron. Then there exists an  $x \in P(A, b)$  with  $A^<x < b^<$ .*

*Proof.* Suppose that the inequalities in  $A^<x \leq b^<$  are  $a_1^T x \leq \beta_1, \dots, a_k^T x \leq \beta_k$ . For each  $1 \leq i \leq k$  there exists an  $x_i \in P$  with  $a_i^T x_i < \beta_i$ . Thus the point  $x = 1/k(x_1 + \dots + x_k)$  is a point of  $P(A, b)$  satisfying  $A^<x < b^<$ .

**Lemma 3.2.** *Let  $Ax \leq b$  be a system of inequalities. One has*

$$\text{affine.hull}(P(A, b)) = \{x \in \mathbb{R}^n \mid A^=x = b^=\} = \{x \in \mathbb{R}^n \mid A^=x \leq b^=\}.$$

*Proof.* Let  $x_1, \dots, x_t \in P(A, b)$  and suppose that  $a^T x \leq \beta$  is an implicit equality. Then since  $a^T x_i = \beta$  one has  $a^T (\sum_{j=1}^t \lambda_j x_j) = \beta$ . Therefore the inclusions  $\subseteq$  follow.

Suppose now that  $x_0$  satisfies  $A^=x \leq b^=$ . Let  $x_1 \in P(A, b)$  with  $A^{\leq} x_1 < b^{\leq}$ . If  $x_0 = x_1$  then  $x_0 \in P(A, b) \subseteq \text{affine.hull}(P(A, b))$ . Otherwise the line segment between  $x_0$  and  $x_1$  contains more than one point in  $P$  and thus  $x_0 \in \text{affine.hull}(P)$ .

A nonempty set  $C \subseteq \mathbb{R}^n$  is a *cone* if  $\lambda x + \mu y \in C$  for each  $x, y \in C$  and  $\lambda, \mu \in \mathbb{R}_{\geq 0}$ . A cone  $C$  is *polyhedral* if  $C = \{x \in \mathbb{R}^n \mid Ax \leq 0\}$ . A cone *generated by* vectors  $x_1, \dots, x_m \in \mathbb{R}^n$  is a set of the form  $C = \{\sum_{i=1}^m \lambda_i x_i \mid \lambda_i \in \mathbb{R}_{\geq 0}, i = 1, \dots, m\}$ . A point  $x = \sum_{i=1}^m \lambda_i x_i$  with  $\lambda_i \in \mathbb{R}_{\geq 0}, i = 1, \dots, m$  is called a *conic combination* of the  $x_1, \dots, x_m$ . The set of conic combinations of  $X$  is denoted by  $\text{cone}(X)$ .

**Theorem 3.12 (Farkas-Minkowski-Weyl theorem).** *A convex cone is polyhedral if and only if it is finitely generated.*

*Proof.* Suppose that  $a_1, \dots, a_m$  span  $\mathbb{R}^n$  and consider the cone  $C = \{\sum_{i=1}^m \lambda_i a_i \mid \lambda_i \geq 0, i = 1, \dots, m\}$ . Let  $b \notin C$ . Then the system  $A\lambda = b, \lambda \geq 0$  has no solution. By Theorem 3.11 (Farkas' lemma), this implies that there exists a  $y \in \mathbb{R}^n$  such that  $A^T y \leq 0$  and  $b^T y > 0$ .

Suppose that the columns of  $A$  which correspond to inequalities in  $A^T y \leq 0$  that are satisfied by  $y$  with equality have rank  $< n - 1$ . Denote these columns by  $a_{i_1}, \dots, a_{i_k}$ . Then there exists a  $v \neq 0$  which is orthogonal to each of these columns and to  $b$ , i.e.,  $a_{i_j}^T v = 0$  for each  $j = 1, \dots, k$  and  $b^T v = 0$ . There also exists a column  $a^*$  of  $A$  which is not in the set  $\{a_{i_1}, \dots, a_{i_k}\}$  such that  $(a^*)^T v > 0$  since the columns of  $A$  span  $\mathbb{R}^n$ . Therefore there exists an  $\epsilon > 0$  such that

- i)  $A^T(y + \epsilon \cdot v) \leq 0$
- ii) The subspace generated by the columns of  $A$  which correspond to inequalities of  $A^T x \leq 0$  which are satisfied by  $y + \epsilon \cdot v$  with equality strictly contains  $\langle a_{i_1}, \dots, a_{i_k} \rangle$ .

Notice that we have  $b^T y = b^T(y + \epsilon \cdot v) > 0$ .

Continuing this way, we obtain a solution of the form  $y + u$  of  $A^T x \leq 0$  such that one has  $n - 1$  linearly independent columns of  $A$  whose corresponding inequality in  $A^T x \leq 0$  are satisfied with equality. Thus we see that each  $b$  which does not belong to  $C$  can be separated from  $C$  with an inequality of the form  $c^T x \leq 0$  which is uniquely defined by  $n - 1$  linearly independent vectors from the set  $a_1, \dots, a_m$ . This shows that  $C$  is polyhedral.

Suppose now that  $a_1, \dots, a_m$  do not span  $\mathbb{R}^n$ . Then there exist linearly independent vectors  $d_1, \dots, d_k$  such that each  $d_i$  is orthogonal to each of the  $a_1, \dots, a_m$  and  $a_1, \dots, a_m, d_1, \dots, d_k$  spans  $\mathbb{R}^n$ . The cone generated by  $a_1, \dots, a_m, d_1, \dots, d_k$  is polyhedral and thus of the form  $Ax \leq 0$  with some

matrix  $A \in \mathbb{R}^{m \times n}$ . Suppose that  $\langle a_1, \dots, a_m \rangle = \{x \in \mathbb{R}^n \mid Ux = 0\}$ . Now  $C = \{x \in \mathbb{R}^n \mid Ax \leq 0, Ux = 0\}$  and  $C$  is polyhedral.

Now suppose that  $C = \{x \in \mathbb{R}^n \mid a_1^T x \leq 0, \dots, a_m^T x \leq 0\}$ . The cone

$$C' := \text{cone}(a_1, \dots, a_m) = \left\{ \sum_{i=1}^m \lambda_i a_i \mid \lambda_i \geq 0, i = 1, \dots, m \right\}$$

is polyhedral and thus of the form  $C' = \{x \in \mathbb{R}^n \mid b_1^T x \leq 0, \dots, b_k^T x \leq 0\}$ . Clearly,  $\text{cone}(b_1, \dots, b_k) \subseteq C$  since  $b_i^T a_j \leq 0$ . Suppose now that  $y \in C \setminus \text{cone}(b_1, \dots, b_k)$ . Then, since  $\text{cone}(b_1, \dots, b_k)$  is polyhedral, there exists a  $w \in \mathbb{R}^n$  with  $w^T y > 0$  and  $w^T b_i \leq 0$  for each  $i = 1, \dots, k$ . From the latter we conclude that  $w \in C'$ . From  $y \in C$  and  $w \in C'$  we conclude  $w^T y \leq 0$ , which is a contradiction.

A set of vectors  $Q = \text{conv}(X)$ , where  $X \subseteq \mathbb{R}^n$  is finite is called a *polytope*.

**Theorem 3.13 (Decomposition theorem for polyhedra).** *A set  $P \subseteq \mathbb{R}^n$  is a polyhedron if and only if  $P = Q + C$  for some polytope  $Q$  and a polyhedral cone  $C$ .*

*Proof.* Suppose  $P = \{x \in \mathbb{R}^n \mid Ax \leq b\}$  is a polyhedron. Consider the polyhedral cone

$$\left\{ \begin{pmatrix} x \\ \lambda \end{pmatrix} \mid x \in \mathbb{R}^n, \lambda \in \mathbb{R}_{\geq 0}; Ax - \lambda b \leq 0 \right\} \quad (3.10)$$

is generated by finitely many vectors  $\begin{pmatrix} x_i \\ \lambda_i \end{pmatrix}$ ,  $i = 1, \dots, m$ . By scaling with a positive number we may assume that each  $\lambda_i \in \{0, 1\}$ . Let  $Q$  be the convex hull of the  $x_i$  with  $\lambda_i = 1$  and let  $C$  be the cone generated by the  $x_i$  with  $\lambda_i = 0$ . A point  $x \in \mathbb{R}^n$  is in  $P$  if and only if  $\begin{pmatrix} x \\ 1 \end{pmatrix}$  belongs to (3.10) and thus if and only if

$$\begin{pmatrix} x \\ 1 \end{pmatrix} \in \text{cone} \left\{ \begin{pmatrix} x_1 \\ \lambda_1 \end{pmatrix}, \dots, \begin{pmatrix} x_m \\ \lambda_m \end{pmatrix} \right\}.$$

Therefore  $P = Q + C$ .

Suppose now that  $P = Q + C$  for some polytope  $Q$  and a polyhedral cone  $C$  with  $Q = \text{conv}(x_1, \dots, x_m)$  and  $C = \text{cone}(y_1, \dots, y_t)$ . A vector  $x_0$  is in  $P$  if and only if

$$\begin{pmatrix} x_0 \\ 1 \end{pmatrix} \in \text{cone} \left\{ \begin{pmatrix} x_1 \\ 1 \end{pmatrix}, \dots, \begin{pmatrix} x_m \\ 1 \end{pmatrix}, \begin{pmatrix} y_1 \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} y_t \\ 0 \end{pmatrix} \right\} \quad (3.11)$$

By Theorem 3.12 (3.11) is equal to

$$\left\{ \begin{pmatrix} x \\ \lambda \end{pmatrix} \mid Ax - \lambda b \leq 0 \right\} \quad (3.12)$$

for some matrix  $A$  and vector  $b$ . Thus  $x_0 \in P$  if and only if  $Ax_0 \leq b$  and thus  $P$  is a polyhedron.

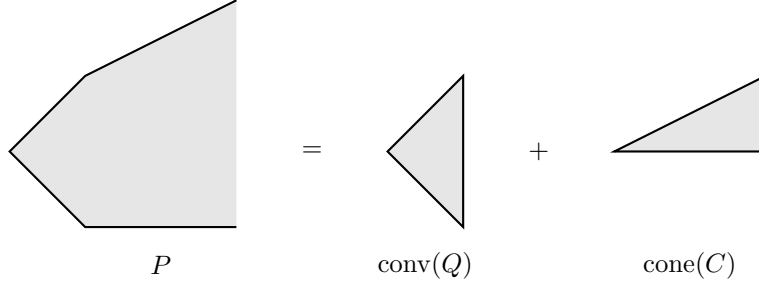


Fig. 3.7: A polyhedron and its decomposition into  $Q$  and  $C$

Let  $P = \{x \in \mathbb{R}^n \mid Ax \leq b\}$ . The *characteristic cone* is  $\text{char.cone}(P) = \{y \mid y + x \in P \text{ for all } x \in P\} = \{y \mid Ay \leq 0\}$ . One has

- i)  $y \in \text{char.cone}(P)$  if and only if there exists an  $x \in P$  such that  $x + \lambda y \in P$  for all  $\lambda \geq 0$
- ii)  $P + \text{char.cone}(P) = P$
- iii)  $P$  is bounded if and only if  $\text{char.cone}(P) = \{0\}$ .
- iv) If the decomposition of  $P$  is  $P = Q + C$ , then  $C = \text{char.cone}(P)$ .

The *lineality space* of  $P$  is defined as  $\text{char.cone}(P) \cap -\text{char.cone}(P)$ . A polyhedron is *pointed*, if its lineality space is  $\{0\}$ .

**Definition 3.7.** A set  $F \subseteq \mathbb{R}^n$  is called a *face* of  $P$  if there exists a valid inequality  $c^T x \leq \delta$  for  $P$  with  $F = P \cap (c^T x = \delta)$ .

**Lemma 3.3.** A set  $\emptyset \neq F \subseteq \mathbb{R}^n$  is a face of  $P$  if and only if  $F = \{x \in P \mid A'x = b'\}$  for a subset  $A'x \leq b'$  of  $Ax \leq b$ .

*Proof.* Suppose that  $F = \{x \in P \mid A'x = b'\}$ . Consider the vector  $c = 1^T A'$  and  $\delta = 1^T b'$ . The inequality  $c^T x \leq \delta$  is valid for  $P$ . It is satisfied with equality by each  $x \in F$ . If  $x' \in P \setminus F$ , then there exists an inequality  $a^T x \leq \beta$  of  $A'x \leq b'$  such that  $a^T x' < \beta$  and consequently  $c^T x' < \delta$ .

On the other hand, if  $c^T x \leq \delta$  defines the face  $F$ , then by the linear programming duality (see chapter 5)

$$\max\{c^T x \mid Ax \leq b\} = \min\{b^T \lambda \mid A^T \lambda = c, \lambda \geq 0\}$$

there exists a  $\lambda \in \mathbb{R}_{\geq 0}^m$  such that  $c = \lambda^T A$  and  $\delta = \lambda^T b$ . Let  $A'x \leq b'$  be the subsystem of  $Ax \leq b$  which corresponds to strictly positive entries in  $Ax \leq b$ . One has  $F = \{x \in P \mid A'x = b'\}$ .

A *facet* of  $P$  is an inclusion-wise maximal face  $F$  of  $P$  with  $F \neq P$ . An inequality  $a^T x \leq \beta$  of  $Ax \leq b$  is called *redundant* if  $P(A, b) = P(A', b')$ , where  $A'x \leq b'$  is the system stemming from  $Ax \leq b$  by deleting  $a^T x \leq \beta$ . A system  $Ax \leq b$  is *irredundant* if  $Ax \leq b$  does not contain a redundant inequality.

**Lemma 3.4.** *Let  $Ax \leq b$  be an irredundant system. Then a set  $F \subseteq P$  is a facet if and only if it is of the form  $F = \{x \in P \mid a^T x = \beta\}$  for an inequality  $a^T x \leq \beta$  of  $Ax \leq b$ .*

*Proof.* Let  $F$  be a facet of  $P$ . Then  $F = \{x \in P \mid c^T x \leq \delta\}$  for a valid inequality  $c^T x \leq \delta$  of  $P$ . There exists a  $\lambda \in \mathbb{R}_{\geq 0}^m$  with  $c = \lambda^T A$  and  $\delta = \lambda^T b$ . There exists an inequality  $a^T x \leq \beta$  of  $Ax \leq b$  whose corresponding entry in  $\lambda$  is strictly positive. Clearly  $F \subseteq \{x \in P \mid a^T x = \beta\} \subset P$ . Since  $F$  is an inclusion-wise maximal face one has  $F = \{x \in P \mid a^T x = \beta\}$ .

Let  $F$  be of the form  $F = \{x \in P \mid a^T x = \beta\}$  for an inequality  $a^T x \leq \beta$  of  $Ax \leq b$ . Clearly  $F \neq \emptyset$  since the system  $Ax \leq b$  is irredundant. If  $F$  is not a facet, then  $F \subseteq F' = \{x \in P \mid a'^T x = \beta'\}$  with another inequality  $a'^T x \leq \beta'$  of  $Ax \leq b$ . Let  $x^* \in \mathbb{R}^n$  be a point with  $a'^T x^* > \beta'$  and which satisfies all other inequalities of  $Ax \leq b$ . Such an  $x^*$  exists, since  $Ax \leq b$  is irredundant. Let  $\tilde{x} \in P$  with  $A\tilde{x} < b$ . There exists a point  $\bar{x}$  on the line-segment  $\tilde{x}x^*$  with  $a^T \bar{x} = \beta$ . This point is then also in  $F'$  and thus  $a'^T \bar{x} = \beta'$  follows. This shows that  $a'^T x^* > \beta'$  and thus  $a'^T x \leq \beta'$  can be removed from the system. This is a contradiction to  $Ax \leq b$  being irredundant.

**Lemma 3.5.** *A face  $F$  of  $P(A, b)$  is inclusion-wise minimal if and only if it is of the form  $F = \{x \in \mathbb{R}^n \mid A'x = b'\}$  for some subsystem  $A'x \leq b'$  of  $Ax \leq b$ .*

*Proof.* Let  $F$  be a minimal face of  $P$  and let  $A'x \leq b'$  a the subsystem of inequalities of  $Ax \leq b$  with  $F = \{x \in P \mid A'x = b'\}$ . Suppose that  $F \subset \{x \in \mathbb{R}^n \mid A'x = b'\}$  and let  $x_1 \in \mathbb{R}^n \setminus P$  satisfy  $A'x_1 = b'$  and  $x_2 \in F$ . There exists “a first” inequality  $a^T x \leq \beta$  of  $Ax \leq b$  which is “hit” by the line-segment  $\overline{x_2 x_1}$ . Let  $x^* = \overline{x_2 x_1} \cap (a^T x = \beta)$ . Then  $x^* \in F$  and thus  $F \cap (a^T x = \beta) \neq \emptyset$ . But  $F \supset F \cap (a^T x = \beta)$  since  $a^T x \leq \beta$  is not an inequality of  $A'x \leq b'$ . This is a contradiction to the minimality of  $F$ .

Suppose that  $F$  is a face with  $F = \{x \in \mathbb{R}^n \mid A'x = b'\} = \{x \in P \mid A'x = b'\}$  for a subsystem  $A'x \leq b'$  of  $Ax \leq b$ . Suppose that there exists a face  $\tilde{F}$  of  $P$  with  $\emptyset \subset \tilde{F} \subset F$ . By Lemma 3.3  $\tilde{F} = \{x \in P \mid A'x = b', A^*x = b^*\}$ , where  $A^*x \leq b^*$  is a sub-system of  $Ax \leq b$  which contains an inequality  $a^T x \leq \beta$  such that there exists an  $x_1, x_2 \in \tilde{F}$  with  $a^T x_1 < \beta$  and  $a^T x_2 = \beta$ . The line  $\ell(x_1, x_2) = \{x_1 + \lambda(x_2 - x_1) \mid \lambda \in \mathbb{R}\}$  is contained in  $\tilde{F}$  but is not contained in  $a^T x \leq \beta$ . This shows that  $\tilde{F}$  is not contained in  $P$  which is a contradiction.

Exercise 21 asks for a proof of the following corollary.

**Corollary 3.3.** *Let  $F_1$  and  $F_2$  be two inclusion-wise minimal faces of  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$ , then  $\dim(F_1) = \dim(F_2)$ .*

We say that a polyhedron contains a line  $\ell(x_1, x_2)$  with  $x_1 \neq x_2 \in P$  if  $\ell(x_1, x_2) = \{x_1 + \lambda(x_2 - x_1) \mid \lambda \in \mathbb{R}\} \subseteq P$ . A *vertex* of  $P$  is a 0-dimensional face of  $P$ . An *edge* of  $P$  is a 1-dimensional face of  $P$ .

*Example 3.2.* Consider a linear program  $\min\{c^T x : Ax = b, x \geq 0\}$ . A basic feasible solution defined by the basis  $B \subseteq \{1, \dots, n\}$  is a vertex of the polyhedron  $P = \{x \in \mathbb{R}^n : Ax = b, x \geq 0\}$ . This can be seen as follows. The inequality  $a^T x \geq 0$  is valid for  $P$ , where  $a_B = \mathbf{0}$  and  $a_{\bar{B}} = \mathbf{1}$ . The inequality is satisfied with equality by a point  $x^* \in P$  if and only if  $x_{\bar{B}}^* = \mathbf{0}$ . Since the columns of  $A_B$  are linearly independent, as  $B$  is a basis, the unique point which satisfies  $a^T x \geq 0$  with equality is the basic feasible solution.

In exercise 23 you are asked to show that the simplex method can be geometrically interpreted as a walk on the graph  $G = (V, E)$ , where  $V$  is the set of basic feasible solutions and  $uv \in E$  if and only if  $\text{conv}\{u, v\}$  is a 1-dimensional face of the polyhedron defined by the linear program.

## Exercises

- 1) Consider the unit ball  $B_n = \{x \in \mathbb{R}^n : \|x\|_2 \leq 1\}$ . Show that the set of extreme points of  $B$  is the sphere  $S^{(n-1)} = \{x \in \mathbb{R}^n : \|x\|_2 = 1\}$ .
- 2) A *line* is a set  $L = \{x \cdot d + t : x \in \mathbb{R}\} \subseteq \mathbb{R}^n$  where  $d, t \in \mathbb{R}^n$   $d \neq 0$ . Show the following.  
A non-empty polyhedron  $P = \{x \in \mathbb{R}^n : Ax \leq b\} \subseteq \mathbb{R}^n$  contains a line if and only if  $\text{rank}(A) < n$ .
- 3) Two different vertices  $v_1 \neq v_2$  of a polyhedron  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  are called *adjacent*, if there exists a subsystem  $A'x \leq b'$  of  $Ax \leq b$  with
  - i)  $A'v_1 = b'$  and  $A'v_2 = b'$  and
  - ii)  $\text{rank}(A') = (n - 1)$ .

Show that there exists a valid inequality  $c^T x \leq \delta$  of  $P$  with

$$(P \cap \{x \in \mathbb{R}^n : c^T x = \delta\}) = \text{conv}\{v_1, v_2\}.$$

- 4) Let  $\{C_i\}_{i \in I}$  be a family of convex subsets of  $\mathbb{R}^n$ . Show that the intersection  $\bigcap_{i \in I} C_i$  is convex.
- 5) Show that the set of feasible solutions of a linear program is convex.
- 6) Prove Carathéodory's Theorem for convex hulls, Corollary 3.2.
- 7) Let  $A \in \mathbb{R}^{n \times n}$  be a non-singular matrix and let  $a_1, \dots, a_n \in \mathbb{R}^n$  be the columns of  $A$ . Show that  $\text{cone}(\{a_1, \dots, a_n\})$  is the polyhedron  $P = \{y \in$

$\mathbb{R}^n: A^{-1}y \geq 0\}$ . Show that  $\text{cone}(\{a_1, \dots, a_k\})$  for  $k \leq n$  is the set  $P_k = \{y \in \mathbb{R}^n: a_i^{-1}x \geq 0, i = 1, \dots, k, a_i^{-1}x = 0, i = k+1, \dots, n\}$ , where  $a_i^{-1}$  denotes the  $i$ -th row of  $A^{-1}$ .

- 8) Prove that for a finite set  $X \subseteq \mathbb{R}^n$  the conic hull  $\text{cone}(X)$  is closed and convex.

*Hint: Use Carathéodory's theorem and exercise 7.*

- 9) Find a countably infinite set  $X \subset \mathbb{R}^2$  such that  $\text{cone}(X)$  is not closed. Are there any cones that are open?
- 10) Prove Theorem 3.5.
- 11) Prove Theorem 3.6.
- 12) Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}^d$  be a linear map.
- Show that  $f(K) = \{f(x): x \in K\}$  is convex if  $K$  is convex. Is the reverse also true?
  - For  $X \subseteq \mathbb{R}^n$  arbitrary, prove that  $\text{conv}(f(X)) = f(\text{conv}(X))$ .
- 13) Using Theorem 3.11, prove the following variant of Farkas' lemma: Let  $A \in \mathbb{R}^{m \times n}$  be a matrix and  $b \in \mathbb{R}^m$  be a vector. The system  $Ax \leq b$ ,  $x \in \mathbb{R}^n$  has a solution if and only if for all  $\lambda \in \mathbb{R}_{\geq 0}^m$  with  $\lambda^T A = 0$  one has  $\lambda^T b \geq 0$ .
- 14) Provide an example of a convex and closed set  $K \subseteq \mathbb{R}^2$  and a linear objective function  $c^T x$  such that  $\inf\{c^T x: x \in K\} > -\infty$  but there does not exist an  $x^* \in K$  with  $c^T x^* \leq c^T x$  for all  $x \in K$ .
- 15) Consider the vectors

$$x_1 = \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix}, x_2 = \begin{pmatrix} 1 \\ 2 \\ 5 \end{pmatrix}, x_3 = \begin{pmatrix} 2 \\ 0 \\ 1 \end{pmatrix}, x_4 = \begin{pmatrix} 2 \\ 4 \\ 3 \end{pmatrix}, x_5 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}.$$

Let  $A = \{x_1, \dots, x_5\}$ . Find two disjoint subsets  $A_1, A_2 \subseteq A$  such that

$$\text{conv}(A_1) \cap \text{conv}(A_2) \neq \emptyset.$$

*Hint: Recall the proof of Radon's lemma*

- 16) Consider the vectors

$$x_1 = \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix}, x_2 = \begin{pmatrix} 1 \\ 2 \\ 5 \end{pmatrix}, x_3 = \begin{pmatrix} 2 \\ 0 \\ 1 \end{pmatrix}, x_4 = \begin{pmatrix} 2 \\ 4 \\ 3 \end{pmatrix}, x_5 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}.$$

The vector

$$v = x_1 + 3x_2 + 2x_3 + x_4 + 3x_5 = \begin{pmatrix} 15 \\ 14 \\ 25 \end{pmatrix}$$

is a conic combination of the  $x_i$ .

Write  $v$  as a conic combination using only three vectors of the  $x_i$ .

*Hint: Recall the proof of Carathéodory's theorem*



- 17) Prove that each nonempty polyhedron  $P \subseteq \mathbb{R}^n$  can be represented as  $P = L + Q$ , where  $L \subseteq \mathbb{R}^n$  is a linear space and  $Q \subseteq \mathbb{R}^n$  is a pointed polyhedron.
- 18) Let  $P \subseteq \mathbb{R}^n$  be a polytope and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  a linear map.
  - i) Show that  $f(P)$  is a polytope.
  - ii) Let  $y \in \mathbb{R}^m$  be a vertex of  $f(P)$ . Show that there is a vertex  $x \in \mathbb{R}^n$  of  $P$  such that  $f(x) = y$ .
- 19) Let  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$  and consider the polyhedron  $P = P(A, b)$ . Show that  $\dim(P) = n - \text{rank}(A^\circ)$ .
- 20)
  - i) Show that the dimension of each minimal face of a polyhedron  $P$  is equal to  $n - \text{rank}(A)$ .
  - ii) Show that a polyhedron has a vertex if and only if the polyhedron does not contain a line.
- 21) Show that the affine dimension of the minimal faces of a polyhedron  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  is invariant.
- 22) In this exercise you can assume that a linear program  $\max\{c^T x \mid Ax \leq b\}$  can be solved in polynomial time. Suppose that  $P(A, b)$  has vertices and that the linear program is bounded. Show how to compute an optimal *vertex* solution of the linear program in polynomial time.
- 23) Let  $P = \{x \in \mathbb{R}^n : Ax = b, x \geq 0\}$  be a polyhedron, where  $A \in \mathbb{R}^{m \times n}$  has full row-rank. Let  $B_1, B_2$  be two bases such that  $|B_1 \cap B_2| = m - 1$  and suppose that the associated basic solutions  $x_1^*$  and  $x_2^*$  are feasible. Show that, if  $x_1 \neq x_2$ , then  $\text{conv}\{x_1^*, x_2^*\}$  is a 1-dimensional face of  $P$ .



## Chapter 4

# The simplex method

In this chapter we describe the simplex method. The task is to solve a linear program

$$\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\}. \quad (4.1)$$

We make the following assumption.

The matrix  $A \in \mathbb{R}^{m \times n}$  is of full column-rank. In other words, the columns of  $A$  are linearly independent.

This assumption is not a restriction, since we can solve the following equivalent linear program instead, where each  $x_i$  is represented as the difference of two positive values  $x_i = x_i^+ - x_i^-$ .

$$\max\{c^T x^+ - c^T x^- : x^+, x^- \in \mathbb{R}^n, Ax^+ - Ax^- \leq b, x^+ \geq 0, x^- \geq 0\}. \quad (4.2)$$

The constraint matrix of the linear program (4.2) in inequality standard form is

$$\begin{pmatrix} A & -A \\ -I_n & \mathbf{0} \\ \mathbf{0} & -I_n \end{pmatrix},$$

where  $I_n$  is the  $n \times n$  identity matrix and  $\mathbf{0}$  is the  $n \times n$  all-zero matrix. Clearly this matrix has linearly independent columns.

### 4.1 Adjacent vertices

Let  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  be the polyhedron of feasible solutions of (4.1).

**Definition 4.1.** Two extreme points  $x_1 \neq x_2$  of  $P$  are *adjacent*, if there exists a valid inequality  $d^T x \leq \delta$  of  $P$  such that

$$P \cap \{x \in \mathbb{R}^n : d^T x = \delta\} = \{\lambda x_1 + (1 - \lambda)x_2 : \lambda \in [0, 1]\}.$$

In other words,  $x_1 \neq x_2$  are adjacent if there exists a valid inequality for  $P$  such that the points of  $P$  that satisfy this inequality with equality are exactly the line-segment spanned by  $x_1$  and  $x_2$ .

Similar to the characterization of extreme points in Theorem 3.1, we can state and prove the following theorem.

**Theorem 4.1.** *Two distinct vertices  $x_1$  and  $x_2$  of  $P$  are adjacent if and only if there exists a sub-system  $A'x \leq b'$  of  $Ax \leq b$  such that*

- i)  $A' \in \mathbb{R}^{(n-1) \times n}$  and the rows of  $A'$  are linearly independent.
- ii)  $x_1$  and  $x_2$  satisfy  $A'x \leq b'$  with equality.

*Proof.* Suppose that  $x_1 \neq x_2$  are adjacent and suppose that  $d^T x \leq \delta$  is a valid inequality that asserts this fact. Consider the sub-system  $\tilde{A}x \leq \tilde{b}$  of inequalities of  $Ax \leq b$  that are satisfied by  $1/2(x_1 + x_2)$  with equality. These are the inequalities that are satisfied by all points on the line-segment with equality. Since  $\tilde{A}(x_1 - x_2) = 0$ , one has  $\text{rank}(\tilde{A}) \leq n - 1$ . If  $\text{rank}(\tilde{A}) < n - 1$ , then there exists a  $v \in \mathbb{R}^n$  that is linearly independent from  $(x_1 - x_2)$  that satisfies  $\tilde{A}v = 0$ . Consequently there exist a  $\varepsilon > 0$  such that

$$\left\{ \frac{1}{2}(x_1 + x_2) + \mu_1(x_1 - x_2) + \mu_2 v : -\varepsilon \leq \mu_1, \mu_2 \leq \varepsilon \right\} \subseteq P. \quad (4.3)$$

All points of the set (4.3) satisfy  $d^T x \leq \delta$  with equality and they are not a subset of the line-segment spanned by  $x_1$  and  $x_2$ . From this we conclude that  $\text{rank}(\tilde{A}) = n - 1$  which implies that there exists a sub-system  $A'x \leq b'$  satisfying i) and ii).

Suppose on the other hand that there exists a sub-system  $A'x \leq b'$  that satisfies i) and ii). The line spanned by  $x_1$  and  $x_2$  is the set of points of  $\mathbb{R}^n$  that satisfies  $A'x = b'$  and the intersection of this line with  $P$  is, since  $x_1$  and  $x_2$  are vertices, the line-segment spanned by these two points.

The inequality  $\mathbf{1}^T A'x \leq \mathbf{1}^T b'$  is valid for  $P$  and is satisfied by the line-segment spanned by  $x_1$  and  $x_2$  with equality. Let  $y^* \in P$  be a point that does not lie on the line segment. Then one of the inequalities of  $A'x \leq b'$  is satisfied by  $y^*$  with strict inequality and thus  $y^*$  does not satisfy  $d^T x \leq \delta$  with equality.

## 4.2 Bases, feasible bases and vertices

We will frequently use the following notation. Let  $B \subseteq \{1, \dots, m\}$  then  $A_B \in \mathbb{R}^{|B| \times n}$  is the matrix consisting of the rows of  $A$  that are indexed by  $B$  and  $b_B \in \mathbb{R}^{|B|}$  is the vector whose components are the ones of  $b$  indexed by  $B$ .

*Example 4.1.* For  $A = \begin{pmatrix} 3 & 2 \\ 7 & 1 \\ 8 & 4 \end{pmatrix}$ ,  $b = \begin{pmatrix} 3 \\ 2 \\ 6 \end{pmatrix}$  and  $B = \{2, 3\}$ , one has

$$A_B = \begin{pmatrix} 7 & 1 \\ 8 & 4 \end{pmatrix} \text{ and } b_B = \begin{pmatrix} 2 \\ 6 \end{pmatrix}.$$

A *python implementation* based on the `sympy` package of the above is as follows. Notice that the index set represented by a list `B` here is starting at 1, since the first row of a matrix in `sympy` is indexed by 0.

```
1 from sympy import *
2 A = Matrix([[3,2],[7,1],[8,4]])
3 b = Matrix([3,2,6])
4
5 B = [1,2]
6
7 pprint(A[B,:])
8 pprint(b[B,:])
```

**Definition 4.2.** An index set  $B \subseteq \{1, \dots, m\}$  is a *basis* if  $|B| = n$  and  $A_B$  is non-singular. If in addition  $x^* = A_B^{-1}b_B$  is feasible, then  $B$  is called a *feasible basis*.

Theorem 3.1 implies that every vertex  $x^*$  of  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  is *represented* by a basis  $B$ , i.e.  $x^*$  is the unique solution of  $A_B x^* = b_B$ . This representation however must not be unique, see Figure 4.1. We say that a linear program is *degenerate*, if there exists a basic solution  $x^* \in \mathbb{R}^n$  that satisfies  $n+1$  inequalities with equality. Otherwise the linear program is called *non-degenerate*. If the linear program is non-degenerate, then each vertex is represented by *exactly one* basis.

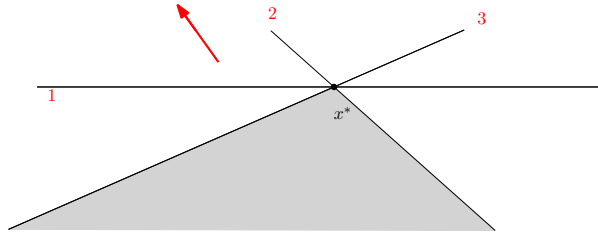


Fig. 4.1: The vertex  $x^*$  is represented by each choice of two of the three tight constraints. The linear program is degenerate. The red vector is the objective function vector and the red labels are the indices of the constraints.

**Definition 4.3.** A basis  $B$  is called *optimal* if it is feasible and the unique  $\lambda \in \mathbb{R}^m$  with

$$\lambda^T A = c^T \text{ and } \lambda_i = 0, i \notin B \quad (4.4)$$

satisfies  $\lambda \geq 0$ .

The basis  $\{1, 2\}$  in Figure 4.1 is not optimal whereas the bases  $\{2, 3\}$  and  $\{1, 3\}$  are optimal bases.

**Theorem 4.2.** *If  $B$  is an optimal basis, then  $x^* = A_B^{-1}b_B$  is an optimal solution of the linear program (4.1).*

*Proof.* The inequality  $\lambda^T Ax \leq \lambda^T b$  is valid for  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$ . But  $\lambda^T A = c^T$  and  $\lambda^T b = \lambda^T A x^* = c^T x^*$ . Consequently  $x^*$  is an optimal solution of the linear program (4.1).

### 4.3 Moving to an improving vertex

Suppose now that  $B$  is a feasible but not optimal basis. Then the unique  $\lambda$  satisfying (4.4) has a negative component  $\lambda_i < 0$  for some  $i \in B$ . The idea is now to move from  $x_B^* = A_B^{-1}b_B$  by remaining tight at all constraints indexed by  $B$  except for  $i$ .

There is only one way this can be achieved. Namely by moving in the unique direction  $d$  with

$$a_j^T d = \begin{cases} 0 & \text{for } j \in B \setminus \{i\} \\ -1 & \text{if } j = i. \end{cases}$$

When we do this, we follow the ray  $x_B^* + \varepsilon \cdot d$  with  $\varepsilon \geq 0$ . What happens to the objective function, as  $\varepsilon$  grows? Since  $c^T d = \lambda^T A d = -\lambda_i > 0$ , the objective function strictly grows with growing  $\varepsilon$ . There are now two cases.

At some point, we hit the boundary of a constraint and further increase of  $\varepsilon$  results in an infeasible point. Let  $K \subseteq \{1, \dots, m\}$  be the set of indices

$$K = \{k : 1 \leq k \leq m, a_k^T d > 0\}. \quad (4.5)$$

Those are the indices of constraints that, at some point, will be violated. We can increase  $\varepsilon$  until

$$\varepsilon^* = \min_{k \in K} \{(b_k - a_k^T x^*) / a_k^T d\}. \quad (4.6)$$

Now pick any  $k \in K$  for which this minimum is achieved and set  $B' = B \setminus \{i\} \cup \{k\}$ . This is a feasible basis,  $d$  is orthogonal to all rows indexed by  $B \setminus \{i\}$  but not to  $a_k$ .

In the case where there is no constraint that puts an upper bound on  $\varepsilon$ , then the linear program is *unbounded*. We have described one iteration of the simplex algorithm. We iterate this procedure until an optimal solution is found.

*Example 4.2.* Consider the linear program  $\max\{c^T x : A \leq b, x \in \mathbb{R}^3\}$  with

$$A = \begin{pmatrix} 1 & 2 & 2 \\ 2 & 1 & 2 \\ 2 & 2 & 1 \\ -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{pmatrix}, b = \begin{pmatrix} 10 \\ 14 \\ 11 \\ 0 \\ 0 \\ 0 \end{pmatrix} \text{ and } c = \begin{pmatrix} 6 \\ 14 \\ 13 \end{pmatrix}$$

with starting basis

$$B = \{1, 2, 3\}.$$

One has

$$A_B = \begin{pmatrix} 1 & 2 & 2 \\ 2 & 1 & 2 \\ 2 & 2 & 1 \end{pmatrix} b_B = \begin{pmatrix} 10 \\ 14 \\ 11 \end{pmatrix} \text{ and } x_B^* = \begin{pmatrix} 4 \\ 0 \\ 3 \end{pmatrix}.$$

For  $\lambda_B$  one obtains

$$\lambda_B = \begin{pmatrix} \frac{36}{5} \\ -\frac{4}{5} \\ \frac{1}{5} \end{pmatrix}.$$

From this it follows that 2 leaves the basis. The direction  $d$  is the solution of the system

$$A_B d = \begin{pmatrix} 0 \\ -1 \\ 0 \end{pmatrix}.$$

One has

$$d = \begin{pmatrix} -\frac{2}{5} \\ \frac{3}{5} \\ -\frac{2}{5} \end{pmatrix}.$$

The vector  $Ad$  is

$$A \cdot d = \begin{pmatrix} 0 \\ -1 \\ 0 \\ \frac{2}{5} \\ -\frac{3}{5} \\ \frac{2}{5} \end{pmatrix}.$$

and one has

$$b - Ax^* = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 4 \\ 0 \\ 3 \end{pmatrix}.$$

This means that 6 enters the basis.

```

1 from sympy import *
2
3 A = Matrix([[1, 2, 2],
4             [2, 1, 2],
5             [2, 2, 1],
6             [-1, 0, 0],
7             [0, -1, 0],
8             [0, 0, -1]])
9
10 b = Matrix([10, 14, 11, 0, 0, 0])
11 c = Matrix([6, 14, 13])
12 r = Matrix([0, -1, 0])
13
14 B = [0, 1, 2]
15
16 A_B = A[B, :]
17 b_B = b[B, :]
18
19 x = A_B.solve(b_B)
20 l = A_B.transpose().solve(c)
21 d = A_B.transpose().solve(r)
22
23 L = {"A_B" : A_B, "b_B" : b_B, "x^*" : x, "\lambda_B" : l, "d" :
24      : d}
25
26 for key, value in L.items():
27     print (key)
28     pprint(value)

```

Fig. 4.2: A partial python code for example 4.2.

Algorithm 4.1 (Simplex algorithm).

Start with feasible basis  $B$

while  $B$  is not optimal

Let  $i \in B$  be index with  $\lambda_i < 0$

Compute  $d \in \mathbb{R}^n$  with  $a_j^T d = 0$ ,  $j \in B \setminus \{i\}$  and  $a_i^T d = -1$

Determine  $K = \{k: 1 \leq k \leq m, a_k^T d > 0\}$

if  $K = \emptyset$

*assert LP unbounded*

else

Let  $k \in K$  index where  $\min_{k \in K} (b_k - a_k^T x^*) / a_k^T d$  is attained

update  $B := B \setminus \{i\} \cup \{k\}$

**Theorem 4.3.** *If the linear program (4.1) is non-degenerate, then the simplex algorithm terminates.*



*Proof.* In the non-degenerate case,  $\varepsilon^* > 0$  and the simplex algorithm makes progress, i.e., the objective function value strictly increases after each iteration. Since there is only a finite number of vertices, the algorithm terminates.

#### 4.4 Termination in the degenerate case

In the case where the linear program (4.1) is degenerate, we cannot argue that the objective function value increases each iteration and that the simplex algorithm terminates. However, the simplex algorithm leaves us some choice. Namely, there can be several indices  $i \in B$  such that  $\lambda_i < 0$ . Also, there could be several indices  $k \in K$  attaining the minimum in (4.6). If one adheres to the *smallest index rule*, then one can prove termination of the simplex algorithm also in the degenerate case. One iteration of the simplex algorithm is now as follows.

Algorithm 4.2 (Simplex algorithm with the smallest index rule).

Start with feasible basis  $B$

while  $B$  is not optimal

Compute  $\lambda \in \mathbb{R}^n$  such that  $\lambda^T A_B = c^T$

Let  $i^* \in B$  be the *smallest index* with  $\lambda_{i^*} < 0$

Compute  $d \in \mathbb{R}^n$  with  $a_j^T d = 0$ ,  $j \in B \setminus \{i^*\}$  and  $a_{i^*}^T d = -1$

Determine  $K = \{k: 1 \leq k \leq m, a_k^T d > 0\}$

if  $K = \emptyset$

*assert LP unbounded*

else

Let  $k^* \in K$  the *smallest index* where  $\min_{k \in K} (b_k - a_k^T x^*) / a_k^T d$  is attained

*update*  $B := B \setminus \{i^*\} \cup \{k^*\}$

**Theorem 4.4.** *The simplex algorithm with the smallest index rule terminates.*

*Proof.* We suppose that the simplex algorithm does not terminate. This means that the simplex algorithm iterates through a sequence of bases

$$B_0, B_1, \dots, B_k$$

with  $B_k = B_0$ . Inspecting two succeeding bases  $B_\ell$  and  $B_{\ell+1}$  for  $0 \leq \ell \leq k-1$ , one has  $B_{\ell+1} = B_\ell \setminus \{i\} \cup \{j\}$ , i.e.,  $i$  leaves and  $j$  enters  $B_\ell$ . Now let  $j$  be the *largest index* that leaves on that sequence. Since  $B_0 = B_k$ ,  $j$  also enters again at some point. Let  $p$  and  $q$  be the indices of bases,  $0 \leq p, q < k$  where  $j$  leaves and enters respectively.

Let  $\lambda^{(p)}$  and  $d^{(q)}$  be the corresponding  $\lambda$  and  $d$  vectors from the iteration  $p$  and  $q$  of the simplex algorithm respectively. Since  $\lambda^{(p)T} A = c^T$  and  $c^T d^{(q)} > 0$

we conclude

$$\lambda^{(p)T} A d^{(q)} > 0. \quad (4.7)$$

Let  $i \in B_p$  be an index with

$$\lambda_i^{(p)} a_i d^{(q)} > 0 \quad (4.8)$$

where  $a_i$  denotes the  $i$ -th row of  $A$ .

We now distinguish three cases. Let us suppose that  $i > j$ . Then, since  $j$  is the largest index that ever leaves or enters one has that  $i$  is also an element of  $B_q$  implying that  $a_i d^{(q)}$  is 0 or  $-1$ . It cannot be  $-1$ , since this would mean that  $i$  leaves  $B_q$  contradictory to the choice of  $j$ .

Suppose then that  $i < j$ . Then  $\lambda_i^{(p)} > 0$  since  $j$  is the smallest index that can leave the basis  $B_p$ . But  $a_i d^{(q)} > 0$  is not possible, otherwise index  $i$  is an index where the minimum ( $\varepsilon^* = 0$ ) in (4.6) is attained and  $j > i$  is the smallest index where this minimum is attained. Thus this case can also be ruled out.

Finally, if  $i = j$ , then  $\lambda_i^{(p)} < 0$  and  $a_i d^{(q)} > 0$  which also contradicts (4.8).

## 4.5 Finding an initial basic feasible solution

The simplex algorithm starts with a feasible basis. How can such a feasible basis be determined? In fact, this can be done with an auxiliary linear program.

Any linear program has an equivalent form

$$\max\{c^T x : Ax \leq b, x \geq 0\}. \quad (4.9)$$

We want to find an initial feasible basis for this linear program. Suppose first that we re-write the constraints  $Ax \leq b$  as  $A_1 x \leq b_1$  and  $A_2 x \leq b_2$  with  $b_1 \geq 0$  and  $b_2 < 0$ . Now consider the following linear program

$$\min\{\mathbf{1}^T y : A_1 x \leq b_1, A_2 x \leq b_2 + y, x, y \geq 0, y \leq -b_2\} \quad (4.10)$$

with constraints in matrix form

$$\begin{pmatrix} A_1 & & \\ A_2 & -I_k & \\ -I_n & & \\ & -I_k & \\ & & I_k \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \leq \begin{pmatrix} b_1 \\ b_2 \\ 0 \\ -b_2 \\ 0 \end{pmatrix}. \quad (4.11)$$

where  $k$  is the number of negative entries in  $b$ . An initial basic feasible solution is  $x = 0$  and  $y = -b_2$  with the basis corresponding to the inequalities  $x \geq 0$  and  $y \leq -b_2$ , i.e., the 3-rd and 5-th block of rows of the matrix (4.11).

The simplex algorithm applied to this linear program terminates. It finds an optimal solution with objective value 0 if and only if the original linear program is feasible. Let  $B$  be the optimal basis in this case. Then  $B$  without the indices corresponding to the constraints  $y \geq 0$  contains a feasible basis for the linear program (4.9). This follows from the fact that the rows of

$$\begin{pmatrix} A_1 \\ A_2 \\ -I_n \end{pmatrix} \quad (4.12)$$

indexed by  $B$  must correspond to a submatrix of rank  $n$ . If the corresponding rank was smaller than  $n$ , then the rank of the rows of the matrix (4.11) indexed by  $B$  would be smaller than  $n + k$ . Let  $B'$  be any subset of  $B$  of size  $n$  corresponding to a basis of 4.12. Then  $B'$  is an initial feasible basis of (4.9).

## 4.6 Removing degeneracy by perturbation

In the following, we derive an alternative pivoting rule that also ensures termination of the simplex algorithm. We begin by *perturbing* the right-hand-sides of our constraints

$$Ax \leq b \quad (4.13)$$

by adding a vector  $p_\varepsilon$

$$Ax \leq b + p_\varepsilon \quad (4.14)$$

where  $p_\varepsilon$  is the vector

$$p_\varepsilon = \begin{pmatrix} \varepsilon \\ \varepsilon^2 \\ \vdots \\ \varepsilon^m \end{pmatrix}.$$

**Lemma 4.1.** *If  $Ax \leq b$  is feasible, then  $Ax \leq b + p_\varepsilon$  is feasible for each  $\varepsilon > 0$ . If  $B$  is an infeasible basis of  $Ax \leq b$ , then  $B$  is an infeasible basis of  $Ax \leq b + p_\varepsilon$  for  $\varepsilon > 0$  sufficiently small.*

*Proof.* Clearly, the feasible solutions of (4.14) contain the feasible solutions of (4.13). Suppose now that  $B$  is infeasible, then there exists an index  $i$  and some  $\delta > 0$  such that

$$a_i A_B^{-1} b_B \geq b_i + \delta,$$

where  $a_i$  and  $b_i$  are the  $i$ -th row of  $A$  and the  $i$ -th component of  $b$  respectively. Now

$$a_i A_B^{-1} (b_B + p_{\varepsilon B}) \geq b_i + \delta + A_B^{-1} p_{\varepsilon B} > b_i$$

if  $\varepsilon > 0$  is sufficiently small.

Furthermore, we can show that the constraint system (4.14) is non-degenerate for  $\varepsilon > 0$  small enough.

**Lemma 4.2.** *If  $\varepsilon > 0$  is small enough, then (4.14) is non-degenerate.*

*Proof.* If the set of inequalities (4.14) is degenerate, then there exists a basis  $B$  and an index  $i \notin B$  such that

$$a_i x_{B,\varepsilon}^* = b_i + \varepsilon^i \quad (4.15)$$

where  $x_{B,\varepsilon}^* = A_B^{-1}(b_B + (p_\varepsilon)_B)$ . Notice that

$$a_i x_{B,\varepsilon}^* - b_i - \varepsilon^i$$

is a nonzero polynomial in  $\varepsilon$ . It is non-zero, since the coefficient of  $\varepsilon^i$  is  $-1$  as  $\varepsilon^i$  is not a component of  $(p_\varepsilon)_B$ . A nonzero polynomial has only a finite number of roots. Thus, if  $\varepsilon > 0$  is small enough, no equation of the form (4.15) can hold.

Suppose now that we want to solve the linear program

$$\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\} \quad (4.16)$$

The idea is run the simplex algorithm on the perturbed linear program

$$\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b + p_\varepsilon\} \quad (4.17)$$

where  $\varepsilon > 0$  is sufficiently small. This linear program is non-degenerate. What is nice is that this perturbation does not have to be computed explicitly. We can formulate a pivot rule for the non-perturbed linear program (4.16) that is conform with the pivoting that is performed on the perturbed linear program (4.15). For this, assume that  $B$  is a feasible basis of the perturbed linear program (4.15). By Lemma 4.1,  $B$  is also feasible for the unperturbed linear program. Let us now consider the iteration of the simplex algorithm at the basis  $B$ . As before, we choose  $i \in B$  with  $\lambda_i < 0$  arbitrary and compute  $d \in \mathbb{R}^n$ . Now we have to determine the unique index of the inequality of (4.17) that is hit first, when moving in the direction of  $d$  in the perturbed linear program.

This index can be determined as follows. As before, we determine  $K = \{k : 1 \leq k \leq m, a_k^T d > 0\}$  and we let  $\tilde{K} \subseteq K$  be those indices of  $K$  where the minimum  $\min_{k \in K} (b_k - a_k^T x^*)/a_k^T d$  is attained. In the perturbed program we would have to determine the unique minimum

$$b_k + \varepsilon^k - a_k^T A_B^{-1}(b_B + (p_\varepsilon)_B)/a_k^T d, \quad k \in \tilde{K}, \quad (4.18)$$

where we imagine that  $\varepsilon$  tends to zero from above. Each of the expressions in (4.18) is a polynomial in  $\varepsilon$ . The index that will leave in the perturbed

linear program is the one that corresponds to the lexicographically minimal polynomial in (4.18).

Algorithm 4.3 (Simplex with largest index leaving rule).

Start with feasible basis  $B$  of the *perturbed* linear program (4.17)

**while**  $B$  is not optimal

    Compute  $\lambda \in \mathbb{R}^n$  such that  $\lambda^T A_B = c^T$

    Let  $i \in B$  be index with  $\lambda_i < 0$

    Compute  $d \in \mathbb{R}^n$  with  $a_j^T d = 0$ ,  $j \in B \setminus \{i\}$  and  $a_i^T d = -1$

    Determine  $K = \{k: 1 \leq k \leq m, a_k^T d > 0\}$

**if**  $K = \emptyset$

        assert *LP unbounded*

**else**

        Let  $k \in K$  be the index corresponding to the lexicographically smallest polynomial of the form (4.18)

        update  $B := B \setminus \{i\} \cup \{k\}$

Exercise 6 explains how to convert a feasible basis of (4.16) into a feasible basis of (4.17).

**Theorem 4.5.** *The variant of the simplex method described in Algorithm 4.3 terminates.*

## Exercises

1. For each of the following assertion, provide a proof or a counterexample.
  - i) An index that has just left the basis  $B$  in the simplex algorithm cannot enter in the very next iteration.
  - ii) An index that has just entered the basis  $B$  in the simplex algorithm cannot leave again in the very next iteration.
2. Consider the auxiliary linear program to find an initial feasible basis (4.10). The constraint matrix of this linear program is of the form

$$\begin{pmatrix} A & 0 \\ -I_n & 0 \\ 0 & -I_{m_2} \\ 0 & I_{m_2} \end{pmatrix},$$

where  $m_2$  is the number of rows of  $A_2$ . This matrix has  $m+n+2 \cdot m_2$  rows. Describe an initial feasible basis that corresponds to the basic feasible solution  $x = 0$  and  $y = 0$ .

- Suppose that the optimal value of the auxiliary linear program is 0 and let  $B'$  be an optimal basis found by the simplex algorithm. Prove that  $B' \setminus \{m+n+1, \dots, m+n+m_2\}$  is a feasible *basis* of the linear program (4.9).
3. Suppose that the linear program  $\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\}$  is non-degenerate and  $B$  is an optimal basis. Show that the linear program has a unique optimal solution if and only if  $\lambda_B > 0$ .
  4. Let  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  be a polyhedron. Show that the following are equivalent for a feasible  $x^*$ :
    - i)  $x^*$  is a vertex of  $P$ .
    - ii) There exists a set  $B \subseteq \{1, \dots, m\}$  such that  $|B| = n$ ,  $A_B$  is invertible and  $A_B x^* = b_B$ . Here the matrix  $A_B$  and the vector  $b_B$  consists of the rows of  $A$  indexed by  $B$  and the components of  $b$  indexed by  $B$  respectively.
    - iii) For every feasible  $x_1, x_2 \neq x^* \in P$  one has  $x^* \notin \text{conv}\{x_1, x_2\}$ .
  5. A polyhedron  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  *contains a line*, if there exists a nonzero  $v \in \mathbb{R}^n$  and an  $x^* \in P$  such that for all  $\lambda \in \mathbb{R}$ , the point  $x^* + \lambda \cdot v \in P$ . Show that a nonempty polyhedron  $P$  contains a line if and only if  $A$  does not have full column-rank.
  6. Let  $x^*$  be a basic feasible solution of the non-perturbed linear program (4.16) and let  $C \subseteq \{1, \dots, m\}$  be the indices of inequalities that are tight at  $x^*$ . Prove that the following *greedy algorithm* produces a feasible basis  $B \subseteq C$  of the perturbed linear program.

Initialize  $B = \emptyset$

**while**  $B$  is not a basis

Let  $i \in C$  be the largest index such that the rows indexed by  $B \cup \{i\}$  are linearly independent

Update  $B = B \cup \{i\}$ .

## Chapter 5

# Duality

Via the termination argument for the simplex algorithm, we can now prove the duality theorem. We are given a linear program

$$\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\}, \quad (5.1)$$

called the *primal* and its *dual*

$$\min\{b^T y : y \in \mathbb{R}^m, A^T y = c, y \geq 0\}. \quad (5.2)$$

We again formulate the theorem of weak duality in this setting.

**Theorem 5.1 (Weak duality).** *If  $x^*$  and  $y^*$  are primal and dual feasible solutions respectively, then  $c^T x^* \leq b^T y^*$ .*

*Proof.* We have  $c^T x^* = y^{*T} A x^* \leq y^{*T} b$ . □

The strong duality theorem tells us that if there exist feasible primal and dual solutions, then there exist feasible primal and dual solutions which have the same objective value. We can prove it with the simplex algorithm.

**Theorem 5.2.** *If the primal linear program is feasible and bounded, then so is the dual linear program. Furthermore in this case, both linear programs have an optimal solution and the optimal values coincide.*

*Proof.* Suppose first that  $A$  has full column rank. The simplex method finds an optimal basis  $B$  of (5.1) with  $x_B^*$  being an optimal feasible solution. At the same time, we have a  $\lambda \in \mathbb{R}_{\geq 0}^m$  with  $\lambda_i = 0$  if  $i \notin B$  and  $\lambda^T A = c^T$ . Notice that  $\lambda$  is a feasible solution of the dual linear program (5.2). One has

$$c^T x_B^* = \lambda^T A x_B^* = \lambda^T b,$$

which shows the theorem in this case.

If  $A$  does not have full column rank, then we re-write the linear program (5.1) as

$$\max\{c^T(x_1 - x_2) : A(x_1 - x_2) \leq b, x_1 \geq 0, x_2 \geq 0\}. \quad (5.3)$$

There is a dual solution that we partition into three parts  $\lambda_1, \lambda_2, \lambda_3 \geq 0$ . Its dual objective function value is  $\lambda_1^T b$ . Furthermore

$$\lambda_1^T A - \lambda_2 = c^T, \quad -\lambda_1^T A - \lambda_3 = -c^T,$$

which together with  $\lambda_2, \lambda_3 \geq 0$  implies  $\lambda_1^T A = c^T$  and  $\lambda_2 = \lambda_3 = 0$ . This means that  $\lambda_1$  is an optimal dual solution.

We can formulate dual linear programs also if the linear program is not in inequality standard form. The procedure above can be described as follows. We transform a linear program into a linear program in inequality standard form and construct its dual linear program. This dual is then transformed into an equivalent linear program again which is conveniently described.

Let us perform such operations on the dual linear program

$$\min\{b^T y : y \in \mathbb{R}^m, A^T y = c, y \geq 0\}$$

of the primal  $\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\}$ . We transform it into inequality standard form

$$\begin{aligned} \max \quad & -b^T y \\ & A^T y \leq c \\ & -A^T y \leq -c \\ & -Iy \leq 0. \end{aligned}$$

The dual linear program of this is

$$\begin{aligned} \min \quad & c^T x_1 - c^T x_2 \\ & Ax_1 - Ax_2 - x_3 = -b \\ & x_1, x_2, x_3 \geq 0 \end{aligned}$$

This is equivalent to

$$\begin{aligned} \max \quad & c^T(x_2 - x_1) \\ & A(x_2 - x_1) + x_3 = b \\ & x_1, x_2, x_3 \geq 0 \end{aligned}$$

which is equivalent to the primal linear program

$$\begin{aligned} \max \quad & c^T x \\ & Ax \leq b. \end{aligned}$$

Loosely formulated one could say that “The dual of the dual is the primal”. But this, of course, is not to be understood as a mathematical statement. In any case we can state the following corollary.

**Corollary 5.1.** *If the dual linear program has an optimal solution, then so does the primal linear program and the objective values coincide.*



We present another example of duality that we will need later on. Consider a linear program

$$\begin{aligned} \max \quad & c^T x \\ Bx &= b \\ Cx &\leq d. \end{aligned} \tag{5.4}$$

After re-formulation, we obtain

$$\begin{aligned} \max \quad & c^T x \\ Bx &\leq b \\ -Bx &\leq -b \\ Cx &\leq d \end{aligned}$$

We can form the dual of the latter problem and obtain

$$\begin{aligned} \min \quad & b^T y_1 - b^T y_2 + d^T y_3 \\ B^T y_1 - B^T y_2 + C^T y_3 &= c \\ y_1, y_2, y_3 &\geq 0. \end{aligned}$$

But this linear program is equivalent to the linear program

$$\begin{aligned} \min \quad & b^T y_1 + d^T y_2 \\ B^T y_1 + C^T y_2 &= c \\ y_2 &\geq 0. \end{aligned} \tag{5.5}$$

This justifies to say that (5.5) is the dual of (5.4).

## 5.1 Zero sum games

Consider the following two-player game defined by a matrix  $A \in \mathbb{R}^{m \times n}$ . The *row-player* chooses a row  $i \in \{1, \dots, m\}$  and the *column-player* chooses a column  $j \in \{1, \dots, n\}$ . Both players make this choice at the same time. The *payoff* for the row-player is then the matrix-element  $A(i, j)$  whereas  $A(i, j)$  also determines the *loss* of the column player. In other words, the column player pays  $A(i, j)$  to the row-player. If this number is negative, then the row-player actually pays the absolute value of  $A(i, j)$  to the column player.

Consider for example the matrix

$$A = \begin{pmatrix} 5 & 1 & 3 \\ 3 & 2 & 4 \\ -3 & 0 & 1 \end{pmatrix}. \tag{5.6}$$

If the row-player chooses the second row and the column player chooses the second-column, then the payoff for the row-player is 2, whereas this is the loss of the column player.

The row-player is now interested in finding a strategy that maximizes his *guaranteed* payoff. For example, if he chooses row 1, then the best choice of the column player would be column 2, since the second element of the first row is the smallest element of that row. Thus the strategy that maximizes the minimal possible payoff would be to choose row 2. In other words

$$\max_i \min_j A(i, j) = 2.$$

What would be the column-player's best hedging strategy? He wants to choose a column such that the largest element in this column is minimized. This column would be the second one. In other words

$$\min_j \max_i A(i, j) = 2.$$

Is it always the case that  $\max_i \min_j A(i, j) = \min_j \max_i A(i, j)$ ? The next example shows that the answer is no:

$$\begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}. \quad (5.7)$$

Here we have  $\max_i \min_j A(i, j) = -1$  and  $\min_j \max_i A(i, j) = 1$ . This can be interpreted as follows. If the column player knows beforehand, the row to be chosen by the row-player, then he would choose a column that results in a gain for him. Similarly, if the row-player knows beforehand the column to be chosen by the column-player, then he can guarantee him a gain of one.

The idea is thus not to stick with a *pure* strategy, but to play with a *random* or *mixed* strategy. If the row-player chooses each of the two rows above uniformly at random, then his expected payoff is zero. Similarly, if the column player chooses each of his two columns with probability  $1/2$ , then his expected payoff is zero as well.

**Definition 5.1 (Mixed strategy).** Let  $A \in \mathbb{R}^{m \times n}$  define a two-player matrix game. A mixed strategy for the row-player is a vector  $x \in \mathbb{R}_{\geq 0}^m$  with  $\sum_{i=1}^m x_i = 1$ . A mixed strategy for the column player is a vector  $y \in \mathbb{R}_{\geq 0}^n$  with  $\sum_{j=1}^n y_j = 1$ .

Such mixed strategies define a probability distribution on the row and column indices respectively. If the row-player and column-player choose a row and column according to this distribution respectively, then the *expected payoff* for the row-player is

$$E[\text{Payoff}] = x^T A y. \quad (5.8)$$

For the game defined by (5.7) and  $x^T = (1/2, 1/2)$  and  $y^T = (1/2, 1/2)$  the expected payoff is 0.

**Lemma 5.1.** Let  $A \in \mathbb{R}^{m \times n}$ , then

$$\max_{x \in X} \min_{y \in Y} x^T A y \leq \min_{y \in Y} \max_{x \in X} x^T A y,$$

where  $X$  and  $Y$  denote the set of mixed row and column-strategies respectively.

*Proof.* Let  $x'$  and  $y'$  be some fixed mixed strategies of the row and column-player respectively. Clearly

$$\min_y x'^T A y \leq x'^T A y' \leq \max_x x^T A y',$$

which implies the assertion.  $\square$

The next theorem is one of the best-known results in the field of *game theory*. It states that there are mixed strategies  $x'$  and  $y'$  from above such that equality holds. It is proved with the theorem of strong duality.

**Theorem 5.3 (Minimax-Theorem).**

$$\max_{x \in X} \min_{y \in Y} x^T A y = \min_{y \in Y} \max_{x \in X} x^T A y,$$

where  $X$  and  $Y$  denote the set of mixed row and column-strategies respectively.

*Proof.* Let us inspect the value  $\max_{x \in X} \min_{y \in Y} x^T A y$ . This can be understood as to maximize the function

$$f(x) = \min\{(x^T A) \cdot y : \sum_{j=1}^n y_j = 1, y \geq 0\}.$$

Thus the value  $f(x)$  is the optimal solution of a bounded and feasible linear program. The dual of this linear program (for fixed  $x$ ) has only one variable  $x_0$  and reads

$$\max\{x_0 : x_0 \in \mathbb{R}, \mathbf{1}x_0 \leq A^T x\}.$$

But this shows that the maximum value of  $f(x)$ , where  $x$  ranges over all mixed row-strategies is the linear program

$$\begin{aligned} & \max x_0 \\ & \mathbf{1}x_0 - A^T x \leq 0 \\ & \sum_{i=1}^m x_i = 1 \\ & x \geq 0. \end{aligned} \tag{5.9}$$

Let us now inspect the value  $\min_{y \in Y} \max_{x \in X} x^T A y$ . Again, by applying duality this can be computed with the linear program

$$\begin{aligned} & \min y_0 \\ & \mathbf{1}y_0 - A y \geq 0 \\ & \sum_{j=1}^n y_j = 1 \\ & y \geq 0. \end{aligned} \tag{5.10}$$

It follows from the duality of (5.5) and (5.4) that the linear programs (5.9) and (5.10) are duals of each other. This proves the Minimax-Theorem.  $\square$

## 5.2 A proof of the duality theorem via Farkas' lemma

Remember Farkas' lemma (Theorem 3.11) which states that  $Ax = b, x \geq 0$  has a solution if and only if for all  $\lambda \in \mathbb{R}^m$  with  $\lambda^T A \geq 0$  one also has  $\lambda^T b \geq 0$ . In fact the duality theorem follows from this. First, we derive another variant of Farkas' lemma.

**Theorem 5.4 (Second variant of Farkas' lemma).** *Let  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ . The system  $Ax \leq b$  has a solution if and only if for all  $\lambda \geq 0$  with  $\lambda^T A = 0$  one has  $\lambda^T b \geq 0$ .*

*Proof.* Necessity is clear: If  $x^*$  is a feasible solution,  $\lambda \geq 0$  and  $\lambda^T A = 0$ , then  $\lambda^T Ax^* \leq \lambda^T b$  implies  $0 \leq \lambda^T b$ .

On the other hand,  $Ax \leq b$  has a solution if and only if

$$Ax^+ - Ax^- + z = b, x^+, x^-, z \geq 0 \quad (5.11)$$

has a solution. So, if  $Ax \leq b$  does not have a solution, then also (5.11) does not have a solution. By Farkas' lemma, there exists a  $\lambda \in \mathbb{R}^m$  with  $\lambda^T (A \ -A \ I_m) \geq 0$  and  $\lambda^T b < 0$ . For this  $\lambda$  one also has  $\lambda^T A = 0$  and  $\lambda \geq 0$ .  $\square$

We are now ready to prove the theorem of strong duality via the second variant of Farkas' lemma.

*Proof (of strong duality via Farkas' lemma).* Let  $\delta$  be the objective function value of an optimal solution of the dual  $\max\{b^T y : y \in \mathbb{R}^m, A^T y \leq c\}$ . For all  $\varepsilon > 0$ , the system  $A^T y \leq c, -b^T y \leq -\delta - \varepsilon$  does not have a solution. By the second variant of Farkas' lemma, there exists a  $\lambda \geq 0$  with  $\lambda^T \begin{pmatrix} -b^T \\ A^T \end{pmatrix} = 0$  and  $\lambda^T \begin{pmatrix} -\delta - \varepsilon \\ c \end{pmatrix} < 0$ . Write  $\lambda$  as  $\lambda = \begin{pmatrix} \lambda_1 \\ \lambda' \end{pmatrix}$  with  $\lambda' \in \mathbb{R}^n$ . If  $\lambda_1$  were zero, we could apply the second variant of Farkas' lemma to the system  $A^T y \leq c$  and  $\lambda'$ , since we know that  $A^T y \leq c$  has a solution. Therefore, we can conclude  $\lambda_1 > 0$ . Furthermore, by scaling, we can assume  $\lambda_1 = 1$ . One has  $\lambda'^T A^T = b^T$  and  $\lambda'^T c < \delta + \varepsilon$ . The first equation implies that  $\lambda'$  is a feasible solution of the primal (recall  $\lambda' \geq 0$ ). The second equation shows that the objective function value of  $\lambda'$  is less than  $\delta + \varepsilon$ . This means that the optimum value of the primal linear program is also  $\delta$ , since the primal has an optimal solution and  $\varepsilon$  can be chosen arbitrarily small.  $\square$

### Exercises

1. Formulate the dual linear program of

$$\begin{aligned} \max \quad & 2x_1 + 3x_2 - 7x_3 \\ \text{s.t.} \quad & x_1 + 3x_2 + 2x_3 = 4 \\ & x_1 + x_2 \leq 8 \\ & x_1 - x_3 \geq -15 \\ & x_1, x_2 \geq 0 \end{aligned}$$

2. Consider the following linear program

$$\begin{aligned} \max \quad & x_1 + x_2 \\ \text{s.t.} \quad & 2x_1 + x_2 \leq 6 \\ & x_1 + 2x_2 \leq 8 \\ & 3x_1 + 4x_2 \leq 22 \\ & x_1 + 5x_2 \leq 23 \end{aligned}$$

Show that  $(4/3, 10/3)$  is an optimal solution by providing a suitable feasible dual solution.

3. Show that for  $A \in \mathbb{R}^{m \times n}$ , one has

$$\max_i \min_j A(i, j) \leq \min_j \max_i A(i, j).$$

4. In the lecture you have seen the simplex algorithm for linear programs of the form

$$\max \{c^T x : Ax \leq b\}.$$

We will now derive a simplex algorithm for linear programs of the form

$$\min \{c^T x : Ax = b, x \geq 0\} \quad (5.12)$$

with  $c \in \mathbb{R}^n$  and  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ . Throughout the exercise we assume that (5.12) is feasible and bounded, and that  $A$  has full row rank.

For  $i \in \{1, \dots, n\}$  we define  $A_i$  as the  $i$ -th column of  $A$ . Moreover, for some subset  $B \subseteq \{1, \dots, n\}$ ,  $A_B$  is the matrix  $A$  restricted to the columns corresponding to elements of  $B$ .

A subset  $B \subseteq \{1, \dots, n\}$  with  $|B| = m$  such that  $A_B$  has full rank is called a *basis*. The vector  $x \in \mathbb{R}^n$  defined as  $x_i := 0$  for all  $i \notin B$  and  $x_B := A_B^{-1}b$  is called the *basic solution* associated to  $B$ . Note that  $x$  is a feasible solution to (5.12) if and only if  $x \geq 0$ .

Given a basis  $B$  and let  $j \in \{1, \dots, n\}$ ,  $j \notin B$ . The vector  $d \in \mathbb{R}^n$  defined as  $d_j = 1$ ,  $d_i = 0$  for all  $i \notin B$  and  $d_B := -A_B^{-1}A_j$  is called the  *$j$ -th basic direction*.

Assume that the solution  $x$  associated to  $B$  is feasible. Moreover assume that  $x_B > 0$ .

- a. Show that there is a  $\theta > 0$  such that  $x + \theta d$  is a feasible solution. Give a formula to compute the largest  $\theta$  such that  $x + \theta d$  is feasible.
- b. Let  $\theta^*$  be maximal. Show that there is a basis  $B'$  such that  $x + \theta^* d$  is the basic solution associated to  $B'$ .
- c. Let  $x' = x + \theta d$ . Show that the objective value of  $x'$  changes by  $\theta (c_j - c_B^T A_B^{-1} A_j)$ .
- d. Consider a basis  $B$  with basic feasible solution  $x$ . Show that if  $c - c_B^T A_B^{-1} A \geq 0$ , then  $x$  is an optimal solution to (5.12).

This suggests the following algorithm: Start with some basis  $B$  whose associated basic solution is feasible. Compute  $\bar{c} := c - c_B^T A_B^{-1} A$ . If  $\bar{c} \geq 0$ , we have an optimal solution (see 4d). Otherwise, let  $j$  be such that  $\bar{c}_j < 0$ . Part 4b and 4c show that if we change the basis, we find a feasible solution with an improved objective value. We repeat these steps until the vector  $\bar{c}$  is nonnegative.

This is the way the simplex algorithm usually is introduced in the literature. This algorithm is exactly the same as the one you learned in the lecture. To get an intuition why this is true, show the following:

- a. Given a basis  $B$ , show that its associated basic solution is feasible if and only if  $B$  is a *basis* of the LP dual to (5.12).
- b. Consider a basis  $B$  and its associated feasible basic solution  $x$ . As seen before,  $B$  is also a basis in the dual LP. Let  $y$  be the vertex of that basis. Show that for any  $j \in \{1, \dots, n\}$  we have  $\bar{c}_j < 0$  if and only if  $A_j^T y > c_j$ .

## Chapter 6

# Algorithms and running time analysis

An *algorithm* executes a set of *instructions* used in common programming languages like arithmetic operations, comparisons or read/write instructions. The sequence of these instructions is controlled by *loops and conditionals* like **if**, **while**, **for** etc.

Each of these instructions requires *time*. The *running time* of an algorithm is the *number of instructions* that the algorithm performs. This number depends on the *input* of the algorithm.

*Example 6.1.* Consider the following algorithm to compute the product of two  $n \times n$  matrices  $A, B \in \mathbb{Q}^{n \times n}$ :

```
for  $i = 1, \dots, n$ 
  for  $j = 1, \dots, n$ 
     $c_{ij} := 0$ 
    for  $k = 1, \dots, n$ 
       $c_{ij} := c_{ij} + a_{ik} \cdot a_{kj}$ 
```

The number of additions that are carried out is  $n^2 \cdot (n - 1)$  and the number of multiplications is  $n^3$ . The number of store-instructions is  $n^2 \cdot (n + 1)$ . The number of read-instructions is of similar magnitude.

The above example shows that an *exact counting* is sometimes tedious. Looking at the algorithm however, you quickly agree that there exists *some constant*  $c \in \mathbb{R}_{>0}$  such that the algorithm performs at most  $c \cdot n^3$  instructions.

In the analysis of algorithms, one does usually not care so much about the constant  $c$  above in the beginning. There are sub-fields of algorithms where this constant however matters. Especially for algorithms on large data sets, where access to external data is costly. However, this is another story that does not concern us here. When we analyze algorithms, we are interested in the *asymptotic running time*.

### Definition 6.1 (*O*-notation).

Let  $T, f : \mathbb{N} \rightarrow \mathbb{R}_{\geq 0}$  be functions. We say

- $T(n) = O(f(n))$ , if there exist positive constants  $n_0 \in \mathbb{N}$  and  $c \in \mathbb{R}_{>0}$  with

$$T(n) \leq c \cdot f(n) \text{ for all } n \geq n_0.$$

- $T(n) = \Omega(f(n))$ , if there exist constants  $n_0 \in \mathbb{N}$  and  $c \in \mathbb{R}_{>0}$  with

$$T(n) \geq c \cdot f(n) \text{ for all } n \geq n_0.$$

- $T(n) = \Theta(f(n))$  if

$$T(n) = O(f(n)) \text{ and } T(n) = \Omega(f(n)).$$

*Example 6.2.* The function  $T(n) = 2n^2 + 3n + 1$  is in  $O(n^2)$ , since for all  $n \geq 1$  one has  $2n^2 + 3n + 1 \leq 6n^2$ . Here  $n_0 = 1$  and  $c = 6$ . Similarly  $T(n) = \Omega(n^2)$ , since for each  $n \geq 1$  one has  $2n^2 + 3n + 1 \geq n^2$ . Thus  $T(n)$  is in  $\Theta(n^2)$ .

We measure the running time of algorithms in terms of the *length of the input*. The matrices  $A$  and  $B$  that are the input of the matrix-multiplication algorithm of Example 6.1 consist of  $n^2$  numbers each. The algorithm runs in time  $O(n^3 = (n^2)^{3/2})$ .

What does it mean for an algorithm to be efficient? For us, this will mean that it runs in *polynomial time*. As a first definition of polynomial time algorithm, we say that an algorithm runs in *polynomial time*, if there exists a constant  $k$  such that the algorithm runs in time  $O(n^k)$ , where  $n$  is the length of the input of the algorithm.

However, we recall the *binary* representation of natural numbers  $n \in \mathbb{N}$ . A sequence of *bits*  $a_0, \dots, a_{k-1}$  with  $a_j \in \{0, 1\}$  for  $0 \leq j \leq k-1$  represents the number

$$\sum_{j=0}^{k-1} a_j \cdot 2^j.$$

Conversely, each positive natural number  $n \geq 1$  has the binary representation that is found recursively by the following process. If  $n = 1$ , then its representation is  $a_0 = 1$ . If  $n > 1$  and is even, then the sequence representing  $n$  is

$$0, b_0, \dots, b_{k-1},$$

where  $b_0, \dots, b_{k-1}$  is the representation of  $n/2$ . If  $n > 1$  and  $n$  is odd, then the sequence representing  $n$  is

$$1, b_0, \dots, b_{k-1},$$

where  $b_0, \dots, b_{k-1}$  is the representation of  $\lfloor n/2 \rfloor$ . This creates a representation with leading bit one, i.e.  $a_{k-1} = 1$ . By deleting *leading zeros*, i.e., ensuring  $a_{k-1} = 1$ , the representation of a natural positive number is unique, (see exercise 6.2).



*Example 6.3.* This example will make us revise our first definition of a polynomial-time algorithm. The input of this algorithm is a list of characters. Let us say that the list has  $n$  elements.

Input: A list  $L$  of characters

$s := 2$

**for**  $c \in L$ :

$s := s \cdot c$

**return**  $s$

Then clearly, the algorithm carries out a polynomial, even linear, number of operations in the input, if we consider an arithmetic operation also as a basic operation that can be carried out in constant time. However, the algorithm squares 2 repeatedly,  $n$  times to be precise. Thus, at the end, the variable  $s$  holds the number  $2^{2^n}$ . The number of *bits* in the binary representation that the return value is  $2^n$  which is *exponential* in the input length.

**Definition 6.2.** The *size* of an integer  $x$  is  $\text{size}(x) = \lceil \log(|x| + 1) \rceil$  and for  $x \in \mathbb{Q}$ ,  $\text{size}(x) = \text{size}(p) + \text{size}(q)$ , where  $x = p/q$  with  $p, q \in \mathbb{Z}$ ,  $q \geq 1$  and  $\gcd(p, q) = 1$ .

Here  $\gcd(p, q)$  denotes the *greatest common divisor* of  $p$  and  $q$  which we define precisely below. Thus the size of a number is asymptotically equal to the number of bits that are needed to store the number. We now provide the definition of what a polynomial-time algorithm is.

**Definition 6.3.** An algorithm is *polynomial time*, if there exists a constant  $k$  such that the algorithm performs  $O(n^k)$  operations on rational numbers whose size is bounded by  $O(n^k)$ . Here  $n$  is the number of bits that encode the input of the algorithm. We say that the algorithm runs in time  $O(n^k)$ .

We now use this definition to analyze the famous *Euclidean* algorithm that computes the *greatest common divisor* of two integers.

For  $a, b \in \mathbb{Z}$ ,  $b \neq 0$  we say  $b$  *divides*  $a$  if there exists an  $x \in \mathbb{Z}$  such that  $a = b \cdot x$ . We write  $b \mid a$ . For  $a, b, c \in \mathbb{Z}$ , if  $c \mid a$  and  $c \mid b$ , then  $c$  is a *common divisor* of  $a$  and  $b$ . If at least one of the two integers  $a$  and  $b$  is non-zero, then there exists a *greatest common divisor* of  $a$  and  $b$ . It is denoted by  $\gcd(a, b)$ .

How do we compute the greatest common divisor efficiently? The following is called *division with remainder*. For  $a, b \in \mathbb{Z}$  with  $b > 0$  there exist unique integers  $q, r \in \mathbb{Z}$  with

$$a = q \cdot b + r, \text{ and } 0 \leq r < b.$$

Now clearly, for  $a, b \in \mathbb{Z}$  with  $b > 0$  and  $q, r \in \mathbb{Z}$  as above one has  $\gcd(a, b) = \gcd(b, r)$ . This gives rise to the famous *Euclidean algorithm*.

Algorithm 6.1 (Euclidean algorithm).

Input: Integers  $a \geq b \geq 0$  not both equal to zero  
Output: The greatest common divisor  $\gcd(a, b)$   
**if**  $(b = 0)$  **return**  $a$   
**else**  
    Compute  $q, r \in \mathbb{N}$  with  $b > r \geq 0$  and  $a = q \cdot b + r$   

(division with remainder)

  
    **return**  $\gcd(b, r)$

**Theorem 6.1.** *The Euclidean algorithm runs in time  $O(n)$ .*

*Proof.* Suppose that  $a$  and  $b$  have at most  $n$  bits each. Clearly, the numbers in the course of the algorithm have at most  $n$  bits. Furthermore, if  $a \geq b$ , then  $r \leq a/2$ , where  $r$  is the remainder of the division of  $a$  by  $b$ . Thus each second iteration, the first parameter of the input has one bit less. Thus the number of operations is bounded by  $O(n)$ .

*Example 6.4.* The *determinant* of a matrix  $A \in \mathbb{R}^{n \times n}$  can be computed by the recursive formula

$$\det(A) = \sum_{i=1}^n (-1)^{1+i} a_{1i} \det(A_{1i}),$$

where  $A_{1i}$  is the  $(n-1) \times (n-1)$  matrix that is obtained from  $A$  by deleting its first row and  $i$ -th column. This yields the following algorithm.

Input:  $A \in \mathbb{Q}^{n \times n}$   
Output:  $\det(A)$   
**if**  $(n = 1)$   
    **return**  $a_{11}$   
**else**  
     $d := 0$   
    **for**  $j = 1, \dots, n$   
         $d := (-1)^{1+j} \cdot \det(A_{1j}) + d$   
    **return**  $d$

This algorithm is *recursive*. This basically means that a tree is constructed, where nodes of the tree correspond to recursive calls to the algorithm  $\det(\cdot)$  including the *root* call. Any node that corresponds to a matrix with  $k > 1$  rows and columns is expanded by adding  $k$  *child nodes* corresponding to the recursive calls that are executed. Once the lower-level nodes of the tree cannot be expanded anymore, one can evaluate the tree, in this case the determinant, in a bottom-up fashion. See Let  $T(n)$  denote the number of basic operations that the algorithm performs. Then  $T(1) \geq 1$  and

$$T(n) \geq n \cdot T(n-1),$$

which shows that  $T(n) \geq n! = 2^{\Omega(n \log n)}$  which is *exponential* in  $n$ .

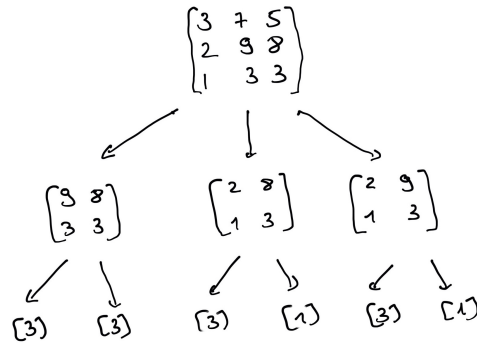


Fig. 6.1: An example of the recursion tree of the algorithm from Example 6.4. The tree corresponds to the run of the algorithm on input  $\begin{pmatrix} 3 & 7 & 5 \\ 2 & 9 & 8 \\ 1 & 3 & 3 \end{pmatrix}$ .

## Exercises

- Find the binary representation of 134.
- Show that the binary representation with leading bit one of a positive natural number is unique.
- Show that there are  $n$ -bit numbers  $a, b \in \mathbb{N}$  such that the Euclidean algorithm on input  $a$  and  $b$  performs  $\Omega(n)$  arithmetic operations. *Hint: Fibonacci numbers*
- Show  $n! = 2^{\Omega(n \log n)}$ .
- Let  $A \in \mathbb{R}^{n \times n}$  and suppose that the  $n^2$  components of  $A$  are pairwise different. Suppose that  $B$  is a matrix that can be obtained from  $A$  by deleting the first  $k$  rows and some of the  $k$  columns of  $A$ . How many (recursive) calls of the form  $\det(B)$  does the algorithm of Example 6.4 create?
- Let  $A \in \mathbb{R}^{n \times n}$  and suppose that the  $n^2$  components of  $A$  are pairwise different. How many different submatrices can be obtained from  $A$  by deleting the first  $k$  rows and some set of  $k$  columns? Conclude that the algorithm of Example 6.4 remains exponential, even if it does not expand repeated subcalls.
- Complete the algorithm below such that it adds two natural numbers in binary representation  $a_0, \dots, a_{l-1}, b_0, \dots, b_{l-1}$ . What is the asymptotic running time (number of basic operations) of your algorithm? Can there be an asymptotically faster algorithm?

Input: Two natural numbers  $a$  and  $b$  in their binary representation

$a_0, \dots, a_{l-1}, b_0, \dots, b_{l-1}$ .

Output: The binary representation  $c_0, \dots, c_l$  of  $a + b$

```

carry := 0
for  $i = 0, \dots, l - 1$ 
     $c_i = \text{carry} + a_i + b_i \pmod{2}$ 
    carry :=
 $c_l :=$ 
return  $c_0, \dots, c_l$ 

```

## 6.1 Analysis of Gaussian elimination

We recall Gaussian elimination.

Algorithm 6.2 (Gaussian elimination).

Input:  $A \in \mathbb{Q}^{m \times n}$

Output:  $A'$  in row echelon form such that there exists an invertible  $Q \in \mathbb{Q}^{m \times m}$  such that  $Q \cdot A = A'$ .

```

 $A' := A$ 
 $i := 1$ 
while  $(i \leq m)$ 
    find minimal  $1 \leq j \leq n$  such that there exists  $k \geq i$  such that  $a'_{kj} \neq 0$ 
    If no such element exists, then stop
    swap rows  $i$  and  $k$  in  $A'$ 
    for  $k = i + 1, \dots, m$ 
        subtract  $(a'_{kj}/a'_{ij})$  times row  $i$  from row  $k$  in  $A'$ 
     $i := i + 1$ 

```

We can easily prove correctness of the algorithm. First of all, the algorithm does only perform elementary row-operations of the form

- i) swap two rows
- ii) subtract a multiple of one row from *another* row.

This means that the resulting matrix  $A'$  can be obtained via

$$A' = Q \cdot A$$

with a non-singular  $Q \in \mathbb{Q}^{m \times m}$ . On the other hand we have the following invariant.

After each iteration of the while-loop, the matrix  $H$  obtained from the first first  $j$  columns of  $A'$  is in row-echelon form and rows  $i, i + 1, \dots, m$  of  $H$  are entirely zero, see Figure 6.2.

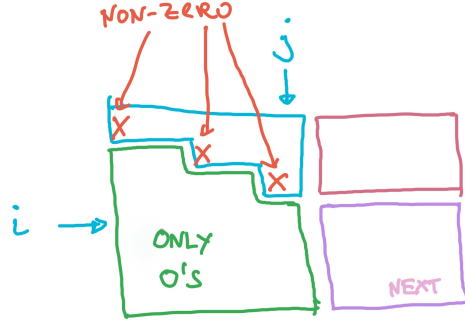


Fig. 6.2: Gaussian elimination: The matrix  $A'$  before the  $i$ -th iteration of the while loop.

How many arithmetic operations does Gaussian elimination perform? Subtracting a multiple of one row from another row in  $A'$  can be done in time  $O(n)$ . Thus the number of operations that are performed within the **for**-loop are  $O(m \cdot n)$  in total. There are  $O(m)$  iterations through the while-loop. Altogether this shows that Gaussian elimination performs  $O(m^2 \cdot n)$  iterations. But how large can the numbers grow in the course of the Gaussian algorithm? Could it be that numbers have to be manipulated, whose binary encoding length is not polynomial in the total encoding length of the matrix  $A$ ? Luckily the answer to this question is “No”. To provide this answer, we have to show the following.

**Theorem 6.2 (Hadamard bound).** *Let  $A \in \mathbb{R}^{n \times n}$  be non-singular. Then*

$$|\det(A)| \leq \prod_{i=1}^n \|a_i\|_2 \leq n^{n/2} \cdot B^n,$$

where  $B$  is upper bound on absolute values of entries of  $A$ .

*Proof.* The Gram-Schmidt orthogonalization of  $A$  yields a factorization

$$A = Q \cdot R,$$

where  $R$  is an upper triangular matrix with ones on the diagonal. The matrix  $Q$  has orthogonal columns, where the length of the  $i$ -th column  $q^{(i)}$  is upper bounded by the length of the  $i$ -th column of  $A$ . The assertion follows from

$$\det(A)^2 = \det(Q)^2 = \det(Q^T) \det(Q) = \prod_i \|q^{(i)}\|^2.$$

**Corollary 6.1.** *If  $A \in \mathbb{Z}^{n \times n}$  is integral and each entry in absolute value is bounded by  $B$ , then  $\text{size}(\det(A)) = O(n \log n + n \cdot \text{size}(B))$ .*

**Corollary 6.2.** *If  $A \in \mathbb{Q}^{n \times n}$  is a rational matrix and  $\phi$  is an upper bound on the size of each component of  $A$ , then  $\text{size}(\det(A)) = O(n^3 \cdot \phi)$ .*

*Proof.* Suppose that  $a_{ij} = p_{ij}/q_{ij}$ , where  $p_{ij}$  and  $q_{ij}$  are integers with  $\gcd(p_{ij}, q_{ij}) = 1$  for each  $i, j$ . Then  $(\prod_{ij} q_{ij})A$  is an integer matrix and

$$\det(A) = \left(\prod_{ij} q_{ij}\right)^{-n} \det\left(\left(\prod_{ij} q_{ij}\right)A\right).$$

The size of  $(\prod_{ij} q_{ij})^{-n}$  is  $O(n^3 \phi)$  and the size of  $\det((\prod_{ij} q_{ij})A)$  is also  $O(n^3 \phi)$  thanks to Corollary 6.1.

Now that we have shown that the determinant of a rational matrix is a number of polynomial encoding length, we can prove that Gaussian elimination is indeed a polynomial time algorithm.

**Theorem 6.3.** *The Gaussian algorithm runs in polynomial time on input  $A \in \mathbb{Z}^{m \times n}$ . More precisely, the rational numbers produced in the algorithm can be maintained to be ratios of sub-determinants of  $A'$  and are thus of polynomial binary encoding length.*

*Proof.* For the proof of this theorem, we assume that we have performed row and column swaps on  $A \in \mathbb{Z}$  beforehand such that the pivot element in iteration  $i \geq 1$  is in row  $i$  and column  $i$ . This means that we never have to swap rows in  $A'$ , after the  $i$ -th iteration of the while-loop the matrix  $A'$  is of the form

$$A' = \begin{pmatrix} U & E \\ 0 & D \end{pmatrix}$$

where  $U \in \mathbb{Q}^{i \times i}$  is upper triangular and non-singular and  $D \in \mathbb{Q}^{(m-i) \times (n-i)}$  is the part that still has to be eliminated.

Consider an element  $d_{k,j}$  of  $D$  and define the index sets  $K = \{1, \dots, i\} \cup \{i+k\}$  and  $J = \{1, \dots, i\} \cup \{i+j\}$  and the matrix  $A_{KJ}$  as the one induced from  $A$  by the rows and columns indexed by  $K$  and  $J$  respectively. A crucial observation is that

$$\det(A_{KJ}) = \det(A'_{KJ})$$

holds. Likewise one has  $\det(A_i) = \det(U)$ , where  $A_i$  is the matrix stemming from the first  $i$  rows and columns of  $A$ .

From this it follows that

$$d_{k,j} \cdot \det(A_i) = \det(A_{KJ})$$

and thus

$$d_{k,j} = \det(A_{KJ}) / \det(A_i).$$

We also write

$$d_{k,j} = r_{ij} / \det(A_i)$$

with the suitable integer  $r_{ij} = \det(A_{KJ})$ . Now the pivot element is  $d_{11}$  and the new element in position  $k, j$  of  $D$  becomes with  $r_{11} = \det(A_{i+1})$

$$\begin{aligned} d_{kj} - d_{k1}/d_{11} \cdot d_{1j} &= \frac{r_{kj}}{\det(A_i)} - \frac{r_{k1}}{r_{11}} \cdot \frac{r_{1j}}{\det(A_i)} \\ &= \frac{(r_{kj}r_{11} - r_{k1}r_{1j})/\det(A_i)}{\det(A_{i+1})} \end{aligned}$$

The numerator is an integer that can be computed with basic integer arithmetic operations. All integers are bounded by  $n^{n/2}B^n$ , where  $B$  is the largest binary encoding length of an entry of  $A$ . This shows the claim.

## Exercises

1. Show that the matrix  $Q \in \mathbb{Q}^{m \times m}$  that transforms  $A \in \mathbb{Q}^{m \times n}$  into  $A'$  in the Gaussian algorithm via  $Q \cdot A = A'$  has entries that are of polynomial size in the binary encoding length of  $A$ .

## 6.2 Fast matrix multiplication

We conclude this chapter on the analysis of algorithms with a result of Volker Strassen [18] who showed that two  $n \times n$  matrices can be multiplied in time (number of arithmetic operations)  $O(n^{2.805})$ . This algorithm was published in 1969 and it showed as well that a matrix can be inverted within the same timebound, see [1].

Our task is to compute the product

$$C = A \cdot B \tag{6.1}$$

for two  $n \times n$  matrices. We have seen that straightforward matrix-multiplication (Example 6.1) requires  $O(n^3)$  arithmetic operations. We can assume, by padding  $A$  and  $B$  with zeroes, that  $n = 2^\ell$  for some  $\ell \in \mathbb{N}_0$ .

If we split the matrices  $A$  and  $B$  into 4  $n/2 \times n/2$  matrices

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \text{ and } B = \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix} \tag{6.2}$$

Then

$$\begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix} = \begin{pmatrix} A_{11} \cdot B_{11} + A_{12} \cdot B_{21} & A_{11} \cdot B_{12} + A_{12} \cdot B_{22} \\ A_{21} \cdot B_{11} + A_{22} \cdot B_{21} & A_{21} \cdot B_{12} + A_{22} \cdot B_{22} \end{pmatrix}.$$

Thus we can reduce *one* multiplication of two  $n \times n$  matrices to 8 multiplications of two  $n/2 \times n/2$  matrices. For the running time, one obtains the recurrence

$$T(n) = 8 \cdot T(n/2) + \Theta(n^2) \quad (6.3)$$

which leads again to  $T(n) = \Theta(n^3)$ , see exercise 2).

Strassen discovered that one can preprocess the input matrices in such a way that the correct result can be retrieved from the result 7 multiplications of  $n/2 \times n/2$  matrices. The preprocessing and retrieval time is  $O(n)$  which yields the recursion

$$T(n) = 7 \cdot T(n/2) + O(n^2). \quad (6.4)$$

The idea is to compute the 7 matrices

$$\begin{aligned} M_1 &= (A_{11} + A_{22}) \cdot (B_{11} + B_{22}) \\ M_2 &= (A_{21} + A_{22}) \cdot B_{11} \\ M_3 &= A_{11} \cdot (B_{12} - B_{22}) \\ M_4 &= A_{22} \cdot (B_{21} - B_{11}) \\ M_5 &= (A_{11} + A_{12}) \cdot B_{22} \\ M_6 &= (A_{21} - A_{11}) \cdot (B_{11} + B_{12}) \\ M_7 &= (A_{12} - A_{22}) \cdot (B_{21} + B_{22}) \end{aligned}$$

This amounts to  $O(n^2)$  arithmetic operations and 7 multiplications of  $n/2 \times n/2$  matrices. From  $M_1, \dots, M_7$  one can retrieve

$$\begin{aligned} C_{11} &= M_1 + M_4 - M_5 + M_7 \\ C_{12} &= M_3 + M_5 \\ C_{21} &= M_2 + M_4 \\ C_{22} &= M_1 - M_2 + M_3 + M_6. \end{aligned}$$

again with  $O(n^2)$  arithmetic operations.

Algorithm 6.3 (Fast matrix multiplication (FMM)).

Input: Two  $n \times n$  matrices  $A$  and  $B$

Output:  $C = \text{FMM}(A, B)$ , the product  $A \cdot B$

**if**  $n = 1$  **return**  $a_{11} \cdot b_{11}$

**else**

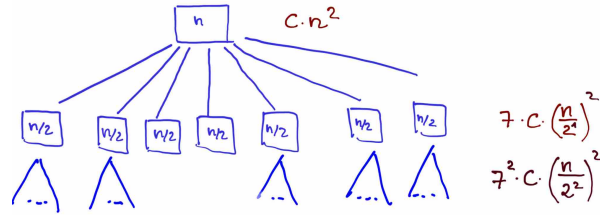
$M_1 = \text{FMM}(A_{11} + A_{22}, B_{11} + B_{22})$

$M_2 = \text{FMM}(A_{21} + A_{22}, B_{11})$

$M_3 = \text{FMM}(A_{11}, B_{12} - B_{22})$



$M_4 = \text{FMM}(A_{22}, B_{21} - B_{22})$   
 $M_5 = \text{FMM}(A_{11} + A_{12}, B_{22})$   
 $M_6 = \text{FMM}(A_{21} - A_{11}, B_{11} + B_{12})$   
 $M_7 = \text{FMM}(A_{12} - A_{22}, B_{21} + B_{22})$   
 Compute the matrices  $C_{11}, C_{12}, C_{21}, C_{22}$  from  $M_1, \dots, M_7$   
**return**  $C$



$$\begin{aligned}
 \text{Total: } & C \sum_{i=0}^{\log_2 n} 7^i \cdot \frac{n^2}{4^i} = C \cdot n^2 \cdot \sum_{i=0}^{\log_2 n} \left(\frac{7}{4}\right)^i \\
 & \leq C \cdot n^2 \cdot \left(\frac{7}{4}\right)^{\log_2 n} = C \cdot n^{2 + \log_2 \left(\frac{7}{4}\right)} = C \cdot n^{2.807}
 \end{aligned}$$

Fig. 6.3: The analysis of the Strassen algorithm.

**Theorem 6.4 (Strassen).** *Two  $n \times n$  matrices can be multiplied in time (number of arithmetic operations)  $O(n^{2+\log_2(7/4)})$ .*

*Proof.* See Figure 6.3.

## Exercises

1. Suppose we are given three  $n \times n$  matrices  $A, B, C \in \mathbb{Z}^{n \times n}$  and we want to test whether  $A \cdot B = C$  holds. We could multiply  $A$  and  $B$  and then compare the result with  $C$ . This would amount to running time (number of arithmetic operations) of  $O(n^3)$  with the standard matrix-multiplication algorithm.

We now show how to perform an efficient *randomized test*. Suppose that you can draw a vector  $v \in \{0, 1\}^n$  i.i.d. at random in time  $O(n)$ . The idea is then to compute the product  $B \cdot v$  and then the product  $A \cdot (B \cdot v)$  and afterwards  $C \cdot v$ , all in time  $O(n^2)$ . Show the following.

- a. If  $A \cdot B \neq C$ , then  $P(A \cdot (B \cdot v) = C \cdot v) \leq 1/2$ .
  - b. Let  $v_1, \dots, v_k \in \{0, 1\}^n$  be i.i.d. at random and suppose that  $A \cdot B \neq C$ . The probability of the event:  $A \cdot (B \cdot v_i) = C \cdot v_i$  for each  $i = 1, \dots, k$  is bounded by  $1/2^k$ .
  - c. Conclude that there is an algorithm that runs in time  $O(k \cdot n^2)$  which tests whether  $A \cdot B = C$  holds. The probability that the algorithm gives the wrong result is bounded by  $1/2^k$ .
2. Show that the recursion (6.3) has the solution  $T(n) = \Theta(n^3)$ .
  3. Describe an algorithm that multiplies two  $n$ -bit integers in time  $O(n^2)$ . You may use the algorithm to add two  $n$ -bit integers from exercise 6.7.
  4. Suppose  $n = 2^\ell$  and  $a, b \in \mathbb{N}$  are two  $n$ -bit integers. Consider the numbers  $a_h$  and  $a_l$  which are represented by the first  $n/2$  bits and the last  $n/2$  bits of  $a$  respectively. Likewise the numbers  $b_h$  and  $b_l$  are the numbers represented by the first half and the second half of the bit-representation of  $b$ .
    - i) Show  $a = a_h \cdot 2^{n/2} + a_l$  and  $b = b_h \cdot 2^{n/2} + b_l$
    - ii) Show  $a \cdot b = a_h \cdot b_h \cdot 2^n + (a_h \cdot b_l + a_l \cdot b_h) \cdot 2^{n/2} + a_l \cdot b_l$
    - iii) Conclude very carefully that two  $n$ -bit numbers can be multiplied by resorting to three multiplications of  $n/2$ -bit numbers and  $O(n)$  basic operations.
    - iv) Conclude that two  $n$ -bit numbers can be computed in time  $O(n^{\log_2(3)})$  elementary bit operations.

### 6.3 One iteration of the simplex algorithm

In the following we will analyze the complexity of one iteration of the simplex algorithm. We suppose that the input data  $A \in \mathbb{Z}^{m \times n}$ ,  $c \in \mathbb{Z}^n$ ,  $b \in \mathbb{Z}^m$  is integral.

Now if  $A_B^{-1}$  has been computed, then  $\lambda_B^T = c^T \cdot A_B^{-1}$  and  $d$ , which is the negative of a column of  $A_B^{-1}$ , can be computed with  $O(n^2)$  operations. To compute  $K$  we have to compute  $A \cdot d$ . This can be done with  $O(m \cdot n)$  operations. The index of element entering the basis can be determined by computing  $x^* = A_B^{-1}b_B$ ,  $b - Ax^*$  and  $Ad$ . Thus, if  $A_B^{-1}$  is known, this amounts to a total of

$$O(m \cdot n)$$

arithmetic operations.

In order to argue that one iteration of the simplex algorithm runs in polynomial time, we have to show that each of the numbers of  $A_B^{-1}$  has size that is polynomial in the size of the input and that  $A_B^{-1}$  can be quickly computed.

Let us first see how large the size of the numbers in  $A_B^{-1}$  can be. Suppose that

$$A = \begin{pmatrix} p_{11}/q_{11} & \cdots & p_{1n}/q_{1n} \\ & \cdots & \\ p_{n1}/q_{n1} & \cdots & p_{nn}/q_{nn} \end{pmatrix} \in \mathbb{Q}^{n \times n}.$$

is invertible. The size of the product of denominators  $\prod_{i=1, j=1}^n q_{ij}$  is clearly linear in the size of the input. Now write  $A = 1/Q \cdot A'$  where  $Q$  is product of denominators and  $A' \in \mathbb{Z}^{n \times n}$ . Since  $A^{-1} = Q \cdot (A')^{-1}$  we only have to answer this question for  $A'$  instead of  $A$ . In other words, we can assume that  $A$  is integral.

For  $A \in \mathbb{R}^{m \times n}$  and  $1 \leq i \leq m$  and  $1 \leq j \leq n$ ,  $A_{ij}$  denotes the matrix obtained from  $A$  by deleting the  $i$ -th row and  $j$ -th column. The following matrix inversion formula is known as Cramer's rule.

$$A^{-1} = \frac{1}{\det(A)} \begin{pmatrix} \det(A_{11}) & -\det(A_{21}) & \det(A_{31}) & \cdots \\ -\det(A_{12}) & \det(A_{22}) & -\det(A_{32}) & \cdots \\ \det(A_{13}) & -\det(A_{23}) & \det(A_{33}) & \cdots \\ \vdots & \vdots & \vdots & \cdots \\ \vdots & \vdots & \vdots & \cdots \end{pmatrix}$$

**Theorem 6.5 (Hadamard bound).** *Let  $A \in \mathbb{R}^{n \times n}$  be non-singular. Then*

$$|\det(A)| \leq \prod_{i=1}^n \|a_i\|_2 \leq n^{n/2} \cdot B^n,$$

where  $B$  is upper bound on absolute values of entries of  $A$ .

*Proof.* The Gram-Schmidt orthogonalization of  $A$  yields a factorization

$$A = Q \cdot R,$$

where  $R$  is an upper triangular matrix with ones on the diagonal. The matrix  $Q$  has orthogonal columns, where the length of the  $i$ -th column  $q^{(i)}$  is upper bounded by the length of the  $i$ -th column of  $A$ . The assertion follows from

$$\det(A)^2 = \det(Q)^2 = \det(Q^T) \det(Q) = \prod_i \|q^{(i)}\|^2.$$

**Corollary 6.3.** *If  $A \in \mathbb{Z}^{n \times n}$  is integral and each entry in absolute value is bounded by  $B$ , then  $\text{size}(\det(A)) = O(n \log n + n \cdot \text{size}(B))$ .*

**Corollary 6.4.** *Let  $A \in \mathbb{Q}^{n \times n}$  be an invertible matrix. The size of  $A^{-1}$  is polynomial in the size of  $A$ .*

Now we know that the size of  $A_B^{-1}$  is polynomial in the size of the input  $(A, b, c)$ ? Now, how expensive is it to compute  $A_B^{-1}$ ? Suppose basis  $B$  is preceded by  $B'$  with

$$\begin{aligned} B' &= \{b_1, \dots, b_{k-1}, b'_k, b_{k+1}, \dots, b_n\} \\ B &= \{b_1, \dots, b_{k-1}, b_k, b_{k+1}, \dots, b_n\} \end{aligned}$$

Then each row of  $A_B \cdot A_{B'}^{-1}$ , except for row  $k$ , is the corresponding row of the  $n \times n$  identity matrix except. Let the  $k$ -th row be  $(v_1, v_2, \dots, v_n)$ . We now only have to perform the elementary column operations that turn this row into the  $k$ -th unit vector on  $A_{B'}^{-1}$  to obtain  $A_B^{-1}$ . In other words, the following algorithm computes  $A_B^{-1}$  given  $A_{B'}^{-1}$ .

- Compute  $a_{b_k}^T \cdot A_{B'}^{-1} = (v_1, \dots, v_k, \dots, v_n)$
- For each column  $i \neq k$ : Subtract  $v_i/v_k$  times column  $k$  from column  $i$
- Divide column  $k$  by  $v_k$

This amounts to a total number of  $O(n^2)$  arithmetic operations for the update. We can conclude with the following theorem.

**Theorem 6.6.** *One iteration of the simplex algorithm requires a total number of  $O(m \cdot n)$  operations on rational numbers whose size is polynomial in the input size.*

## Chapter 7

# Integer Programming

An *integer program* is a problem of the form

$$\begin{aligned} \max \quad & c^T x \\ \text{s.t.} \quad & Ax \leq b \\ & x \in \mathbb{Z}^n, \end{aligned}$$

where  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ .

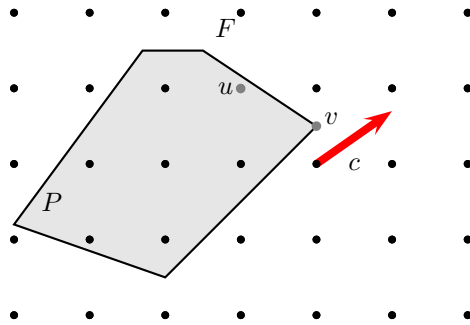


Fig. 7.1: This picture illustrates a polyhedron  $P$ , an objective function vector  $c$  and optimal points  $u, v$  of the integer program and the relaxation respectively.

The difference to linear programming is the *integrality constraint*  $x \in \mathbb{Z}^n$ . This powerful constraint allows to model discrete choices but, at the same time, makes an integer program much more difficult to solve than a linear program. In fact one can show that integer programming is NP-hard, which means that it is *in theory* computationally intractable. However, integer pro-

gramming has nowadays become an important tool to solve difficult industrial optimization problems efficiently. In this chapter, we characterize some integer programs which are easy to solve, since the *linear programming relaxation*  $\max\{c^T x : Ax \leq b\}$  yields already an optimal integer solution. The following observation is crucial.

**Theorem 7.1.** *Suppose that  $x^*$  is an integral optimal solution of the linear programming relaxation  $\max\{c^T x : Ax \leq b\}$ , i.e.,  $x^* \in \mathbb{Z}^n$ , then  $x^*$  is also an optimal solution of the integer programming problem  $\max\{c^T x : Ax \leq b, x \in \mathbb{Z}^n\}$*

Before we present an example for the power of integer programming we recall the definition of an undirected graph.

**Definition 7.1 (Undirected graph, matching).** An *undirected graph* is a tuple  $G = (V, E)$  where  $V$  is a finite set of elements, called the *vertices* or the *nodes*, and  $E \subseteq \binom{V}{2}$  is the set of *edges* of  $G$ . A *matching* of  $G$  is a subset  $M \subseteq E$  such that for all  $e_1 \neq e_2 \in M$  one has  $e_1 \cap e_2 = \emptyset$ .

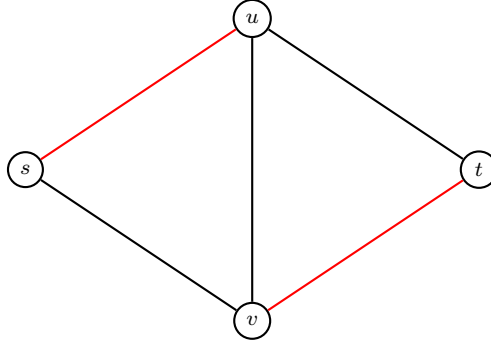


Fig. 7.2: A graph with 4 nodes  $V = \{s, u, v, t\}$  and 5 edges  $E = \{\{s, u\}, \{s, v\}, \{u, v\}, \{u, t\}, \{v, t\}\}$ . The red edges are a matching of the graph

We are interested in the solution of the following problem, which is called *maximum weight matching* problem. Given a graph  $G = (V, E)$  and a weight function  $w : E \rightarrow \mathbb{R}$ , compute a matching with maximum weight  $w(M) = \sum_{e \in M} w(e)$ .

For a vertex  $v \in V$ , the set  $\delta(v) = \{e \in E : v \in e\}$  denotes the *incident* edges to  $v$ . The maximum weight matching problem can now be modeled as an integer program as follows.

$$\begin{aligned}
 & \max \sum_{e \in E} w(e)x(e) \\
 & v \in V : \sum_{e \in \delta(v)} x(e) \leq 1 \\
 & e \in E : x(e) \geq 0 \\
 & x \in \mathbb{Z}^{|E|}.
 \end{aligned} \tag{7.1}$$

Clearly, if an integer vector  $x \in \mathbb{Z}^n$  satisfies the constraints above, then this vector is the *incidence vector* of a matching of  $G$ . In other words, the integral solutions to the constraints above are the vectors  $\{\chi^M : M \text{ matching of } G\}$ , where  $\chi_e^M = 1$  if  $e \in M$  and  $\chi_e^M = 0$  otherwise.

## 7.1 Integral Polyhedra

In this section we derive sufficient conditions on an integer program to be solved easily by an algorithm for linear programming. A central notion is the one of an integral polyhedron.

**Definition 7.2 (Valid inequality, face, vertex).** Let  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  be a polyhedron. An inequality  $c^T x \leq \beta$  is *valid* for  $P$  if  $c^T x^* \leq \beta$  for all  $x^* \in P$ . A *face* of  $P$  is a set of the form  $P \cap \{x \in \mathbb{R}^n : c^T x = \beta\}$  for a valid inequality  $c^T x \leq \beta$  of  $P$ . If a face consist of one point, then it is called a *vertex* of  $P$ .

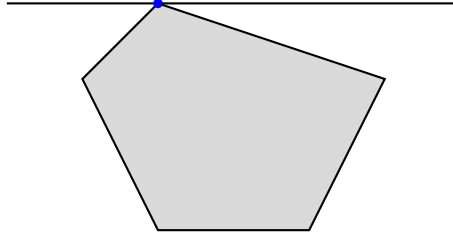


Fig. 7.3: A polyhedron with a valid inequality defining a vertex.

**Definition 7.3.** A rational polyhedron is called *integral* if each nonempty face of  $P$  contains an integer vector.

**Lemma 7.1.** Let  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  be an integral polyhedron with  $A \in \mathbb{R}^{m \times n}$  full-column rank. If the linear program

$$\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\} \quad (7.2)$$

is feasible and bounded, then the simplex method computes an optimal integral solution to the linear program.

*Proof.* Recall that, if the linear program (7.2) is bounded, the simplex method finds an optimal basis  $B \subseteq \{1, \dots, m\}$  of (7.2) and the vertex of the basis

$x_B^*$  is an optimal solution to (7.2). We have to show that  $x_B^*$  is integral. This will follow from the fact that  $\{x_B^*\}$  is a face of  $P$ .

Theorem 3.1 implies that  $x_B^*$  is the unique optimum solution of the linear program  $\max\{\tilde{c}^T x : x \in \mathbb{R}^n, a_i^T x \leq b_i, i \in B\}$ , where  $\tilde{c} = \sum_{i \in B} a_i$ . Consequently  $x_B^*$  is the unique solution of the linear program

$$\max\{\tilde{c}^T x : x \in P\}$$

which implies that  $\{x_B^*\}$  is a face defined by the valid inequality  $\tilde{c}^T x \leq \tilde{c}^T x_B^*$ .  $\square$

**Lemma 7.2.** *Let  $A \in \mathbb{Z}^{n \times n}$  be an integral and invertible matrix. One has  $A^{-1}b \in \mathbb{Z}^n$  for each  $b \in \mathbb{Z}^n$  if and only if  $\det(A) = \pm 1$ .*

*Proof.* Recall Cramer's rule which says  $A^{-1} = \tilde{A}/\det(A)$ , where  $\tilde{A}$  is the adjoint matrix of  $A$ . Clearly  $\tilde{A}$  is integral. If  $\det(A) = \pm 1$ , then  $A^{-1}$  is an integer matrix.

If  $A^{-1}b$  is integral for each  $b \in \mathbb{Z}^n$ , then  $A^{-1}$  is an integer matrix. We have  $1 = \det(A \cdot A^{-1}) = \det(A) \cdot \det(A^{-1})$ . Since  $A$  and  $A^{-1}$  are integral it follows that  $\det(A)$  and  $\det(A^{-1})$  are integers. The only divisors of one in the integers are  $\pm 1$ .  $\square$

**Definition 7.4 (Total unimodularity).** An integral matrix  $A \in \{0, \pm 1\}^{m \times n}$  is called *totally unimodular* if each of its square sub-matrices has determinant  $0, \pm 1$ .

**Theorem 7.2 (Hoffman-Kruskal Theorem).** *Let  $A \in \mathbb{Z}^{m \times n}$  be an integral matrix. The polyhedron  $P = \{x \in \mathbb{R}^n \mid Ax \leq b, x \geq 0\}$  is integral for each integral  $b \in \mathbb{Z}^m$  if and only if  $A$  is totally unimodular.*

[Cambristi Lemani] Moment Musical au bord du lac le 10 mai

*Proof.* Let  $A \in \mathbb{Z}^{m \times n}$  be totally unimodular and  $b \in \mathbb{Z}^m$ . Let  $x^*$  be vertex of  $P$  and suppose that this vertex is defined by the valid inequality  $c^T x \leq \delta$ . Notice that the matrix  $\begin{pmatrix} A \\ -I \end{pmatrix}$  has full column-rank. If one applies the simplex algorithm to the problem

$$\max\{c^T x : x \in \mathbb{R}^n, \begin{pmatrix} A \\ -I \end{pmatrix} x \leq \begin{pmatrix} b \\ 0 \end{pmatrix}\},$$

it finds an optimal basis  $B \subseteq \{1, \dots, m+n\}$  with  $x_B^* = x^*$ . If  $A_B$  denotes the matrix whose rows are those rows of  $\begin{pmatrix} A \\ -I \end{pmatrix}$  indexed by  $B$  and if  $b_B$  denotes the vector whose components are those of  $\begin{pmatrix} b \\ 0 \end{pmatrix}$  indexed by  $B$ , then  $x^* = A_B^{-1}b_B$ . We are done, once we conclude that  $\det(A_B) = \pm 1$ , since then  $A_B^{-1}$  is an integer matrix and since  $b_B$  is an integer vector  $x^* = A_B^{-1}b$  is integral as well. We can permute the columns of  $A_B$  in such a way that one obtains a matrix of the form

$$\begin{pmatrix} \bar{A} & \tilde{A} \\ 0 & -I_k \end{pmatrix}$$



where  $\bar{A}$  is a  $(n-k) \times (n-k)$  sub-matrix of  $A$  and  $I_k$  is the  $k \times k$  identity matrix. Here  $k = |B \cap \{m+1, \dots, m+n\}|$ . Clearly  $0 \neq \det(A_B) = \pm \det(\bar{A}) = \pm 1$ .

For the converse, suppose that  $A$  is not totally unimodular. Then there exists an index set  $B \subseteq \{1, \dots, m+n\}$  with  $|B| = n$  such that the matrix  $A_B$  defined as above satisfies  $|\det(A_B)| \geq 2$ . We can suppose w.l.o.g. that  $B = \{1, \dots, n\}$ . By Lemma 7.2 there exists choices for the components of  $b_B$  making  $A_B^{-1}b_B$  non-integral. In fact, if we split  $B$  into components  $L \subseteq B$  corresponding to lines of  $A$  and  $C$  corresponding to lines of  $-I$  we can choose those components of  $b_B$  corresponding to  $L$  being equal to zero. Now let  $v$  be the vector with  $v_i = 1$  for all  $i \in C$  and  $v_i = 0$  for all  $i \in L$ . By choosing  $\gamma \in \mathbb{N}$  large enough the point  $x_B^* = A_B^{-1}(b_B + \gamma A_B v)$  is non-integral and positive. Notice that starting from now we will consider a new vector  $\tilde{b}$  instead of  $b$ , where  $\tilde{b}_B = b_B$ . In the next lines we will say  $\tilde{b}_{\{1, \dots, m+n\} \setminus B}$  has to be to finish the proof. The set  $B$  is a basis of the linear program

$$\max\{\bar{c}^T x : x \in \mathbb{R}^n, \begin{pmatrix} A \\ -I \end{pmatrix} x \leq \begin{pmatrix} \tilde{b} \\ 0 \end{pmatrix}\},$$

where  $\bar{c} = \sum_{i \in B} a_i$  and  $a_i$  denotes the  $i$ -th row of  $\begin{pmatrix} A \\ -I \end{pmatrix}$ . If we define for  $j \in \{1, \dots, m\} \setminus B$ ,  $\tilde{b}_j = \lceil a_j^T x_B^* \rceil$ , then  $x_B^*$  is feasible and thus a vertex of  $P$  that is non-integral.  $\square$

A direct consequence of theorem 7.2 is the following corollary.

**Corollary 7.1.** *If  $A \in \mathbb{Z}^{m \times n}$  is totally unimodular,  $b \in \mathbb{Z}^m$  and if  $\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b, x \geq 0\}$  is bounded, then*

$$\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b, x \geq 0\} = \max\{c^T x : x \in \mathbb{Z}^n, Ax \leq b, x \geq 0\}.$$

## 7.2 Applications of total unimodularity

### 7.2.1 Bipartite matching

A graph is *bipartite*, if  $V$  has a partition into sets  $A$  and  $B$  such that each edge  $uv$  satisfies  $u \in A$  and  $v \in B$ . Recall that  $\delta(v)$  is the set of edges incident to the vertex  $v \in V$ , that is  $\delta(v) = \{e \in E \mid v \in e\}$ .

The *node-edge* incidence matrix of a graph  $G = (V, E)$  is the matrix  $A \in \{0, 1\}^{|V| \times |E|}$  with

$$A(v, e) = \begin{cases} 1, & \text{if } v \in e, \\ 0 & \text{otherwise.} \end{cases}$$

The integer program (7.1) can thus be formulated as

$$\max\{w^T x : Ax \leq 1, x \geq 0, x \in \mathbb{Z}^E\}. \quad (7.3)$$

The next lemma implies that the simplex algorithm can be used to compute a maximum-weight matching of a bipartite graph.

**Lemma 7.3.** *If  $G$  is bipartite, the node-edge incidence matrix of  $G$  is totally unimodular.*

*Proof (By induction).* Let  $G = (V, E)$  be a bipartite graph with bi-partition  $V = V_1 \cup V_2$ .

The case where  $A'$  is a  $1 \times 1$  sub-matrix of  $A$  is trivial. Suppose the lemma is proven for  $k - 1 \geq 1$  and let  $A'$  be a  $k \times k$  sub-matrix of  $A$ . We are interested in the determinant of  $A$ . Clearly, we can assume that  $A$  does not contain a column which contains no 1 or only one 1, since we simply consider the  $(k - 1) \times (k - 1)$  sub-matrix  $A''$  of  $A'$ , which emerges from developing the determinant of  $A'$  along this column. By the induction hypothesis the determinant of  $A'$  would be zero or  $\pm 1 \cdot \det(A'')$ .

Thus we can assume that each column contains exactly two ones. Now we can order the rows of  $A'$  such that the first rows correspond to vertices of  $V_1$  and then follow the rows corresponding to vertices in  $V_2$ . This re-ordering only affects the sign of the determinant. By summing up the rows of  $A'$  in  $V_1$  we obtain exactly the same row-vector as we get by summing up the rows of  $A'$  corresponding to  $V_2$ . This shows that  $\det(A') = 0$ .  $\square$

### 7.2.2 Bipartite vertex cover

A *vertex cover* of a graph  $G = (V, E)$  is a subset  $C \subseteq V$  of the nodes such  $e \cap C \neq \emptyset$  for each  $e \in E$ . Let us formulate an integer program for the *minimum-weight vertex-cover* problem. Here, one is given a graph  $G = (V, E)$  and weights  $w \in \mathbb{R}^V$ . The goal is to find a vertex cover  $C$  with minimum weight  $w(C) = \sum_{v \in V} w(v)$ .

$$\begin{aligned} \min \quad & \sum_{v \in V} w(v)x_v \\ \text{subject to} \quad & uv \in E : x_u + x_v \geq 1 \\ & v \in V : x_v \geq 0 \\ & x \in \mathbb{Z}^V. \end{aligned} \tag{7.4}$$

Clearly, this is the integer program

$$\min\{w^T x : A^T x \geq 1, x \geq 0, x \in \mathbb{Z}^V\}, \tag{7.5}$$

where  $A$  is the node-edge incidence matrix of  $G$ . A matrix  $A$  is totally unimodular if and only if  $A^T$  is totally unimodular. Thus the simplex algorithm can be used to compute a minimum weight vertex-cover of a bipartite graph. Furthermore we have the following theorem.

**Theorem 7.3 (König's theorem).** *In any bipartite graph, the number of edges in a maximum matching equals the number of vertices in a minimum vertex cover.*

*Proof.* Let  $A$  be the node-edge incidence-matrix of the bipartite graph  $G = (V, E)$ . The linear programs  $\max\{1^T x : Ax \leq 1, x \geq 0\}$  and  $\min\{1^T x : Ax \geq 1, x \geq 0\}$  are duals of each other. Since  $A$  is totally unimodular, the value of the linear programs are the cardinality of a maximum matching and minimum vertex-cover respectively. Thus the theorem follows from strong duality.  $\square$

### 7.2.3 Flows

Let  $G = (V, A)$  be a directed graph. The *node-edge incidence matrix* of a directed graph is a matrix  $A \in \{0, \pm 1\}^{V \times E}$  with

$$A(v, a) = \begin{cases} 1 & \text{if } v \text{ is the starting-node of } a, \\ -1 & \text{if } v \text{ is the end-node of } a, \\ 0 & \text{otherwise.} \end{cases} \quad (7.6)$$

A *feasible flow*  $f$  of  $G$  with capacities  $u$  and in-out-flow  $b$  is then a solution  $f \in \mathbb{R}^A$  to the system  $Af = b, 0 \leq f \leq u$ .

**Lemma 7.4.** *The node-edge incidence matrix  $A$  of a directed graph is totally unimodular.*

*Proof (By induction).* The case where  $A'$  is a  $1 \times 1$  sub-matrix of  $A$  is trivial. Let  $A'$  be a  $k \times k$  sub-matrix of  $A$  and suppose we have proven the lemma for every  $(k-1) \times (k-1)$  sub-matrix with  $k-1 \geq 1$ . Again, we can assume that in each column we have exactly one 1 and one  $-1$ . Otherwise, we develop the determinant along a column which does not have this property. But then, the matrix  $A'$  is singular, since adding up all rows of  $A'$  yields the 0-vector.

A consequence is that, if the  $b$ -vector and the capacities  $u$  are integral and an optimal flow exists, then there exists an integer optimal flow.

### 7.2.4 Doubly stochastic matrices

A matrix  $A \in \mathbb{R}^{n \times n}$  is *doubly stochastic* if it satisfies the following linear constraints

$$\begin{aligned} \sum_{i=1}^n A(i, j) &= 1, \forall j = 1, \dots, n \\ \sum_{j=1}^n A(i, j) &= 1, \forall i = 1, \dots, n \\ A(i, j) &\geq 0, \forall 1 \leq i, j \leq n. \end{aligned} \quad (7.7)$$

A permutation matrix is a matrix which contains exactly one 1 per row and column, where the other entries are all 0.

**Theorem 7.4.** *A matrix  $A \in \mathbb{R}^{n \times n}$  is doubly stochastic if and only if  $A$  is a convex combination of permutation matrices.*

*Proof.* Since a permutation matrix satisfies the constraints (7.7), then so does a convex combination of these constraints.

On the other hand it is enough to show that each vertex of the polytope defined by the system (7.7) is integral and thus a permutation matrix. However, the matrix defining the system (7.7) is the node-edge incidence matrix of the complete bipartite graph having  $2n$  vertices. Since such a matrix is totally unimodular, the theorem follows.

### 7.3 The matching polytope

We now come to a deeper theorem concerning the convex hull of matchings. We mentioned several times in the course that the maximum weight matching problem can be solved in polynomial time. We are now going to show a theorem of Edmonds [2] which provides a complete description of the matching polytope and present the proof by Lovász [11].

Before we proceed let us inspect the symmetric difference  $M_1 \Delta M_2$  of two matchings of a graph  $G$ . If a vertex is adjacent to two edges of  $M_1 \cup M_2$ , then one of the two edges belongs to  $M_1$  and one belongs to  $M_2$ . Also, a vertex can never be adjacent to three edges in  $M_1 \cup M_2$ . Edges which are both in  $M_1$  and  $M_2$  do not appear in the symmetric difference. We therefore have the following lemma.

**Lemma 7.5.** *The symmetric difference  $M_1 \Delta M_2$  of two matchings decomposes into node-disjoint paths and cycles, where the edges on these paths and cycles alternate between  $M_1$  and  $M_2$ .*

The *Matching polytope*  $P(G)$  of an undirected graph  $G = (V, E)$  is the convex hull of incidence vectors  $\chi^M$  of matchings  $M$  of  $G$ .

The incidence vectors of matchings are exactly the 0/1-vectors that satisfy the following system of equations.

$$\begin{aligned} \sum_{e \in \delta(v)} x_e &\leq 1 \quad \forall v \in V \\ x_e &\geq 0 \quad \forall e \in E. \end{aligned} \tag{7.8}$$

However the triangle (Figure 7.4) shows that the corresponding polytope is not integral. The objective function  $\max 1^T x$  has value 1.5. However, one can show that a maximum weight matching of an undirected graph can be computed in polynomial time which is a result of Edmonds [3].

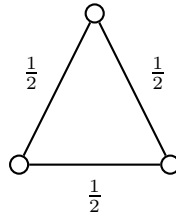


Fig. 7.4: Triangle

The following (Figure 7.5) is an illustration of an Edmonds inequality. Suppose that  $U$  is an odd subset of the nodes  $V$  of  $G$  and let  $M$  be a matching of  $G$ . The number of edges of  $M$  with both endpoints in  $U$  is bounded from above by  $\lfloor |U|/2 \rfloor$ .

Thus the following inequality is valid for the integer points of the polyhedron defined by (7.8).

$$\sum_{e \in E(U)} x_e \leq \lfloor |U|/2 \rfloor, \quad \text{for each } U \subseteq V, \quad |U| \equiv 1 \pmod{2}. \quad (7.9)$$

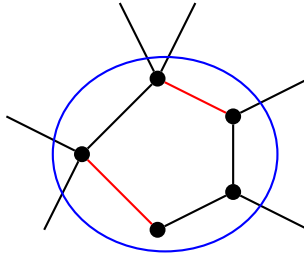


Fig. 7.5: Edmonds inequality.

The goal of this lecture is a proof of the following theorem.

**Theorem 7.5 (Edmonds 65).** *The matching polytope is described by the following inequalities:*

- i)  $x_e \geq 0$  for each  $e \in E$ ,
- ii)  $\sum_{e \in \delta(v)} x_e \leq 1$  for each  $v \in V$ ,
- iii)  $\sum_{e \in E(U)} x_e \leq \lfloor |U|/2 \rfloor$  for each  $U \subseteq V$

**Lemma 7.6.** Let  $G = (V, E)$  be connected and let  $w : E \rightarrow \mathbb{R}_{>0}$  be a weight-function. Denote the set of maximum weight matchings of  $G$  w.r.t.  $w$  by  $\mathcal{M}(w)$ . Then one of the following statements must be true:

- i)  $\exists v \in V$  such that  $\delta(v) \cap M \neq \emptyset$  for each  $M \in \mathcal{M}(w)$
- ii)  $|M| = \lfloor |V|/2 \rfloor$  for each  $M \in \mathcal{M}(w)$  and  $|V|$  is odd.

*Proof.* Suppose both i) and ii) do not hold. Then there exists  $M \in \mathcal{M}(w)$  leaving two exposed nodes  $u$  and  $v$ . Choose  $M$  such that the minimum distance between two exposed nodes  $u, v$  is minimized.

Now let  $t$  be on shortest path from  $u$  to  $v$ . The vertex  $t$  cannot be exposed.

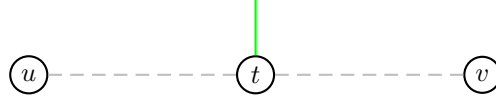


Fig. 7.6: Shortest path between  $u$  and  $v$ .

Let  $M' \in \mathcal{M}(w)$  leave  $t$  exposed. Both  $u$  and  $v$  are covered by  $M'$  because the distance to  $u$  or  $v$  from  $t$  is smaller than the distance of  $u$  to  $v$ .

Consider the symmetric difference  $M \Delta M'$  which decomposes into node disjoint paths and cycles. The nodes  $u, v$  and  $t$  have degree one in  $M \Delta M'$ . Let  $P$  be a path with endpoint  $t$  in  $M \Delta M'$



Fig. 7.7: Swapping colors.

If we swap colors on  $P$ , see Figure 7.7, we obtain matchings  $\widetilde{M}$  and  $\widetilde{M}'$  with  $w(M) + w(M') = w(\widetilde{M}) + w(\widetilde{M}')$  and thus  $\widetilde{M} \in \mathcal{M}(w)$ .

The node  $t$  is exposed in  $\widetilde{M}$  and  $u$  or  $v$  is exposed in  $\widetilde{M}$ . This is a contradiction to  $u$  and  $v$  being shortest distance exposed vertices

*Proof (Proof of Theorem 7.5).*

Let  $w^T x \leq \beta$  be a facet of  $P(G)$ , we need to show that this facet is of the form

- i)  $x_e \geq 0$  for some  $e \in E$
- ii)  $\sum_{e \in \delta(v)} x_e \leq 1$  for some  $v \in V$
- iii)  $\sum_{e \in E(U)} x_e \leq \lfloor |U|/2 \rfloor$  for some  $U \in P_{\text{odd}}$

To do so, we use the following method: One of the inequalities i), ii), iii) is satisfied with equality by each  $\chi^M$ ,  $M \in \mathcal{M}(w)$ . This establishes the claim since the matching polytope is full-dimensional and a facet is a maximal face.

If  $w(e) < 0$  for some  $e \in E$ , then each  $M \in \mathcal{M}(w)$  satisfies  $e \notin M$  and thus satisfies  $x_e \geq 0$  with equality.

Thus we can assume that  $w \geq 0$ .

Let  $G^* = (V^*, E^*)$  be the graph induced by edges  $e$  with  $w(e) > 0$ . Each  $M \in \mathcal{M}(w)$  contains maximum weight matching  $M^* = M \cap E^*$  of  $G^*$  w.r.t.  $w^*$ .

If  $G^*$  is not *connected*, suppose that  $V^* = V_1 \cup V_2$ , where  $V_1 \cap V_2 = \emptyset$  and  $V_1, V_2 \neq \emptyset$  and there is no edge connecting  $V_1$  and  $V_2$ , then  $w^T x \leq \beta$  can be written as the sum of  $w_1^T x \leq \beta_1$  and  $w_2^T x \leq \beta_2$ , where  $\beta_i$  is the maximum weight of a matching in  $V_i$  w.r.t.  $w_i$ ,  $i = 1, 2$ , see Figure 7.8. This would also contradict the fact that  $w^T x \leq \beta$  is a facet, since it would follow from the previous inequalities and thus would be a redundant inequality.



Fig. 7.8:  $G^*$  is connected.

Now we can use Lemma 7.6 for  $G^*$ .

- i)  $\exists v$  such that  $\delta(v) \cap M = \emptyset$  for each  $M \in \mathcal{M}(w)$ . This means that each  $M$  in  $\mathcal{M}(w)$  satisfies

$$\sum_{e \in \delta(v)} x_e \leq 1 \quad \text{with equality}$$

- ii)  $|M \cap E^*| = \lfloor |V^*|/2 \rfloor$  for each  $M \in \mathcal{M}(w)$  and  $|V^*|$  is odd. This means that each  $M$  in  $\mathcal{M}(w)$  satisfies

$$\sum_{e \in E(V^*)} x_e \leq \lfloor |V^*|/2 \rfloor \quad \text{with equality}$$

### Exercises

- Let  $M \in \mathbb{Z}^{n \times m}$  be totally unimodular. Prove that the following matrices are totally unimodular as well:

- $M^T$
- $\begin{pmatrix} M & I_n \end{pmatrix}$
- $\begin{pmatrix} M & -M \end{pmatrix}$
- $M \cdot (I_n - 2e_j e_j^T)$  for some  $j$

$I_n$  is the  $n \times n$  identity matrix, and  $e_j$  is the vector having a 1 in the  $j^{th}$  component, and 0 in the other components.

- A family  $\mathcal{F}$  of subsets of a finite groundset  $E$  is *laminar*, if for all  $C, D \in \mathcal{F}$ , one of the following holds:

$$(i) C \cap D = \emptyset, \quad (ii) C \subseteq D, \quad (iii) D \subseteq C.$$

Let  $\mathcal{F}_1$  and  $\mathcal{F}_2$  be two laminar families of the same groundset  $E$  and consider its union  $\mathcal{F}_1 \cup \mathcal{F}_2$ . Define the  $|\mathcal{F}_1 \cup \mathcal{F}_2| \times |E|$  adjacency matrix  $A$  as follows: For  $F \in \mathcal{F}_1 \cup \mathcal{F}_2$  and  $e \in E$  we have  $A_{F,e} = 1$ , if  $e \in F$  and  $A_{F,e} = 0$  otherwise.

Show that  $A$  is totally unimodular.

- Consider the following scheduling problem: Given  $n$  tasks with periods  $p_1, \dots, p_n \in \mathbb{N}$ , we want to find offsets  $x_i \in \mathbb{N}_0$ , such that every task  $i$  can be executed periodically at times  $x_i + p_i \cdot k$  for all  $k \in \mathbb{N}_0$ . In other words, for all pairs  $i, j$  of tasks we require  $x_i + k \cdot p_i \neq x_j + l \cdot p_j$  for all  $k, l \in \mathbb{N}_0$ .

Formulate the problem of finding these offsets as an integer program (with zero objective function).

- Show the following: A polyhedron  $P \subseteq \mathbb{R}^n$  with vertices is integral, if and only if each vertex is integral.
- Consider the polyhedron  $P = \{x \in \mathbb{R}^3 : x_1 + 2x_2 + 4x_3 \leq 4, x \geq 0\}$ . Show that this polyhedron is integral.
- Which of these matrices is totally unimodular? Justify your answer.

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

- Consider the complete graph  $G_n$  with 3 vertices, i.e.,  $G = (\{1, 2, 3\}, \binom{3}{2})$ . Is the polyhedron of the linear programming relaxation of the vertex-cover integer program integral?



## Chapter 8

# Paths, cycles and flows in graphs

Suppose you want to find a shortest path from a given starting point to a given destination. This is a common scenario in driver assistance systems (GPS) and can be modeled as one of the most basic combinatorial optimization problems, the *shortest path problem*. In this chapter, we introduce directed graphs, shortest paths and flows in networks. We focus in particular on the maximum-flow problem, which is a linear program that we solve with direct methods, versus the simplex method, and analyze the running time of these direct methods.

### 8.1 Graphs

**Definition 8.1.** A *directed graph* is a tuple  $G = (V, A)$ , where  $V$  is a finite set of elements, called the *vertices* of  $G$  and  $A \subseteq (V \times V)$  is the set of *arcs* of  $G$ . We denote an arc by its two defining nodes  $(u, v) \in A$ . The nodes  $u$  and  $v$  are called *tail* and *head* of the arc  $(u, v)$  respectively.

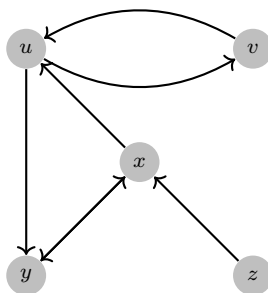


Fig. 8.1: Example of a directed graph with 5 nodes and 7 arcs.

**Definition 8.2 (Walk, path, distance).** A *walk* is a sequence of the form

$$P = (v_0, a_1, v_1, \dots, v_{m-1}, a_m, v_m),$$

where  $a_i = (v_{i-1}, v_i) \in A$  for  $i = 1, \dots, m$ . If the nodes  $v_0, \dots, v_m$  are all different, then  $P$  is a *path*. The *length* of  $P$  is  $m$ . The *distance* of two nodes  $u$  and  $v$  is the length of a shortest path from  $u$  to  $v$ . It is denoted by  $d(u, v)$ .

*Example 8.1.* The following is a walk and a path of the graph in Figure 8.1.

$$\begin{aligned} & z, (z, x), x, (x, u), u, (u, v), v, (v, u), u, (u, y), y, (y, x), x \\ & y, (y, x), x, (x, u), u, (u, v), v \end{aligned}$$

## 8.2 Representing graphs and computing the distance of two nodes

We represent a graph with  $n$  vertices  $v_1, \dots, v_n$  as an array  $A[v_1, \dots, v_n]$ , where the entry  $A[v_i]$  is a pointer to a linked list of vertices, the *neighbours* of  $v_i$ .  $N(v_i) = \{u \in V : (v_i, u) \in A\}$ .

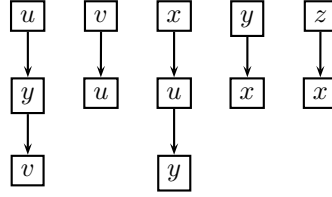


Fig. 8.2: Adjacency list representation of the graph in Figure 8.1.

### 8.2.1 Breadth-first search

We next describe a very basic algorithm that computes the distances from a designated node  $s \in V$  to all other nodes. The *distance* from  $s$  to  $v$  is denoted by  $d(s, v)$ . It is the smallest integer  $i$  such that there exists a path from  $s$  to  $v$  of length  $i$ . If there does not exist such a path, then  $s$  and  $v$  are *not connected* and we define  $d(s, v) = \infty$ . For  $i \in \mathbb{N}_0$ ,  $V_i \subseteq V$  denotes the set of vertices that have distance  $i$  from  $s$ . Notice that  $V_0 = \{s\}$ .

**Lemma 8.1.** *For  $i = 1, \dots, n-1$ , the set  $V_i$  is equal to the set of vertices  $v \in V \setminus (V_0 \cup \dots \cup V_{i-1})$  such that there exists an arc  $(u, v) \in A$  with  $u \in V_{i-1}$ .*

*Proof.* Suppose that  $v \notin V_0 \cup \dots \cup V_{i-1}$  and there exists an arc  $uv \in A$  with  $u \in V_{i-1}$ . Since  $u \in V_{i-1}$ , there exists a path  $s, a_1, v_1, a_2, v_2, \dots, a_{i-1}, u$  of length  $i-1$  from  $s$  to  $u$ . The sequence  $s, a_1, v_1, a_2, v_2, \dots, a_i, u, uv, v$  is a path of length  $i$  from  $s$  to  $v$  and thus  $v \in V_i$ .

If, on the other hand,  $v \in V_i$ , then there exists a path

$$s, a_1, v_1, \dots, a_{i-1}, u, a_i, v$$

of length  $i$  from  $s$  to  $v$ . We need to show that  $u \in V_{i-1}$  holds. Clearly, since there exists a path of length  $i-1$  from  $s$  to  $u$ , one has  $u \in V_j$  with  $j \leq i-1$ . If  $j < i-1$ , then there exists a path  $s, a'_1, v'_1, \dots, a'_j, u$  of length  $j$  which can be extended to a path of length  $j+1 < i$  from  $s$  to  $v$

$$s, a'_1, v'_1, \dots, a'_j, u, a_i, v$$

which contradicts  $v \in V_i$ .  $\square$

The *breadth-first search algorithm* is an implementation of Lemma 8.1. The algorithm maintains arrays

$$\begin{aligned} D[v_1 = s, v_2, \dots, v_n] \\ \pi[v_1 = s, v_2, \dots, v_n] \end{aligned}$$

and a queue  $Q$  that contains only  $s$  in the beginning. The array  $D$  contains at termination of the algorithm the distances from  $s$  to all other nodes and is initialized with  $[0, \infty, \dots, \infty]$ . The array  $\pi$  contains predecessor information for shortest paths, in other words, when the algorithm terminates,  $\pi[v] = u$ , where  $uv$  is an arc and  $D[u] + 1 = D[v]$ . The array  $\pi$  is initialized with  $[0, \dots, 0]$ .

After this initialization, the algorithm proceeds as follows.

```

while  $Q \neq \emptyset$ 
   $u := \text{head}(Q)$ 
  for each  $v \in \delta^+(u)$ 
    if  $(D[v] = \infty)$ 
       $\pi[v] := u$ 
       $D[v] := D[u] + 1$ 
       $\text{enqueue}(Q, v)$ 
   $\text{dequeue}(Q)$ 

```

Here the function  $\text{head}(Q)$  returns the next element in the queue and  $\text{dequeue}(Q)$  removes the first element of  $Q$ , while  $\text{enqueue}(Q, v)$  adds  $v$  to the queue  $Q$  as last element.

**Lemma 8.2.** *The breadth-first search algorithm assigns distance labels  $D$  correctly.*

*Proof.* We show the following claim by induction on  $i \in \{0, \dots, n-1\}$ .

For each  $i \in \{1, \dots, n-1\}$  there exists a point in time where:

- i)  $Q$  contains precisely the elements of  $V_i$
- ii) for each  $v \in V_i$ ,  $D[v] = d(s, v)$
- iii) for each  $v \in V_i$  one has  $\pi[v]v$  is an arc and  $\pi[v] \in V_{i-1}$ .

Once this claim is shown, the lemma follows, because the labels  $D[v]$  and  $\pi[v]$  are only changed once, if at all, from  $\infty$  or 0 to an integer or a vertex respectively.

Since  $V_0 = \{s\}$  and since  $Q = [s]$  and  $D[s] = 0$  after the initialization, the claim holds for  $i = 0$ . Suppose  $i > 0$ . By the induction hypothesis, there is a point in time, where  $Q$  contains precisely  $V_{i-1}$ . By Lemma 8.1, after the last element of  $V_{i-1}$  is dequeued  $Q$  contains precisely the elements in  $V_i$ . Also, since  $D[u] = d(s, u) = i-1$  for all  $u \in V_{i-1}$ , we have for each  $v \in V_i$  that  $D[v] = d(s, v) = i$ . Also  $\pi[v]v$  is an arc, by virtue of the algorithm, and  $\pi[v] \in V_{i-1}$ .  $\square$

**Definition 8.3 (Directed tree).** A *directed tree* is a directed graph  $T = (V, A)$  with  $|A| = |V| - 1$  and containing a node  $r \in V$  such that there exists a path from  $r$  to all other nodes of  $T$ .

**Lemma 8.3.** Consider the arrays  $D$  and  $\pi$  after the termination of the breadth-first-search algorithm. The graph  $T = (V', A')$  with  $V' = \{v \in V : D[v] < \infty\}$  and  $A' = \{\pi(v)v : 1 \leq D[v] < \infty\}$  is a tree.

*Proof.* Clearly,  $|A'| = |V'| - 1$ . For any  $i \in \{1, \dots, n-1\}$ , by backtracking the  $\pi$ -labels from any  $v \in V_i$ , we will eventually reach  $s$ .

**Definition 8.4.** The tree  $T$  from lemma 8.3 is the *shortest-path-tree* of the (unweighted) directed graph  $G = (V, A)$ .

**Theorem 8.1.** The breath-first-search algorithm runs in time  $O(|V| + |A|)$ .

*Proof.* Each vertex is queued and dequeued at most once. These queuing operations take constant time each. Thus queuing and dequeuing costs  $O(|V|)$  in total.

When a vertex  $u$  is dequeued, its neighbors are inspected and the operations in the **if** statement cost constant time each. Thus one has an additional cost of  $O(|A|)$ , since these constant-time operations are carried out for each arc  $a \in A$ .  $\square$

### 8.3 Shortest Paths

**Definition 8.5 (Cycle).** A walk in which starting node and end-node agree is called a *cycle*.

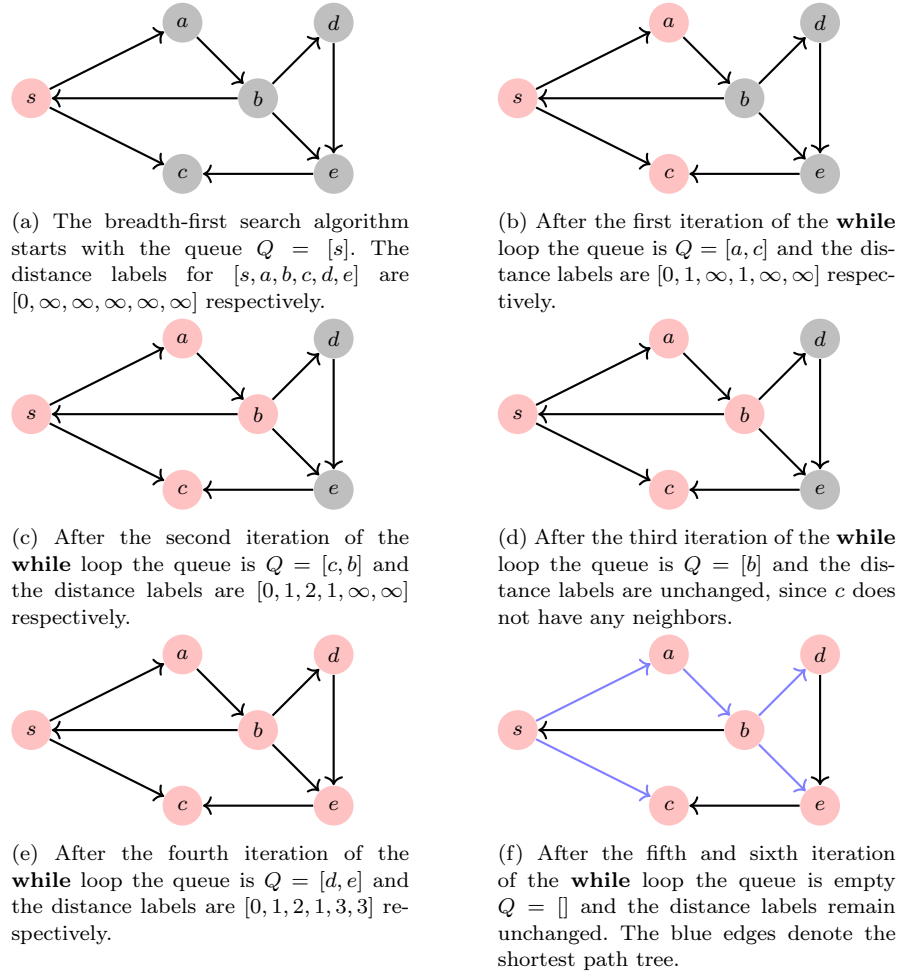


Fig. 8.3: An example-run of breadth-first search

Suppose we are given a directed graph  $D = (V, A)$  and a length function  $c : A \rightarrow \mathbb{R}$ . The *length* of a walk  $W$  is defined as

$$c(W) = \sum_{\substack{a \in A \\ a \in W}} c(a).$$

We now study how to determine a shortest path in a weighted directed graph  $G$  efficiently, in case of the absence of cycles of negative length (such cycles are called *negative cycles*).

**Theorem 8.2.** *Suppose that each cycle in  $D$  has non-negative length and suppose there exists an  $s - t$ -walk in  $D$ . Then there exists a path connecting  $s$  with  $t$  which has minimum length among all walks connecting  $s$  and  $t$ .*

*Proof.* If there exists an  $s - t$ -walk, then there exists an  $s - t$ -path. Since the number of arcs in a path is at most  $|V| - 1$ , there must exist a shortest path  $P$  connecting  $s$  and  $t$ . We claim that  $c(P) \leq c(W)$  for all  $s - t$ -walks  $W$ . Suppose that there exists an  $s - t$ -walk  $W$  with  $c(W) < c(P)$ . Then let  $W$  be such a walk with a minimum number of arcs. Clearly  $W$  contains a cycle  $C$ . Since the cycle has non-negative length, then it can be removed from  $W$  to obtain a walk whose length is at most  $c(W)$  and whose number of arcs is strictly less than  $|W|$ . This is a contradiction to the minimality of the number of arcs in  $W$ .  $\square$

We use the notation  $|W|, |C|, |P|$  to denote the number of arcs in a walk  $W$ , a cycle  $C$  or a path  $P$ .

As a conclusion we can note here:

If there do not exist negative cycles in  $D$ , and  $s$  and  $t$  are connected, then there exists a shortest walk traversing at most  $|V| - 1$  arcs.

### *The Bellman-Ford algorithm*

Let  $n = |V|$ . We calculate functions  $f_0, f_1, \dots, f_n : V \rightarrow \mathbb{R} \cup \{\infty\}$  successively by the following rule.

- i)  $f_0(s) = 0, f_0(v) = \infty$  for all  $v \neq s$
- ii) For  $k < n$  if  $f_k$  has been found, compute

$$f_{k+1}(v) = \min\{f_k(v), \min_{(u,v) \in A} \{f_k(u) + c(u, v)\}\}$$

for all  $v \in V$ .

The following theorem shows us that the Bellman-Ford algorithm is a method to compute minimum length walks.

**Theorem 8.3.** *For each  $k = 0, \dots, n$  and for each  $v \in V$*

$$f_k(v) = \min\{c(P) : P \text{ is an } s - v\text{-walk traversing at most } k \text{ arcs}\}.$$

**Theorem 8.4.** *Given a directed graph  $D = (V, A)$ ,  $s \in V$  and a length function  $c : A \rightarrow \mathbb{R}$ , one has  $f_n = f_{n-1}$  if and only if  $D$  does not have a cycle of negative length that is reachable from  $s$ .*

*Proof.* ( $\Leftarrow$ ) Suppose there exists  $t \in V$  such that  $f_n(t) < f_{n-1}(t) < \infty$ . This implies that the shortest  $s - t$ -walk traversing at most  $n$  arcs (call it  $W$ )

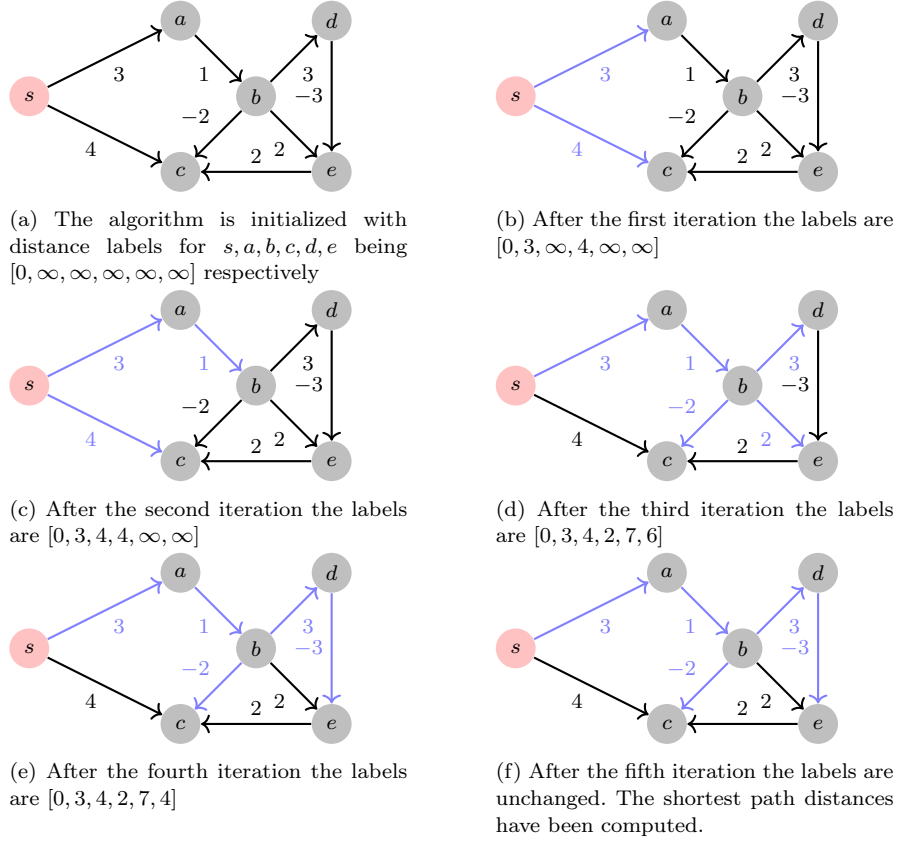


Fig. 8.4: An example-run of the Bellman-Ford algorithm. The blue edges represent the tree whose paths have the corresponding lengths.

traverses exactly  $n$  arcs and thus contains a cycle (call it  $C$ ). Consider the  $s - t$ -walk  $W'$  obtained by eliminating  $C$  from  $W$ .  $W'$  traverses at most  $n - 1$  arcs and thus we have  $c(W') > c(W)$ . Since  $c(W) = c(W') + c(C)$ , we have that  $c(C) < 0$ .

( $\Rightarrow$ ) Let  $C = v_0, v_1, \dots, v_k, v_0$  be a cycle reachable from  $s$ . Notice that  $f_{n-1}(v_i) < \infty \forall 0 \leq i \leq k$ . We can show that  $C$  is a non-negative cycle by these simple calculations (notice that the node indices  $i$  and  $i + 1$  are considered modulo  $k$  in the following sums).

$$0 = \sum_{i=0}^k f_n(v_{i+1}) - f_n(v_i) \leq \sum_{i=0}^k c(v_i, v_{i+1}) = c(C)$$

Theorem 8.4 can be generalized for every  $n \geq |V|$ . This allows us to obtain the following corollary.

**Corollary 8.1.** *If  $D = (V, A)$  does not contain negative cycles w.r.t.  $c$ , then  $f_n(v)$  is equal to the length of a shortest  $s - v$ -path. The numbers  $f_n(v)$  can be computed in time  $O(|V| \cdot |A|)$ .*

Notice that by using theorem 8.4 we can see if a directed graph  $D(V, A)$  has negative cycles in the following way. We obtain a new graph  $D' = (V', A')$  add a node  $s$  to  $D$  and we connect  $s$  to all other vertices with an outgoing arc of weight 0. Then, we apply Bellman-Ford algorithm to  $D'$  with starting node  $s$ . Notice that  $D$  contains a negative cycle, if and only if  $D'$  contains a negative cycle (since every cycle in  $D$  is reachable from  $s$  in  $D'$ ). Thus  $D$  contains a negative cycle if and only if there exists  $t \in V$  such that  $f_n(t) < f_{n-1}(t)$  in  $D'$ . This gives us the following corollary.

**Corollary 8.2.** *In time  $O(|V| \cdot |A|)$  one can test whether  $D = (V, A)$  has a negative cycle w.r.t.  $c$  and eventually return one.*

*Proof.* By our previous words we know that we simply have to apply Bellman-Ford algorithm to the graph  $D'$ . The corollary follows by the fact that  $|V'| = |V| + 1$  and  $|A'| = |A| + |V|$  and by using corollary 8.1.

## 8.4 Maximum $s - t$ -flows

We now turn our attention to a linear programming problem which we will solve by direct methods, motivated by the nature of the problem. We often use the following notation. If  $f : A \rightarrow B$  denotes a function and if  $U \subseteq A$ , then  $f(U)$  is defined as  $f(U) = \sum_{a \in U} f(a)$ .

**Definition 8.6 (Network,  $s - t$ -flow).** A network with capacities consists of a directed simple graph  $D = (V, A)$  and a *capacity function*  $u : A \rightarrow \mathbb{R}_{\geq 0}$ . We also require that if there is an arc  $uv \in A$ , then there is no reverse arc  $vu$ , and we disallow self-loops. A function  $f : A \rightarrow \mathbb{R}_{\geq 0}$  is called an  $s - t$ -flow, if

$$\sum_{e \in \delta^{out}(v)} f(e) = \sum_{e \in \delta^{in}(v)} f(e), \text{ for all } v \in V - \{s, t\}, \quad (8.1)$$

where  $s, t \in V$ . These two vertices are called *source* and *sink* respectively. The flow is *feasible*, if  $f(e) \leq u(e)$  for all  $e \in A$ . The *value* of  $f$  is defined as  $value(f) = \sum_{e \in \delta^{out}(s)} f(e) - \sum_{e \in \delta^{in}(s)} f(e)$ . The *maximum  $s - t$ -flow problem* is the problem of determining a maximum feasible  $s - t$ -flow.

Here, for  $U \subseteq V$ ,  $\delta^{in}(U)$  denotes the arcs which are entering  $U$  and  $\delta^{out}(U)$  denotes the arcs which are leaving  $U$ . Arc sets of the form  $\delta^{out}(U)$  are called a *cut* of  $D$ . The *capacity of a cut*  $u(\delta^{out}(U))$  is the sum of the capacities of its arcs.

Thus the maximum flow problem is a linear program of the form



$$\max \sum_{e \in \delta^{out}(s)} x(e) - \sum_{e \in \delta^{in}(s)} x(e) \quad (8.2)$$

$$\sum_{e \in \delta^{out}(v)} x(e) = \sum_{e \in \delta^{in}(v)} x(e), \text{ for all } v \in V - \{s, t\} \quad (8.3)$$

$$x(e) \leq u(e), \text{ for all } e \in A \quad (8.4)$$

$$x(e) \geq 0, \text{ for all } e \in A \quad (8.5)$$

**Definition 8.7 (excess function).** For any  $f : A \rightarrow \mathbb{R}$ , the excess function is the function  $excess_f : 2^V \rightarrow \mathbb{R}$  defined by  $excess_f(U) = \sum_{e \in \delta^{in}(U)} f(e) - \sum_{e \in \delta^{out}(U)} f(e)$ .

**Theorem 8.5.** Let  $D = (V, A)$  be a digraph, let  $f : A \rightarrow \mathbb{R}$  and let  $U \subseteq V$ , then

$$excess_f(U) = \sum_{v \in U} excess_f(v). \quad (8.6)$$

*Proof.* An arc which has both endpoints in  $U$  is counted twice with different parities on the right, and thus cancels out. An arc which has its tail in  $U$  is subtracted once on the right and once on the left. An arc which has its head in  $U$  is added once on the right and once on the left.  $\square$

A cut  $\delta^{out}(U)$  with  $s \in U$  and  $t \notin U$  is called an  $s - t$ -cut.

**Theorem 8.6 (Weak duality).** Let  $f$  be a feasible  $s - t$ -flow and let  $\delta^{out}(U)$  be an  $s - t$ -cut, then  $value(f) \leq u(\delta^{out}(U))$ .

*Proof.*  $value(f) = -excess_f(s) = -excess_f(U) = f(\delta^{out}(U)) - f(\delta^{in}(U)) \leq f(\delta^{out}(U)) \leq u(\delta^{out}(U))$ .  $\square$

For an arc  $a = (u, v) \in A$  the arc  $a^{-1}$  denotes the arc  $(v, u)$ .

**Definition 8.8 (Residual graph).** Let  $f : A \rightarrow \mathbb{R}$ , and  $u : A \rightarrow \mathbb{R}$  where  $0 \leq f \leq u$ . Consider the sets of arcs

$$A_f = \{a \mid a \in A, f(a) < u(a)\} \cup \{a^{-1} \mid a \in A, f(a) > 0\}. \quad (8.7)$$

The digraph  $D(f) = (V, A_f)$  is called the *residual graph* of  $f$  (for capacities  $u$ ).

**Corollary 8.3.** Let  $f$  be a feasible  $s - t$ -flow and suppose that  $D(f)$  has no path from  $s$  to  $t$ , then  $f$  has maximum value.

*Proof.* Let  $U$  be the set of nodes which are reachable in  $D(f)$  from  $s$ . Clearly  $\delta^{out}(U)$  is an  $s - t$ -cut. Now  $value(f) = f(\delta^{out}(U)) - f(\delta^{in}(U))$ . Each arc leaving  $U$  is not an arc of  $D(f)$  and thus  $f(\delta^{out}(U)) = u(\delta^{out}(U))$ . Each arc entering  $U$  does not carry any flow and thus  $f(\delta^{in}(U)) = 0$ . It follows that  $value(f) = u(\delta^{out}(U))$  and  $f$  is optimal by Theorem 8.6.  $\square$

**Definition 8.9 (undirected walk).** An *undirected walk* is a sequence of the form  $P = (v_0, a_1, v_1, \dots, v_{m-1}, a_m, v_m)$ , where  $a_i \in A$  for  $i = 1, \dots, m$  and  $a_i = (v_{i-1}, v_i)$  or  $a_i = (v_i, v_{i-1})$ . If the nodes  $v_0, \dots, v_m$  are all different, then  $P$  is an *undirected path*.

Any directed path  $P$  in  $D(f)$  yields an undirected path in  $D$ . Define for such a path  $P$  the vector  $\chi^P \in \{0, \pm 1\}^A$  as

$$\chi^P(a) = \begin{cases} 1 & \text{if } P \text{ traverses } a, \\ -1 & \text{if } P \text{ traverses } a^{-1}, \\ 0 & \text{if } P \text{ traverses neither } a \text{ or } a^{-1}. \end{cases} \quad (8.8)$$

**Theorem 8.7 (max-flow min-cut theorem, strong duality).** *The maximum value of a feasible  $s - t$ -flow is equal to the minimum capacity of an  $s - t$  cut.*

*Proof.* Let  $f$  be a maximum  $s - t$ -flow. Consider the residual graph  $D(f)$ . If this residual graph contains an  $s - t$ -path  $P$ , then we can route flow along this path. More precisely, there exists an  $\epsilon > 0$  such that  $f + \epsilon \chi^P$  is feasible. We have  $\text{value}(f + \epsilon \chi^P) = \text{value}(f) + \epsilon$ . This contradicts the maximality of  $f$  thus there exists no  $s - t$ -path in  $D(f)$ .

Let  $U$  be the nodes reachable from  $s$  in  $D(f)$ . Then  $\text{value}(f) = u(\delta^{\text{out}}(U))$  and  $\delta^{\text{out}}(U)$  is an  $s - t$ -cut of minimum capacity by the weak duality theorem.

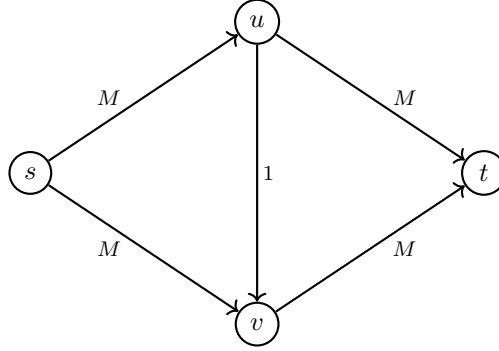
This suggests the algorithm of Ford and Fulkerson to find a maximum flow. Start with  $f = 0$ . Next iteratively apply the following *flow augmentation algorithm*.

Let  $P$  be a directed  $s - t$ -path in  $D(f)$ . Set  $f \leftarrow f + \epsilon \chi^P$ , where  $\epsilon$  is as large as possible to maintain  $0 \leq f \leq u$ .

**Exercise 8.1.** Define a *residual capacity* for  $D(f)$ . Then determine the maximum  $\epsilon$  such that  $0 \leq f \leq u$ .

**Theorem 8.8.** *If all capacities are rational, this algorithm terminates.*

*Proof.* By multiplying all the capacities by  $10^x$  with  $x$  sufficiently large, we can suppose the capacities to be integral. Since the value of the flow is augmented at every step at least by 1 and the capacities are finite, our algorithm finds a maximum flow in finitely many iterations.



The example above shows that, if the augmenting paths are chosen in a disadvantageous way, then the Ford-Fulkerson algorithm may take  $\Omega(M)$  iterations, where  $M$  is the largest capacity in the network. This happens if all augmenting paths use the arc  $uv$  or  $vu$  respectively in the residual network.

**Corollary 8.4 (integrality theorem).** *If  $u(a) \in \mathbb{N}$  for each  $a \in A$ , then there exists an integer maximum flow ( $f(a) \in \mathbb{N}$  for all  $a \in A$ ).*

*Proof.* This follows from the fact that the residual capacities remain integral and thus the augmented flow is always integral.  $\square$

**Theorem 8.9.** *If we choose in each iteration a shortest  $s - t$ -path in  $D(f)$  as a flow-augmenting path, the number of iterations is at most  $|V| \cdot |A|$ .*

**Definition 8.10.** Let  $D = (V, A)$  be a digraph,  $s, t \in V$  and let  $\mu(D)$  denote the length of a shortest path from  $s$  to  $t$ . Let  $\alpha(D)$  denote the set of arcs contained in at least one shortest  $s - t$  path.

**Theorem 8.10.** *Let  $D = (V, A)$  be a digraph and  $s, t \in V$ . Define  $D' = (V, A \cup \alpha(D)^{-1})$ . Then  $\mu(D) = \mu(D')$  and  $\alpha(D) = \alpha(D')$ .*

*Proof.* It suffices to show that  $\mu(D)$  and  $\alpha(D)$  are invariant if we add  $a^{-1}$  to  $D$  for one arc  $a \in \alpha(D)$ . Suppose not, then there is a directed  $s - t$ -path  $P_1$  traversing  $a^{-1}$  of length at most  $\mu(D)$ . As  $a \in \alpha(D)$  there is a path  $P_2$  traversing  $a$  of length  $\mu(D)$ . If we follow  $P_2$  until the tail of  $a$  is reached and from thereon follow  $P_1$ , we obtain another  $s - t$  path  $P_3$  in  $D$ . Similarly if we follow  $P_1$  until the head of  $a$  is reached and then follow  $P_2$ , we obtain a fourth  $s - t$  path  $P_4$  in  $D$ . However  $P_3$  or  $P_4$  has length less than  $\mu(D)$ . This is a contradiction.  $\square$

*Proof (of Theorem 8.9).* Let us augment flow  $f$  along a shortest  $s - t$ -path  $P$  in  $D(f)$  obtaining flow  $f'$ . The residual graph  $D_{f'}$  is a subgraph of  $D' = (V, A_f \cup \alpha(D(f))^{-1})$ . Hence  $\mu(D_{f'}) \geq \mu(D') = \mu(D(f))$ . If  $\mu(D_{f'}) = \mu(D(f))$ , then  $\alpha(D_{f'}) \subseteq \alpha(D') = \alpha(D(f))$ . At least one arc of  $P$  does not belong to  $D_{f'}$ , (the arc of minimum residual capacity!) thus the inclusion is strict. Since  $\mu(D(f))$  increases at most  $|V|$  times and, as long as  $\mu(D(f))$  does not change,  $|\alpha(D(f))|$  decreases at most  $2|A|$  times, we have the theorem.  $\square$

In the following let  $m = |A|$  and  $n = |V|$ .

**Corollary 8.5.** *A maximum flow can be found in time  $O(nm^2)$ .*

## 8.5 Minimum cost network flows, MCNFP

In contrast to the maximum  $s - t$ -flow problem, the goal here is to route a flow, which comes from several sources and sinks through a network with capacities and *costs* in such a way, that the total cost is minimized.

*Example 8.2.* Suppose you are given a directed graph  $D = (V, A)$  with arc weights  $c : A \rightarrow \mathbb{R}_{\geq 0}$  and your task is to compute a shortest path from a particular node  $s$  to all other nodes in the graph and assume that such paths exist. Then one can model this as a MCNFP (minimum cost network flow problem) by sending a flow of value  $|V| - 1$  into the source node and by letting a flow of value 1 leave each node. The costs on the arcs are defined by  $c$ . The arcs have infinite capacities. We will see later, that this MCNFP has an integral solution which corresponds to the shortest paths from  $s$  to all other nodes.

Here is a formal definition of a MCNFP. In this notation, vertices are indexed with the letters  $i, j, k$  and arcs are denoted by their tail and head respectively, for example  $(i, j)$  denotes the arc from  $i$  to  $j$ .

A *network* is now a directed graph  $D = (V, A)$  together with a capacity function  $u : A \rightarrow \mathbb{Q}_{\geq 0}$ , a cost function  $c : A \rightarrow \mathbb{Q}$  and an external flow  $b : V \rightarrow \mathbb{Q}$ . The value of  $b_i$  denotes the amount of flow which comes from the exterior. If  $b_i > 0$ , then there is flow from the outside, entering the network through node  $i$ . If  $b_i < 0$ , there is flow which leaves the network through  $i$ .

In the following we often use the notation  $f(i, j)$  for the flow-value on the arc  $(i, j)$  (instead of  $f((i, j))$ ). Similarly we write  $c(i, j)$  and  $u(i, j)$ .

A *feasible flow* is a function  $f : A \rightarrow \mathbb{Q}_{\geq 0}$  which satisfies the following constraints.

$$\begin{aligned} \sum_{e \in \delta^{out}(i)} f(e) - \sum_{e \in \delta^{in}(i)} f(e) &= b_i \quad \text{for all } i \in V, \\ 0 \leq f(e) \leq u(e) &\quad \text{for all } e \in A. \end{aligned}$$

The goal is to find a feasible flow with minimum cost:

$$\begin{aligned} &\text{minimize} && \sum_{e \in A} c(e)f(e) \\ \text{subject to} &&& \sum_{e \in \delta^{out}(i)} f(e) - \sum_{e \in \delta^{in}(i)} f(e) = b_i \quad \text{for all } i \in V, \\ &&& 0 \leq f(e) \leq u(e) \quad \text{for all } e \in A \end{aligned}$$

*Example 8.3.* Imagine you are a pilot and fly a passenger airplane in hops from airport 1 to airport 2 to airport 3 and so on, until airport  $n$ . At airport

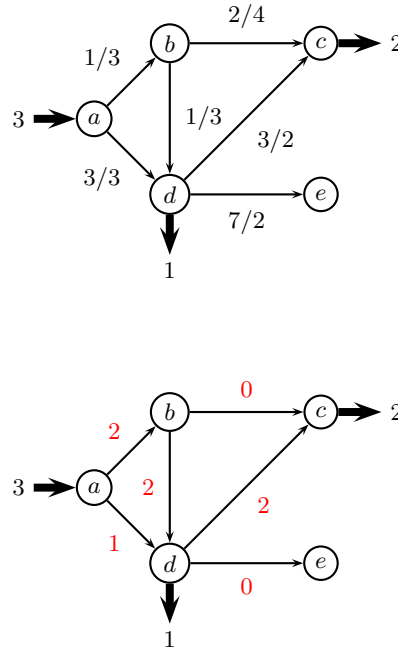


Fig. 8.5: A Network with in/out-flow, costs and capacities and a feasible flow of cost 13.

$i$  there are  $b_{ij}$  passengers that want to travel to airport  $j$ , where  $j > i$ . You may decide how many of the  $b_{ij}$  passengers you will take on board. Each of the passengers will pay  $c_{ij}$  dollars for the trip. The airplane can accommodate  $p$  people.

You are a greedy pilot and think of a plan to pick up and deliver passengers on your hop from 1 to  $n$  which maximizes your revenue.

Finding this plan can be modeled as a MCNFP. Your network has nodes  $1, \dots, n$  and arcs  $(i, i+1), i = 1, \dots, n-1$  with capacities  $p$  and without costs. These nodes do not have in/out-flow from the outside. You furthermore have nodes  $i \rightarrow j$  for  $i < j$  and  $i, j \in \{1, \dots, n\}$  which are excess nodes with in-flow  $b_{ij}$  from the outside. Each node  $i \rightarrow j$  is connected to  $i$  and to  $j$  with a directed arc. The capacities on these arcs are infinite. The cost of the arc  $(i \rightarrow j, i)$  is  $-c_{ij}$ . The cost of the arc  $(i \rightarrow j, j)$  is zero. The outflow on the node  $j$  is the total number of passengers that want to fly to node  $j$ . An integral optimal flow to this problem is an optimal plan for you.

Throughout this chapter we make the following assumptions.

1. All data (cost, supply, demand and capacity) are integral.
2. The network contains an incapacitated directed path between every pair of nodes.
3. The supplies/demands at the nodes satisfy the condition  $\sum_{i \in V} b_i = 0$  and the MCNFP has a feasible solution.
4. All arc costs are nonnegative.
5. The graph does not contain a pair of reverse arcs.

**Exercise 8.2.** Show how to transform a MCNFP on a digraph with pairs of reverse arcs into a MCNFP on a digraph with no pairs of reverse arcs. The number of arcs and nodes should asymptotically remain the same.

An *arc-flow* of  $D$  is a flow vector, that satisfies the nonnegativity and capacity constraints.

$$\sum_{e \in \delta^{in}(i)} f(e) - \sum_{e \in \delta^{out}(i)} f(e) = g(i) \quad \text{for all } i \in V,$$

$$0 \leq f(e) \leq u(e) \quad \text{for all } e \in A.$$

- If  $g(i) > 0$ , then  $i$  is an *excess node* (more inflow than outflow).
- If  $g(i) < 0$ , then  $i$  is a *deficit node* (more outflow than inflow).
- If  $g(i) = 0$  then  $i$  is a *balanced node*.

**Exercise 8.3.** Prove that  $\sum_{i \in V} g(i) = 0$  holds and thus that a feasible flow only exists if the sum of the  $b(i)$  is equal to zero.

Let  $\mathcal{P}$  be the collection of directed paths of  $D$  and let  $\mathcal{C}$  be the collection of directed cycles of  $D$ . A path-flow is a function  $\beta : \mathcal{P} \cup \mathcal{C} \rightarrow \mathbb{R}_{\geq 0}$  which assigns flow values to paths and cycles.

For  $(i, j) \in A$  and  $P \in \mathcal{P}$  let  $\delta_{(i,j)}(P)$  be 1 if  $(i, j) \in P$  and 0 otherwise. For  $C \in \mathcal{C}$  let  $\delta_{(i,j)}(C)$  be 1 if  $(i, j) \in C$  and 0 otherwise.

A path-flow  $\beta$  determines a unique arc-flow

$$f(i, j) = \sum_{P \in \mathcal{P}} \delta_{(i,j)}(P) \beta(P) + \sum_{C \in \mathcal{C}} \delta_{(i,j)}(C) \beta(C).$$

**Theorem 8.11.** *Every path and cycle flow has a unique representation as a nonnegative arc-flow. Conversely, every nonnegative arc-flow  $f$  can be represented as a path and cycle flow with the following properties:*

1. *Every directed path with positive flow connects a deficit node with an excess node.*
2. *At most  $n + m$  paths and cycles have nonzero flow and at most  $m$  cycles have nonzero flow.*

*If the arc-flow  $f$  is integral, then so are the path and cycle flows into which it decomposes.*

*Proof.* “ $\Rightarrow$ ” See discussion above.

“ $\Leftarrow$ ”

Let  $f$  be an arc-flow. Suppose  $i_0$  is a deficit node. Then there exists an incident arc  $(i_0, i_1)$  which carries a positive flow. If  $i_1$  is an excess node, we have found a path from deficit to excess node. Otherwise, the flow balance constraint at  $i_1$  implies that there exists an arc  $(i_1, i_2)$  with positive flow. Repeating this procedure, we finally must arrive at an excess node or revisit a node. This means that we either have constructed a directed path  $P$  from deficit node to excess node or a directed cycle  $C$ , both involving only arcs with strictly positive flow.

In the first case, let  $P = i_0, \dots, i_k$  be the directed path from deficit node  $i_0$  to excess node  $i_k$ . We set  $\beta(P) = \min\{-e_{i_0}, e_{i_k}, \min\{f(i, j) \mid (i, j) \in P\}\}$  and  $f(i, j) = f(i, j) - \beta(P)$ ,  $(i, j) \in P$ . In the second case, set  $\beta(C) = \min\{f(i, j) \mid (i, j) \in C\}$  and  $f(i, j) = f(i, j) - \beta(C)$ ,  $(i, j) \in C$ . Repeat this procedure until all node imbalances are zero.

Now find an arc with positive flow and construct a cycle  $C$  by following only positive arcs from there. Set  $\beta(C) = \min\{f(i, j) \mid (i, j) \in C\}$  and  $f(i, j) = f(i, j) - \beta(C)$ ,  $(i, j) \in C$ . Repeat this process until there are no positive flow-arcs left.

Each time a path or a cycle is identified, the excess/deficit of some node is set to zero or some arc is set to zero. This implies that we decompose into at most  $n + m$  paths and cycles. Since cycle detection sets an arc to zero we have at most  $m$  cycles.  $\square$

An arc flow  $f$  with  $g(i) = 0$  for each  $i \in V$  is called a *circulation*.

**Corollary 8.6.** *A circulation can be decomposed into at most  $m$  cycle-flows.*

Let  $D = (V, A)$  be a network with capacities  $u(i, j)$ ,  $(i, j) \in A$  and costs  $c(i, j)$ ,  $(i, j) \in A$  and let  $f$  be a feasible flow of the network. The *residual network*  $D(f)$  is defined as follows.

- We replace each arc  $(i, j) \in A$  with two arcs  $(i, j)$  and  $(j, i)$ .
- The arc  $(i, j)$  has cost  $c(i, j)$  and *residual capacity*  $r(i, j) = u(i, j) - f(i, j)$ .
- The arc  $(j, i)$  has cost  $-c(i, j)$  and residual capacity  $r(j, i) = f(i, j)$ .
- Delete all arcs which do not have strictly positive residual capacity.

A directed cycle (or path) in  $D(f)$  is called an *augmenting cycle (or path)* of  $f$ .

**Lemma 8.4.** *Suppose that  $f$  and  $f^\circ$  are feasible flows, then  $f - f^\circ$  is a circulation in  $D(f^\circ)$ . Here  $f - f^\circ$  is the flow*

$$(f - f^\circ)(e) = \begin{cases} \max\{0, f(e) - f^\circ(e)\}, & \text{if } e \in A(D) \\ \max\{0, f^\circ(e) - f(e)\}, & \text{if } e^{-1} \in A(D) \\ 0, & \text{otherwise.} \end{cases}$$

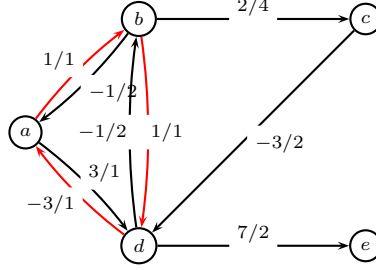


Fig. 8.6: The residual network of the flow in Figure 8.5 and a negative cycle marked by the red edges.

*Proof.* It is very easy to see that the flow  $f - f^\circ$  satisfies the capacity constraints. One also has for each  $v \in V$

$$\sum_{e \in \delta^{out}(v)} (f(e) - f^\circ(e)) - \sum_{e \in \delta^{in}(v)} (f(e) - f^\circ(e)) = 0.$$

If a term  $(f(e) - f^\circ(e))$  is negative, it is replaced by its absolute value and charged as flow on the arc  $e^{-1}$  in  $D(f^\circ)$  which leaves its contribution to the sum above invariant.  $\square$

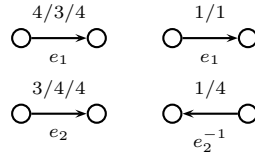


Fig. 8.7: Two arcs  $e_1, e_2 \in A$  labeled with  $f(e)/f^\circ(e)/u(e)$  and the corresponding flow on these arcs (or their reverse) in  $D(f^\circ)$ . Arcs in  $D(f^\circ)$  are labeled with flow and capacity values respectively.



**Theorem 8.12 (Augmenting Cycle Theorem).** *Let  $f$  and  $f^\circ$  be any two feasible flows of a network flow problem. Then  $f$  equals  $f^\circ$  plus the flow of at most  $m$  directed cycles in  $D(f^\circ)$ . Furthermore the cost of  $f$  equals the cost of  $f^\circ$  plus the cost of flow on these augmenting cycles.*

*Proof.* This can be seen by applying flow decomposition on the flow  $f - f^\circ$  in  $D(f^\circ)$ .  $\square$

**Theorem 8.13 (Negative Cycle Optimality Conditions).** *A feasible flow  $f^*$  is an optimal solution of the MCNFP, if and only if it satisfies the negative cycle optimality conditions: the residual network  $D(f^*)$  contains no directed cycle of negative cost.*

*Proof.* “ $\Rightarrow$ ” Suppose that  $f$  is a feasible flow and that  $D(f)$  contains a negative directed cycle. Then  $f$  cannot be optimal, since we can augment positive flow along the corresponding cycle in the network. Therefore, if  $f^*$  is an optimal flow, then  $D(f^*)$  cannot contain a negative directed cycle.

“ $\Leftarrow$ ” Suppose now that  $f^*$  is a feasible flow and suppose that  $D(f^*)$  does not contain a negative cycle. Let  $f^\circ$  be an optimal flow with  $f^\circ \neq f^*$ . The vector  $f^\circ - f^*$  is a circulation in  $D(f^\circ)$  with non-positive cost  $c^T(f^\circ - f^*) \leq 0$ . It follows from Theorem 8.12 that the cost of  $f^\circ$  equals the cost of  $f^*$  plus the cost of directed cycles in the residual network  $D(f^*)$ . The cost of these cycles is nonnegative, and therefore  $c(f^\circ) \geq c(f^*)$  which implies that  $f^*$  is optimal.  $\square$

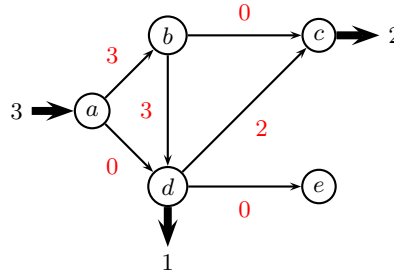


Fig. 8.8: The result of augmenting a flow of one along the negative cycle in Figure 8.6. This flow has cost 12 but is not optimal, since the residual network still contains a negative cycle.

Algorithm 8.1 (Cycle Canceling Algorithm).

1. establish a feasible flow  $f$  in the network
2. WHILE  $D(f)$  contains a negative cycle
  - a. detect a negative cycle  $C$  in  $D(f)$
  - b. let  $\delta = \min\{r(i, j) \mid (i, j) \in C\}$
  - c. augment  $\delta$  units of flow along the cycle  $C$
  - d. update  $D(f)$
3. RETURN  $f$

**Theorem 8.14.** *The cycle canceling algorithm terminates after a finite number of steps if the MCNFP has an optimal solution.*

*Proof.* The cycle canceling algorithm reduces the cost in each iteration. We have assumed that the input data is integral. Thus the cost decreases by at least one unit each iteration. Therefore the number of iterations is finite.  $\square$

**Corollary 8.7.** *If the capacities are integral and if the MCNFP has a optimal flow, then it has an optimal flow with integer values only.*

Let  $\pi : V \rightarrow \mathbb{R}$  be a function (*node potential*). The *reduced cost* of an arc  $(i, j)$  w.r.t.  $\pi$  is  $c_\pi((i, j)) = c((i, j)) + \pi(i) - \pi(j)$ . The potential  $\pi$  is called *feasible* if  $c_\pi((i, j)) \geq 0$  for all arcs  $(i, j) \in A$ .

**Lemma 8.5.** *Let  $D = (V, A)$  be a digraph with arc weights  $c : A \rightarrow \mathbb{R}$ . Then  $D$  does not have a negative cycle if and only if there exists a feasible node potential  $\pi$  of  $D$ .*

*Proof.* Consider a directed path  $P = i_0, i_1, \dots, i_k$ . The cost of this path is

$$c(P) = \sum_{j=1}^k c((i_{j-1}, i_j)).$$

The reduced cost of this path is equal to

$$c_\pi(P) = \sum_{j=1}^k c_\pi((i_{j-1}, i_j)) + \pi(i_0) - \pi(i_k).$$

If  $P$  is a cycle, then  $i_0$  and  $i_k$  are equal, which means that its cost and reduced cost coincide. Thus, if there exists a feasible node potential, then there does not exist a negative cycle.

On the other hand, suppose that  $D$  does not contain a negative cycle w.r.t.  $c$ . Add a vertex  $s$  to  $D$  and the arcs  $(s, i)$  for all  $i \in V$ . The weights (costs) of all these new arcs is 0. Notice that in this way, no new cycles are created, thus still there does not exist a negative cycle. This means we can compute the shortest paths from  $s$  to all other nodes  $i \in V$ . Let  $\pi$  be the function which assigns these shortest paths lengths. Clearly  $c_\pi((i, j)) = \pi(i) - \pi(j) + c((i, j)) \geq 0$ , since the shortest-path length to  $j$  is at most the shortest-path length to  $i$  plus  $c((i, j))$ .  $\square$

This means that we have again a nice way to prove that a flow is optimal. Simply equip the residual network with a feasible node potential.

**Corollary 8.8 (Reduced Cost Optimality Condition).** *A feasible flow  $f^*$  is optimal if and only if there exists a node potential  $\pi$  such that the reduced costs  $c_\pi(i, j)$  of each arch  $(i, j)$  of  $D(f)$  are nonnegative.*

The cycle canceling algorithm is only pseudopolynomial. If we could always chose a minimum cycle (cycle with best improvement) as an augmenting cycle, we would have a polynomial number of iterations. Finding minimum cycles is *NP*-hard. Instead we augment along *minimum mean cycles*. One can find minimum mean cycles in polynomial time.

The *mean cost* of a cycle  $C \in \mathcal{C}$  is the cost of  $C$  divided by the number of arcs in  $C$ :

$$\sum_{(i,j) \in C} c(i, j) / |C|.$$

Algorithm 8.2 (Minimum Mean Cycle Canceling, MMCC).

1. establish a feasible flow  $f$  in the network
2. WHILE  $D(f)$  contains a negative cycle
  - a. detect a minimum mean cycle  $C$  in  $D(f)$
  - b.  $\delta = \min\{r(i, j) \mid (i, j) \in C\}$
  - c. augment  $\delta$  units of flow along the cycle  $C$
  - d. update  $D(f)$
3. RETURN  $f$

We now analyze the MMCC-algorithm. Let  $\mu(f)$  denote the minimum mean-weight of a cycle in  $D(f)$ .

**Lemma 8.6 (See Korte & Vygen [9]).** *Let  $f_1, f_2, \dots$  be a sequence of feasible flows such that  $f_{i+1}$  results from  $f_i$  by augmenting flow along  $C_i$ , where  $C_i$  is a minimum mean cycle of  $D(f_i)$ , then*

1.  $\mu(f_k) \leq \mu(f_{k+1})$  for all  $k$ .
2.  $\mu(f_k) \leq \frac{n}{n-1} \mu(f_l)$ , where  $k < l$  and  $C_k \cup C_l$  contains a pair of reversed arcs.

*Proof.* 1): Suppose  $f_k$  and  $f_{k+1}$  are two subsequent flows in this sequence. Consider the multi-graph  $H$  which results from  $C_k$  and  $C_{k+1}$  by deleting pairs of opposing arcs. The arcs of  $H$  are a subset of the arcs of  $D(f_k)$ , since an arc of  $C_{k+1}$  which is not in  $D(f_k)$  must be a reverse arc of  $C_k$ .

Each node in  $H$  has even degree. Thus  $H$  can be decomposed into cycles, each of mean weight at least  $\mu(f_k)$ . Thus we have  $c(A(H)) \geq \mu(f_k)|A(H)|$ .

Since the total weight of each reverse pair of arcs is zero we have

$$c(A(H)) = c(C_k) + c(C_{k+1}) = \mu(f_k)|C_k| + \mu(f_{k+1})|C_{k+1}|.$$

Since  $|A(H)| \leq |C_k| + |C_{k+1}|$  we conclude

$$\begin{aligned} \mu(f_k)(|C_k| + |C_{k+1}|) &\leq \mu(f_k)|A(H)| \\ &\leq c(A(H)) \\ &= \mu(f_k)|C_k| + \mu(f_{k+1})|C_{k+1}|. \end{aligned}$$

Thus  $\mu(f_k) \leq \mu(f_{k+1})$ .

2): By the first part of the theorem, it is enough to prove the statement for  $k, l$  such that  $C_i \cup C_l$  does not contain a pair of reverse arcs for each  $i, k < i < l$ .

Again, consider the graph  $H$  resulting from  $C_k$  and  $C_l$  by deleting pairs of opposing arcs.  $H$  is a subgraph of  $D(f_k)$ , since any arc of  $C_l$  which does not belong to  $D(f_k)$  must be a reverse arc of  $C_k, C_{k+1}, \dots, C_{l-1}$ . But only  $C_k$  contains a reverse arc of  $C_l$ . So as above we have

$$c(A(H)) = c(C_k) + c(C_l) = \mu(f_k)|C_k| + \mu(f_l)|C_l|.$$

Since  $|A(H)| \leq |C_k| + |C_l| - 2$  we have  $|A(H)| \leq \frac{n-1}{n}(|C_k| + |C_l|)$ . Thus we get

$$\begin{aligned} \mu(f_k) \frac{n-1}{n}(|C_k| + |C_l|) &\leq \mu(f_k)|A(H)| \\ &\leq c(A(H)) \\ &= \mu(f_k)|C_k| + \mu(f_l)|C_l| \\ &\leq \mu(f_l)(|C_k| + |C_l|) \end{aligned}$$

This implies that  $\mu(f_k) \leq \frac{n}{n-1}\mu(f_l)$ . □

**Corollary 8.9.** *During the execution of the MMCC-algorithm,  $|\mu(f)|$  decreases by a factor of  $1/2$  every  $n \cdot m$  iterations.*

*Proof.* Let  $C_1, C_2, \dots$  be the sequence of augmenting cycles. Every  $m^{th}$  iteration, there must be an arc of the cycle, which is reverse to one of the succeeding  $m-1$  cycles, because every iteration, one arc of the residual network will be deleted. Thus after  $nm$  iterations, the absolute value of  $\mu$  has dropped by  $\left(\frac{n-1}{n}\right)^n \leq e^{-1} \leq 1/2$ . □

**Corollary 8.10.** *If all data are integral, then the MMCC-algorithm runs in polynomial time.*

*Proof.* • A lower bound on  $\mu$  is the smallest cost  $c_{min}$

- $|\mu|$  drops by  $1/2$  every  $mn$  iterations.
- After  $mn \log n |c_{min}|$  iterations, absolute value of minimum mean weight cycle drops below  $1/n$ , thus is zero.

- We need to prove that a minimum mean cycle can be found in polynomial time

□

This is a so-called *weakly polynomial* bound, since the binary encoding length of the numbers in the input (here the costs) influences the running time. We now prove that the MMCC-algorithm is *strongly polynomial*.

**Theorem 8.15 (See Korte & Vygen [9]).** *The MMCC-algorithm requires  $O(m^2 n \log n)$  iterations (mean weight cycle cancellations).*

*Proof.* One shows that every  $m n (\lceil \log n \rceil + 1)$  iterations, at least one arc is *fixed*, which means that the flow through this arc does not change anymore.

Let  $f_1$  be some flow at some iteration and let  $f_2$  be the flow  $m n (\lceil \log n \rceil + 1)$  iterations later. It follows from Corollary 8.9 that

$$\mu(f_1) \leq 2 n \mu(f_2) \quad (8.9)$$

holds.

Define the costs  $c'(e) = c(e) - \mu(f_2)$  for the residual network  $D(f_2)$ . There exists no negative cycle in  $D(f_2)$  w.r.t. this cost  $c'$ . (A cycle  $C$  has weight  $c'(C) = \sum_{e \in C} c(e) - |C| \mu(f_2)$  and thus  $c'(C)/|C| = \sum_{e \in C} c(e)/|C| - \mu(f_2) \geq 0$ ). By Lemma 8.5 there exists a feasible node potential  $\pi$  for these weights. One has  $0 \leq c'_\pi(e) = c_\pi(e) - \mu(f_2)$  and thus

$$c_\pi(e) \geq \mu(f_2), \text{ for all } e \in A(D(f_2)). \quad (8.10)$$

Let  $C$  be a minimum mean cycle of  $D(f_1)$ . One has

$$c_\pi(C) = c(C) = \mu(f_1) |C| \leq 2 n \mu(f_2) |C|. \quad (8.11)$$

It follows that there exists an arc  $e_0$  of  $C$  such that

$$c_\pi(e_0) \leq 2 n \mu(f_2) \quad (8.12)$$

holds. The inequalities (8.10) imply that  $e_0 \notin A(D(f_2))$

We now make the following claim:

Let  $f'$  be a feasible flow such that  $e_0 \in D(f')$ , then  $\mu(f') \leq \mu(f_2)$ .

If we have shown this claim, then it follows from Lemma 8.6 that  $e_0$  cannot be anymore in the residual network of a flow after  $f_2$ . Thus the flow along the arc  $e_0$  (or  $e_0^{-1}$ ) is fixed.

Let  $f'$  be a flow such that  $e_0 \in A(D(f'))$ . Recall that  $f' - f_2$  is a circulation in  $D(f_2)$  where  $e_0 \notin D(f_2)$ ,  $e_0^{-1} \in D(f_2)$  and this circulation sends flow over  $e_0^{-1}$ . This circulation can be decomposed into cycles and one of these cycles  $C$  contains  $e_0^{-1}$ . One has  $c_\pi(e_0^{-1}) = -c_\pi(e_0) \geq -2 n \mu(f_2)$  (eq. (8.12)). Using (8.10) one obtains

$$c(C) = \sum_{e \in C} c_\pi(e) \quad (8.13)$$

$$\geq -2n\mu(f_2) + (n-1)\mu(f_2) \quad (8.14)$$

$$= -(n+1)\mu(f_2) \quad (8.15)$$

$$> -n\mu(f_2). \quad (8.16)$$

The reverse of  $C$  is an augmenting cycle for  $f'$  with total weight at most  $n\mu(f_2)$  and thus with mean weight at most  $\mu(f_2)$ . Thus  $\mu(f') \leq \mu(f_2)$ .  $\square$

## 8.6 Computing a minimum cost-to-profit ratio cycle

Given a digraph  $D = (V, A)$  with costs  $c : A \rightarrow \mathbb{Z}$  and profit  $p : A \rightarrow \mathbb{N}_{>0}$ , the task is to compute a cycle  $C \in \mathcal{C}$  with minimum ratio

$$\frac{c(C)}{p(C)}. \quad (8.17)$$

Notice that this is the largest number  $\beta \in \mathbb{Q}$  which satisfies

$$\beta \leq \frac{c(C)}{p(C)}, \text{ for all } C \in \mathcal{C}. \quad (8.18)$$

By rewriting this inequality, we understand this to be the largest number  $\beta \in \mathbb{Q}$  such that

$$c(C) - \beta p(C) \geq 0 \text{ for all } C \in \mathcal{C}. \quad (8.19)$$

In other words, given a digraph  $D = (V, A)$  with costs  $c_\beta : A \rightarrow \mathbb{Q}$ , where  $c_\beta(e) = c(e) - \beta p(e)$ , we search the largest number  $\beta \in \mathbb{Q}$  such that the

$$c_\beta \geq 0. \quad (8.20)$$

We need a routine to check whether  $D$  has a negative cycle for a given weight function  $c$ . For this we assume w.l.o.g. that each vertex is reachable from the vertex  $s$ , if necessary by introducing a new vertex  $s$  from which there is an arc with cost and profit 0 to all other nodes. The minimum cost-to-profit ratio cycle w.r.t. this new graph is then the minimum cost-to-profit ratio cycle w.r.t. the original graph, since  $s$  is not a vertex of any cycle.

Recall the following single-source shortest-path algorithm of Bellman-Ford which we now apply with weights  $c_\beta$ :

Let  $n = |V|$  and  $m = |A|$ . We calculate functions  $f_0, f_1, \dots, f_n : V \rightarrow \mathbb{R} \cup \{\infty\}$  successively by the following rule.

- i)  $f_0(s) = 0, f_0(v) = \infty$  for all  $v \neq s$
- ii) For  $k < n$  if  $f_k$  has been found, compute

$$f_{k+1}(v) = \min\{f_k(v), \min_{(u,v) \in A} \{f_k(u) + c_\beta(u, v)\}\}$$

for all  $v \in V$ .

As seen when we talked about the Bellman-Ford algorithm, there exists a negative cycle w.r.t.  $c_\beta$  if and only if  $f_n(v) < f_k(v)$  for some  $v \in V$  and  $1 \leq k < n$ . Thus we can test in  $O(m \cdot n)$  steps whether  $D$  contains a negative cycle w.r.t.  $c_\beta$ .

We now apply the following idea to search for the correct value of  $\beta$ . We keep an interval  $I = [L, U]$  with the invariant that the value  $\beta$  that we are searching lies in this interval  $I$ . As starting values, we can choose  $L = c_{\min}$  and  $U = c_{\max}$ , where  $c_{\min}$  and  $c_{\max}$  are the smallest and largest cost respectively. In one iteration we compute  $M = (L + U)/2$ . We then check whether  $D$ , together with  $c_M$  contains a negative cycle. If yes, we know that  $\beta$  is at least  $M$  and we set  $L \leftarrow M$ . If not, then  $\beta$  is at most  $M$  and we update the upper bound  $U \leftarrow M$ .

When can we stop this procedure? We can stop it, if we can assure that only one valid cost-to-profit ratio cycle lies in  $[L, U]$ . Suppose that  $C_1$  and  $C_2$  have different cost-to-profit ratios. Then

$$|c(C_1)/p(C_1) - c(C_2)/p(C_2)| = \left| \frac{c(C_1)p(C_2) - c(C_2)p(C_1)}{p(C_1)p(C_2)} \right| \quad (8.21)$$

$$\geq 1/(n^2 p_{\max}^2). \quad (8.22)$$

Thus we can stop our process, if  $U - L < 1/(n^2 p_{\max}^2)$ , since we know then that there can be only one cycle  $c \in \mathcal{C}$  with  $c(C)/p(C) \in [L, U]$ .

Suppose that  $[L, U]$  is the final interval. We know then that

$$L \leq c(C)/p(C) \text{ for all } C \in \mathcal{C}$$

and

$$U > c(C)/p(C) \text{ holds for some } C \in \mathcal{C}.$$

Let  $C$  be a minimum weight cycle w.r.t. the arc costs  $c_L$ . Clearly  $U > c(C)/p(C) \geq L$  holds and thus  $C$  is the minimum cost-to-profit cycle we have been looking for.

Let us analyze the number of required iterations. We need to halve the starting interval-length  $2c$ , where  $c$  is the largest absolute value of a cost, until the length is at most  $1/(n^2 p_{\max}^2)$ . We search the minimal  $i \in \mathbb{N}$  such that

$$(1/2)^i c \leq 1/(n^2 p_{\max}^2). \quad (8.23)$$

This shows us that we need  $O(\log(c p_{\max}^2 n^2))$  iterations which is  $O(\log n \log K)$ , where  $K$  is the largest absolute value of a cost or a profit.

**Theorem 8.16 (Lawler [10]).** *Let  $D$  be a digraph with costs  $c : A \rightarrow \mathbb{Z}$  and profit  $p : A \rightarrow \mathbb{N}_{>0}$  and let  $K \in \mathbb{N}$  such that  $|c(e)| + |p(e)| \leq K$  for all*

$e \in \mathbb{N}$ . A minimum cost-to-profit ratio cycle of  $G$  can be computed in time  $O(mn \log n \log K)$ .

But we knew a weakly polynomial algorithm for MCNFP from the exercises. So you surely ask: Can we do better for minimum cost-to-profit cycle computation? The answer is “Yes”!

### 8.6.1 Parametric search

Let us first roughly describe the idea on how to obtain a strongly polynomial algorithm, see [14]. The Bellman-Ford algorithm tells us whether our current  $\beta$  is too large or too small, depending on whether  $D$  with weights  $c_\beta$  contains a negative cycle or not. Recall that the B-F algorithm computes labels  $f_i(v)$  for  $v \in V$  and  $1 \leq i \leq n$ . If these labels are computed with costs  $c_\beta$ , then they are *piecewise linear* functions in  $\beta$  and we denote them by  $f_i(v)[\beta]$ .

Denote the optimal  $\beta$  that we look for by  $\beta^*$  and suppose that we know an interval  $I$  with such that  $\beta^* \in I$  and each function  $f_i(v)[\beta]$  is linear if it is restricted to this domain  $I$ . Then we can determine  $\beta^*$  as follows.

Let  $I = [L, U]$  be the interval and remember that we are searching for the largest value of  $\beta \in I$  such that  $f_n(v)[\beta] = f_{n-1}(v)[\beta]$  holds for each  $v \in V$ . Clearly this holds for  $\beta = L$ . Thus we only need to check whether  $\beta = U$  by computing the values  $f_n(v)[U]$  and  $f_{n-1}(v)[U]$  for each  $v \in V$  and check whether one of these pairs consists of different numbers.

The idea is now to compute such an interval  $I = [L, U]$  in strongly polynomial time.

Consider the function  $f_1(v)[\beta]$ . Clearly one has

$$f_1(v)[\beta] = \begin{cases} c(s, v) - \beta \cdot p(s, v) & \text{if } (s, v) \in A, \\ \infty & \text{otherwise.} \end{cases}$$

This shows that  $f_1(v)[\beta]$  is a linear function in  $\beta$  for each  $v \in V$ .

Now suppose that  $i \geq 1$  and that we have computed an interval  $I = [L, U]$  with  $\beta^* \in I$  and each function  $f_i(v)[\beta]$  is a linear function if  $\beta$  is restricted to  $I$ .

Now consider the function  $f_{i+1}(v)[\beta]$  for a particular  $v \in V$ . Recall the formula

$$f_{i+1}(v)[\beta] = \min\{f_i(v)[\beta], \min_{(u,v) \in A} \{f_i(u)[\beta] + c(u, v) - \beta \cdot p(u, v)\}\}. \quad (8.24)$$

Each of the functions  $f_i(v)[\beta]$  and  $f_i(u)[\beta] + c(u, v) - \beta \cdot p(u, v)$  are linear on  $I$ . The function  $f_i(v)[\beta]$  can be retrieved by computing a shortest path  $P_i(v)$  from  $s$  to  $v$  with arc weights  $c_\beta$  for some  $\beta$  in  $(L, U)$  which uses at most  $i$  arcs. If  $\beta$  is then allowed to vary, the line which is defined by  $f_i(v)[\beta]$  on  $I$



is then the length of this path  $P$  with parameter  $\beta$ . Similarly we can retrieve the functions (lines)  $f_i(u)[\beta] + c(u, v) - \beta \cdot p(u, v)$  for each  $(u, v) \in A$ . With the Bellman-Ford algorithm, this amounts to a running time of  $O(m \cdot n)$ .

We now have  $n$  lines and can now compute the lower envelope of these lines in time  $O(n \log n)$  alternatively we can also compute all intersection points of these lines and sort them w.r.t. increasing  $\beta$ -coordinate. This would amount to  $O(n^2 \log n)$ . Let  $\beta_1, \dots, \beta_k$  be the sorted list of these  $\beta$ -coordinates. Now  $\beta_{trial} := \beta_{\lfloor k/2 \rfloor}$  and check whether  $\beta^* > \beta_{trial}$ . If yes, we can replace  $L$  by  $\beta_{trial}$  and we can delete the numbers  $\beta_1, \dots, \beta_{\lfloor k/2 \rfloor - 1}$ . Otherwise, we replace  $U$  by  $\beta_{trial}$  and delete  $\beta_{\lfloor k/2 \rfloor + 1}, \dots, \beta_k$ . In any case, we halved the number of possible  $\beta$ -coordinates and continue in this way. Such a check requires a negative cycle test in the graph  $D$  with arc weights  $\beta_{trial}$  and costs  $O(m \cdot n)$ . In the end we have two consecutive  $\beta$ -coordinates and have an interval  $[L, U]$  on which  $f_{i+1}(v)[\beta]$  is linear. To find an interval  $I$  such that  $f_{i+1}(v)[\beta]$  is linear on  $I$  and  $\beta^* \in I$  costs thus  $O(m \cdot n \log n)$  steps.

We now continue to tighten *this interval* such that all functions  $f_{i+1}(v)[\beta]$ ,  $v \in V$  are linear on  $[L, U]$ . Thus in step  $i + 1$  this amounts to a running time of

$$O(n \cdot (m \cdot n \log n)).$$

The total running time is thus

$$O(n^3 \cdot m \cdot \log n).$$

**Theorem 8.17.** *Let  $D = (V, A)$  be a directed graph and let  $c : A \rightarrow \mathbb{R}$  and  $p : A \rightarrow \mathbb{R}_{>0}$  be functions. One can compute a cycle  $C$  of  $D$  minimizing  $c(C)/p(C)$  in time  $O(n^3 \cdot m \cdot \log n)$ .*

### 8.6.1.1 Exercises

- 1) Show that there are no two different paths from  $r$  to another node in a directed tree  $T = (V, A)$ .
- 2) Prove Lemma 8.3.
- 3) Why can we assume without loss of generality that a minimum cost network has a path from  $i$  to  $j$  for all  $i \neq j \in V$  which is incapacitated?
- 4) Provide an example of a MCNFP for which the simple cycle-canceling algorithm from above can require an exponential number of cancels, if the cycles are chosen in a disadvantageous way.
- 5) Provide a proof of Theorem 8.8.
- 6) Let  $Q = \langle u_1, \dots, u_k \rangle$  be the queue before an iteration of the **while** loop of the breadth-first-search algorithm. Show that  $D[u_i]$  is monotonously increasing and that  $D[u_1] + 1 \geq D[u_k]$ . Conclude that the sequence of assigned labels (over time) is a monotonously increasing sequence.
- 7) Extend the definitions and properties of flows and networks in a multi-source and multi-sink network. Show that any flow in a multi-source and

multi-sink network corresponds to a flow of identical value in the single-source and single-sink network obtained by unifying the sources and sinks in a single super-source and super-sink respectively.

## Chapter 9

# The ellipsoid method

It is not known whether the simplex algorithm is an algorithm that runs in polynomial time. For many pivoting rules it was even proved to require an exponential number of iterations [7]. It was long open, whether there exists a polynomial time algorithm for linear programming until Khachiyan [6] showed that the ellipsoid method [17, 15] can solve linear programs in polynomial time. The remarkable fact is that the algorithm is polynomial in the binary encoding length of the linear program. In other words, if the input consists of the problem  $\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\}$ , where  $A \in \mathbb{Q}^{m \times n}$  and  $b \in \mathbb{Q}^m$ , then the algorithm runs in polynomial time in  $m + n + s$ , where  $s$  is the largest binary encoding length of a rational number appearing in  $A$  or  $b$ . The question, whether there exists an algorithm which runs in time polynomial in  $m + n$  and performs arithmetic operations on numbers, whose binary encoding length remains polynomial in  $m + n + s$  is one of the most prominent open problems in theoretical computer science and discrete optimization.

Initially, the ellipsoid method can be used to solve the following problem.

Given a matrix  $A \in \mathbb{Z}^{m \times n}$  and a vector  $b \in \mathbb{Z}^m$ , determine a feasible point  $x^*$  in the polyhedron  $P = \{x \in \mathbb{R}^n \mid Ax \leq b\}$  or assert that  $P$  is *not full-dimensional* or  $P$  is unbounded.

After we understand how the ellipsoid method solves this problem in polynomial time, we discuss why linear programming can be solved in polynomial time.

Clearly, we can assume that  $A$  has full column rank. Otherwise, we can find with Gaussian elimination an invertible matrix  $U \in \mathbb{R}^{n \times n}$  with  $A \cdot U = \begin{pmatrix} A' & 0 \end{pmatrix}$  where  $A'$  has full column rank. The system  $A'x \leq b$  is then feasible if and only if  $Ax \leq b$  is feasible.

**Exercise 9.1.** Let  $x'$  be a feasible solution of  $A'x \leq b$  and suppose that  $U$  from above is given. Show how to compute a feasible solution  $\tilde{x}$  of  $Ax \leq b$ . Also vice versa, show how to compute  $x'$ , if  $\tilde{x}$  is given.

The *unit ball* is the set  $B = \{x \in \mathbb{R}^n \mid \|x\| \leq 1\}$  and an *ellipsoid*  $E(A, b)$  is the image of the unit ball under a linear map  $t : \mathbb{R}^n \rightarrow \mathbb{R}^n$  with  $t(x) = Ax + b$ ,

where  $A \in \mathbb{R}^{n \times n}$  is an invertible matrix and  $b \in \mathbb{R}^n$  is a vector. Clearly

$$E(A, b) = \{x \in \mathbb{R}^n \mid \|A^{-1}x - A^{-1}b\| \leq 1\}. \quad (9.1)$$

**Exercise 9.2.** Consider the mapping  $t(x) = \begin{pmatrix} 1 & 3 \\ 2 & 5 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ . Draw the ellipsoid which is defined by  $t$ . What are the axes of the ellipsoid?

The *volume* of the unit ball is denoted by  $V_n$ , where  $V_n \sim \frac{1}{\pi n} \left(\frac{2e\pi}{n}\right)^{n/2}$ . It follows that the volume of the ellipsoid  $E(A, b)$  is equal to  $|\det(A)| \cdot V_n$ . The next lemma is the key to the development of the ellipsoid method.

**Lemma 9.1 (Half-Ball Lemma).** *The half-ball  $H = \{x \in \mathbb{R}^n \mid \|x\| \leq 1, x_1 \geq 0\}$  is contained in the ellipsoid*

$$E = \left\{x \in \mathbb{R}^n \mid \left(\frac{n+1}{n}\right)^2 \left(x_1 - \frac{1}{n+1}\right)^2 + \frac{n^2-1}{n^2} \sum_{i=2}^n x_i^2 \leq 1\right\} \quad (9.2)$$

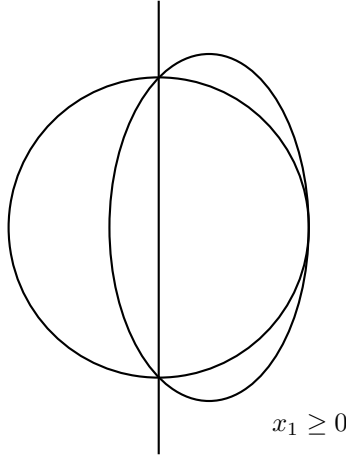


Fig. 9.1: Half-ball lemma.

*Proof.* Let  $x$  be contained in the unit ball, i.e.,  $\|x\| \leq 1$  and suppose further that  $0 \leq x_1$  holds. We need to show that

$$\left(\frac{n+1}{n}\right)^2 \left(x_1 - \frac{1}{n+1}\right)^2 + \frac{n^2-1}{n^2} \sum_{i=2}^n x_i^2 \leq 1 \quad (9.3)$$

holds. Since  $\sum_{i=2}^n x_i^2 \leq 1 - x_1^2$  holds we have

$$\begin{aligned}
\left(\frac{n+1}{n}\right)^2 \left(x_1 - \frac{1}{n+1}\right)^2 + \frac{n^2-1}{n^2} \sum_{i=2}^n x_i^2 \\
\leq \left(\frac{n+1}{n}\right)^2 \left(x_1 - \frac{1}{n+1}\right)^2 + \frac{n^2-1}{n^2} (1 - x_1^2)
\end{aligned} \tag{9.4}$$

This shows that (9.3) holds if  $x$  is contained in the half-ball and  $x_1 = 0$  or  $x_1 = 1$ . Now consider the right-hand-side of (9.4) as a function of  $x_1$ , i.e., consider

$$f(x_1) = \left(\frac{n+1}{n}\right)^2 \left(x_1 - \frac{1}{n+1}\right)^2 + \frac{n^2-1}{n^2} (1 - x_1^2). \tag{9.5}$$

The first derivative is

$$f'(x_1) = 2 \cdot \left(\frac{n+1}{n}\right)^2 \left(x_1 - \frac{1}{n+1}\right) - 2 \cdot \frac{n^2-1}{n^2} x_1. \tag{9.6}$$

We have  $f'(0) < 0$  and since both  $f(0) = 1$  and  $f(1) = 1$  (and since  $f(x_1)$  is a 2-degree polynomial w.r.t.  $x_1$ ), we have  $f(x_1) \leq 1$  for all  $0 \leq x_1 \leq 1$  and the assertion follows.

In terms of a matrix  $A$  and a vector  $b$ , the ellipsoid  $E$  is described as  $E = \{x \in \mathbb{R}^n \mid \|A^{-1}x - A^{-1}b\| \leq 1\}$ , where  $A$  is the diagonal matrix with diagonal entries

$$\frac{n}{n+1}, \sqrt{\frac{n^2}{n^2-1}}, \dots, \sqrt{\frac{n^2}{n^2-1}}$$

and  $b$  is the vector  $b = (1/(n+1), 0, \dots, 0)$ . Our ellipsoid  $E$  is thus the image of the unit sphere under the linear transformation  $t(x) = Ax + b$ . The determinant of  $A$  is thus  $\frac{n}{n+1} \left(\frac{n^2}{n^2-1}\right)^{(n-1)/2}$  which is bounded by

$$e^{-1/(n+1)} e^{(n-1)/(2 \cdot (n^2-1))} = e^{-\frac{1}{2(n+1)}}. \tag{9.7}$$

We can conclude the following theorem.

**Theorem 9.1.** *The half-ball  $\{x \in \mathbb{R}^n \mid x_1 \geq 0, \|x\| \leq 1\}$  is contained in an ellipsoid  $E$ , whose volume is bounded by  $e^{-\frac{1}{2(n+1)}} \cdot V_n$ .*

Recall the following notion from linear algebra. A symmetric matrix  $A \in \mathbb{R}^{n \times n}$  is called positive definite if all its eigenvalues are positive. Recall the following theorem.

**Theorem 9.2.** *Let  $A \in \mathbb{R}^{n \times n}$  be a symmetric matrix. The following are equivalent.*

- i)  $A$  is positive definite.

- ii)  $A = L^T L$ , where  $L \in \mathbb{R}^{n \times n}$  is a uniquely determined upper triangular matrix.
- iii)  $x^T A x > 0$  for each  $x \in \mathbb{R}^n \setminus \{0\}$ .
- iv)  $A = Q^T \text{diag}(\lambda_1, \dots, \lambda_n) Q$ , where  $Q \in \mathbb{R}^{n \times n}$  is an orthogonal matrix and  $\lambda_i \in \mathbb{R}_{>0}$  for  $i = 1, \dots, n$ .

It is now convenient to switch to a different representation of an ellipsoid. An ellipsoid  $\mathcal{E}(A, a)$  is the set  $\mathcal{E}(A, a) = \{x \in \mathbb{R}^n \mid (x - a)^T A^{-1} (x - a) \leq 1\}$ , where  $A \in \mathbb{R}^{n \times n}$  is a symmetric positive definite matrix and  $a \in \mathbb{R}^n$  is a vector, called the *center* of the ellipsoid. Consider the half-ellipsoid  $\mathcal{E}(A, a) \cap \{c^T x \leq c^T a\}$ .

Our goal is a similar lemma as the half-ball-lemma for ellipsoids. Geometrically it is clear that each half-ellipsoid  $\mathcal{E}(A, a) \cap \{c^T x \leq c^T a\}$  must be contained in another ellipsoid  $\mathcal{E}(A', b')$  with  $\text{vol}(\mathcal{E}(A', b')) / \text{vol}(\mathcal{E}(A, a)) \leq e^{-1/(2n)}$ . More precisely this follows from the fact that the half-ellipsoid is the image of the half-ball under a linear transformation. Therefore the image of the ellipsoid  $E$  under the same transformation contains the half-ellipsoid. Also, the volume-ratio of the two ellipsoids is invariant under a linear transformation.

We now record the formula for the ellipsoid  $\mathcal{E}'(A', a')$ . It is defined by

$$a' = a - \frac{1}{n+1} b \tag{9.8}$$

$$A' = \frac{n^2}{n^2 - 1} \left( A - \frac{2}{n+1} b b^T \right), \tag{9.9}$$

where  $b$  is the vector  $b = A c / \sqrt{c^T A c}$ . The proof of the correctness of this formula can be found in [5].

**Lemma 9.2 (Half-Ellipsoid-Theorem).** *The half-ellipsoid  $\mathcal{E}(A, b) \cap \{c^T x \leq c^T a\}$  is contained in the ellipsoid  $\mathcal{E}'(A', a')$  and one has  $\text{vol}(\mathcal{E}') / \text{vol}(\mathcal{E}) \leq e^{-1/(2n)}$ .*

Before talking about the method, we give a useful definition.

**Definition 9.1.** A polyhedron  $P$  is full-dimensional if it has positive volume.

## 9.1 The method

Suppose we know the following things of our polyhedron  $P$ .

- I) We have a number  $L$  such that  $\text{vol}(P) \geq L$  if  $P$  is full-dimensional.
- II) We have an ellipsoid  $\mathcal{E}_{init}$  which contains  $P$  if  $P$  is bounded.

The ellipsoid method is now easily described.

Algorithm 9.1 (Ellipsoid method exact version).

- a) (Initialize): Set  $\mathcal{E}(A, a) := \mathcal{E}_{init}$
- b) If  $a \in P$ , then assert  $P \neq \emptyset$  and stop
- c) If  $\text{vol}(\mathcal{E}) < L$ , then assert that  $P$  is unbounded or  $P$  is not full-dimensional
- d) Otherwise, compute an inequality  $c^T x \leq \beta$  which is valid for  $P$  and satisfies  $c^T a > \beta$  and replace  $\mathcal{E}(A, a)$  by  $\mathcal{E}(A', a')$  computed with formula (9.8) and goto step b).

**Theorem 9.3.** *The ellipsoid method computes a point in the polyhedron  $P$  or asserts that  $P$  is unbounded or not full-dimensional. The number of iterations is bounded by  $2 \cdot n \ln(\text{vol}(\mathcal{E}_{init})/L)$ .*

*Proof.* Unless  $P$  is unbounded, we start with an ellipsoid which contains  $P$ . This then holds for all the subsequently computed ellipsoids. After  $i$  iterations one has

$$\text{vol}(\mathcal{E})/\text{vol}(\mathcal{E}_{init}) \leq e^{-\frac{i}{2n}}. \quad (9.10)$$

Since we stop when  $\text{vol}(\mathcal{E}) < L$ , we stop at least after  $2 \cdot n \ln(\text{vol}(\mathcal{E}_{init})/L)$  iterations. This shows the claim.

## 9.2 Deciding feasibility

Suppose that we have to decide the feasibility of  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  using the ellipsoid method described above. Here  $A \in \mathbb{Z}^{m \times n}$  and  $b \in \mathbb{Z}^m$ . In the following  $B$  is an upper bound on the absolute values of the components of  $A$  and  $b$ .

To apply the ellipsoid method we have to take care of the following items:

- i) We have to find a starting ball that contains a subset of the feasible region.
- ii) We have to deal with the fact that  $P$  might not be full dimensional even if it is feasible.
- iii) We have to bound the volume  $P$  from below, if  $P$  is full dimensional.

We first deal with the issue i). Suppose that there exists a feasible point  $x^*$  and consider the index sets  $I = \{i : x_i^* \geq 0\}$  and  $J = \{j : x_j^* \leq 0\}$ . The polyhedron

$$P' = \{x \in \mathbb{R}^n : Ax \leq b, x_i \geq 0, i \in I, x_j \leq 0, j \in J\}$$

is contained in  $P$  and it has vertices. Clearly  $P'$  is described by a system of inequalities  $Cx \leq d$  where each component of  $C$  and  $d$  is bounded by  $B$  in absolute value.

A vertex  $v^*$  of  $P'$  is of the form  $v^* = C_B^{-1}b_B$  for a basis  $B$  of  $C$ . Since  $C_B^{-1} = \widetilde{C}_B / \det(C_B)$  where  $\widetilde{C}_B \in \mathbb{Z}^{n \times n}$  is the adjoint of  $C_B$ . Each component

of  $\widetilde{C}_B$  is, by the Hadamard bound, bounded by  $B^n \cdot n^{n/2}$ . It follows that  $\|v^*\|_\infty$  is bounded by  $n^{n/2+1} \cdot B^{n+1}$ . We can thus conclude that  $P$  is infeasible if and only if

$$P \cap \{x \in \mathbb{R}^n : \|x\|_\infty \leq n^{n/2+1} \cdot B^{n+1}\}$$

is infeasible. For this polyhedron, we have a starting ellipsoid which is the ball of radius  $n^{n/2+1} \cdot B^{n+1}$  around 0.

Next, we deal with the issue ii).

**Exercise 9.3.** Let  $P = \{x \in \mathbb{R}^n \mid Ax \leq b\}$  be a polyhedron and  $\varepsilon > 0$  be a real number. Show that  $P_\varepsilon = \{x \in \mathbb{R}^n \mid Ax \leq b + \varepsilon \cdot \mathbf{1}\}$  is full-dimensional if  $P \neq \emptyset$ .

The above exercise raises the following question. Is there an  $\varepsilon > 0$  such that  $P_\varepsilon = \emptyset$  if and only if  $P = \emptyset$  and furthermore is the binary encoding length of this  $\varepsilon$  polynomial in the binary encoding length of  $A$  and  $b$ ?

We recall the Farkas' lemma (Theorem 1.2).

**Theorem 9.4.** *The system  $Ax \leq b$  does not have a solution if and only if there exists a nonnegative vector  $\lambda \in \mathbb{R}_{\geq 0}^m$  such that  $\lambda^T A = 0$  and  $\lambda^T b = -1$ .*

If  $Ax \leq b$  is infeasible, then the polyhedron

$$\{\lambda \in \mathbb{R}^m : \lambda^T A = 0, \lambda^T b = -1, \lambda \geq 0\}$$

has vertices. As we have argued before, there exists a vertex  $\lambda^*$  of  $P$  such that  $\|\lambda^*\| \leq n^{n/2+1} \cdot B^{n+1}$ .

Now let  $\varepsilon = (2 \cdot m \cdot n^{n/2+1} \cdot B^{n+1})^{-1}$ . Then  $|\lambda^{*T} \mathbf{1} \cdot \varepsilon| < 1$  and thus

$$\lambda^T(b + \varepsilon \cdot \mathbf{1}) < 0. \quad (9.11)$$

Consequently the system  $Ax \leq b + \varepsilon \cdot \mathbf{1}$  is infeasible if and only if  $Ax \leq b$  is infeasible. Furthermore, if  $Ax \leq b + \varepsilon \cdot \mathbf{1}$  is feasible, then the polyhedron described by this set of inequalities is full-dimensional. Notice again that the encoding length of  $\varepsilon$  is polynomial in the encoding length of  $Ax \leq b$ .

Next we deal with issue iii. We ask ourselves the following question. If  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  is full-dimensional and bounded, what is a lower bound on  $\text{vol}(P)$ ?

If  $P$  is full-dimensional, then  $P$  contains  $n+1$  affinely independent vertices  $v_1, \dots, v_{n+1}$ , see Exercise 1. Each of these vertices is of the form  $v_i = v'_i/q_i$ , where  $v'_i \in \mathbb{Z}^n$  and  $q_i$  is a  $n \times n$ -sub-determinant of  $A$ . Consequently

$$\begin{aligned} \text{vol}(P) &\geq (1/\prod_{i=1}^{n+1} q_i) \cdot \text{vol}(\text{conv}(v'_1, \dots, v'_{n+1})) \\ &\geq (B^n \cdot n^{n/2})^{-n} \cdot 1/n!, \end{aligned}$$

where the last inequality follows from the Hadamard bound and Exercise 2 and Exercise 3. A crude bound on  $\text{vol}(P)$  is thus  $\text{vol}(P) \geq 1/(B^{n^2} \cdot n^{2 \cdot n^2})$ .



We can conclude with the following theorem.

**Theorem 9.5.** *Let  $Ax \leq b$  be a system of linear inequalities with  $A \in \mathbb{Z}^{m \times n}$  and  $b \in \mathbb{Z}^m$ . Let  $B$  be an upper bound on the absolute value of the entries of  $A$  and  $b$ . The ellipsoid method can decide in time polynomial in  $n, m$  and  $\log B$  whether  $Ax \leq b$  is feasible.*

The next theorem follows from binary search. Its proof is an exercise.

**Theorem 9.6.** *Let  $\max\{c^T x : x \in \mathbb{R}^n, Ax \leq b\}$  be a linear program with  $A \in \mathbb{Z}^{m \times n}$ ,  $b \in \mathbb{Z}^m$ ,  $c \in \mathbb{Z}^n$ . Let  $B$  be an upper bound on the absolute value of the entries of  $A$  and  $b$ . The ellipsoid method finds an optimal solution of the linear program if one exists in time polynomial in  $n, m$  and  $\log B$ . Let  $Ax \leq b$  be a system of linear inequalities with  $A \in \mathbb{Z}^{m \times n}$  and  $b \in \mathbb{Z}^m$ . Let  $B$  be an upper bound on the absolute value of the entries of  $A$  and  $b$ . The ellipsoid can decide in time polynomial in  $n, m$  and  $\log B$  whether  $Ax \leq b$  is feasible.*

## Exercises

1. Show the following. If  $P \subseteq \mathbb{R}^n$  is a bounded and full-dimensional polyhedron, then there exist vertices  $v_1, \dots, v_{n+1}$  of  $P$  that are affinely independent, i.e.,  $v_2 - v_1, v_3 - v_1, \dots, v_{n+1} - v_1$  are linearly independent. *Hint: If  $a^T x = \beta$  is some hyperplane, where  $a \in \mathbb{R}^n \setminus \{0\}$ , then there exists a vertex of  $P$  that is not contained in that hyperplane.*
2. Show that  $\text{vol}(\text{conv}(0, e_1, \dots, e_n)) = 1/n!$
3. Let  $a_1, \dots, a_n \in \mathbb{Z}^n$  be linearly independent. Show that

$$\text{vol}(\text{conv}(0, a_1, \dots, a_n)) = |\det(a_1, \dots, a_n)|/n!.$$

## 9.3 The separation problem

At this point we can already notice a very important fact. Inspect step d of the algorithm. What is required here? An inequality which is valid for  $P$  but not for the center  $a$  of  $\mathcal{E}(A, a)$ . Such an inequality is readily at hand if we have the complete inequality description of  $P$  in terms of a system  $Cx \leq d$ . Just pick an inequality which is violated by  $a$ . Sometimes however, it is not possible to describe the polyhedron of a combinatorial optimization problem with an inequality system efficiently, simply because the number of inequalities is too large. An example of such a polyhedron is the matching polytope, see Theorem 7.5.

The great power of the ellipsoid method lies in the fact that we do not have to *write down* the polyhedron entirely. We only have to solve the so-called separation problem for the polyhedron, which is defined as follows.

**Separation Problem**

Given a point  $a \in \mathbb{R}^n$  determine, whether  $a \in P$  and if not, compute an inequality  $c^T x \leq \beta$  which is valid for  $P$  with  $c^T a > \beta$ .

**Exercise 9.4.** We are given an undirected graph  $G = (V, E)$ . A *spanning tree*  $T$  is a subset  $T \subseteq E$  of the edges such that  $T$  does not contain a cycle and  $T$  connects all the vertices  $V$ . Consider the following *spanning tree polytope*  $P_{span}$

$$\sum_{e \in E} x(e) = n - 1 \quad (9.12)$$

$$\sum_{e \in \delta(U)} x(e) \geq 1 \quad \forall \emptyset \subset U \subset V \quad (9.13)$$

$$x(e) \leq 1 \quad \forall e \in E \quad (9.14)$$

$$x(e) \geq 0 \quad \forall e \in E. \quad (9.15)$$

Let  $x$  be an integral solution of  $P_{span}$  and define  $T = \{e \in E \mid x(e) = 1\}$ . The inequality (9.12) ensures that exactly  $n-1$  edges are picked. The inequalities (9.13) ensure that  $T$  connects the vertices of  $G$ . Thus  $T$  must be a spanning tree. Clearly, there are exponentially many inequalities of type (9.13). Nevertheless, a fractional solution of this polytope can be computed using the ellipsoid method.

Show that the separation problem for  $P_{span}$  can be solved in polynomial time.

*Hint:* To verify whether a vector  $x \in \mathbb{R}_{\geq 0}^{|E|}$  fulfills inequalities of type (9.13), it is a good idea to recall the MinCut or MaxFlow problem.

Via binary search even an optimal solution can be computed in polynomial time (in the input length) if we introduce edge costs (you don't have to show that). In the next semester you will see that any optimal basis solution is integral and hence defines an optimal spanning tree w.r.t. the edge costs.

**Exercise 9.5.** Consider the triangle defined by

$$\begin{aligned} -x_1 - x_2 &\leq -2 \\ 3x_1 &\leq 4 \\ -2x_1 + 2x_2 &\leq 3. \end{aligned}$$

Draw the triangle and simulate the ellipsoid method with starting ellipsoid being the ball of radius 6 around 0. Draw each of the computed ellipsoids

with your favorite program (pstrics, maple,...). How many iterations does the ellipsoid method take?

*Ignore the occurring rounding errors!*

## 9.4 The ellipsoid method for optimization

Suppose that you want to solve a linear program

$$\max\{c^T x \mid x \in \mathbb{R}^n, Ax \leq b\} \quad (9.16)$$

and recall that if (9.16) is bounded and feasible, then so is its dual and the two objective values are equal. Thus, we can use the ellipsoid method to find a point  $(x, y)$  with  $c^T x = b^T y$ ,  $Ax \leq b$  and  $A^T y = c, y \geq 0$ .

However, we mentioned that the strength of the ellipsoid method lies in the fact that we do not need to write the system  $Ax \leq b$  down explicitly. The only thing which has to be solvable is the *separation problem*. This is to be exploited in the next exercise.

**Exercise 9.6.** Show how to solve the optimization problem  $\max\{c^T x \mid Ax \leq b\}$  with a polynomial number of calls to an algorithm which solves the separation problem for  $Ax \leq b$ . You may assume that  $A$  has full column rank and the polynomial bound on the number of calls to the algorithm to solve the separation problem can depend on  $n$  and the largest size of a component of  $A, b$  and  $c$ .

## 9.5 Numerical issues

We did not discuss the numerical details on how to implement the ellipsoid method such that it runs in polynomial time. One issue is crucial.

!

We only want to compute with a precision which is polynomial in the input encoding!

In the formula (9.8) the vector  $b$  is defined by taking a square root. The question thus rises on how to round the numbers in the intermediate ellipsoids such that they can be handled on a machine. Also one has to analyze the growth of the numbers in the course of the algorithm. All these issues can be overcome but we do not discuss them in this course. I would like to refer you to the book of Alexander Schrijver [16] for further details. They are not difficult, but a little technical.



## Chapter 10

# Primal-Dual algorithm

In this chapter we talk about matching problems, a very important topic in combinatorial optimization.

### 10.1 Graphs and Matchings

We begin by giving some useful definitions.

**Definition 10.1.** A graph  $G = (V, E)$  (or  $G = (V, A)$ ) is called *bipartite* if we can split  $V$  into two disjoint subsets  $V_1$  and  $V_2$ , such that for every  $e \in E$  (or  $a \in A$ ) one endpoint is in  $V_1$  and the other one is in  $V_2$ .

We recall an important result about bipartite graphs.

**Lemma 10.1.** *A graph is bipartite if and only if it does not contain an odd cycle (that is, a cycle of odd length).*

**Definition 10.2 (Matching).** Let  $G = (V, E)$ , or  $G = (V, A)$ , be a graph. A *matching* in  $G$  is a subset  $M \subseteq E$ , or  $M \subseteq A$ , of pairwise disjoint edges, or arcs (that is, edges, or arcs, which do not share any vertex). A vertex  $v$  that is an endpoint of an edge (or arc) in  $M$  is called a *matched node*, otherwise it is called an *exposed node*. A matching is called *perfect* if there are no exposed nodes.

**Definition 10.3.** Let  $G = (V, E)$  be a graph and  $M \subseteq E$  be a matching in  $G$ . An *alternating path* is a path that alternates between edges in  $M$  and edges in  $E \setminus M$ . An alternating path that starts and ends at an exposed node is called an *augmenting path*.

**Definition 10.4 (Vertex cover).** Let  $G = (V, E)$  be a graph. A *vertex cover* is a set  $C \subseteq V$  such that every  $e \in E$  has at least one endpoint in  $C$ .

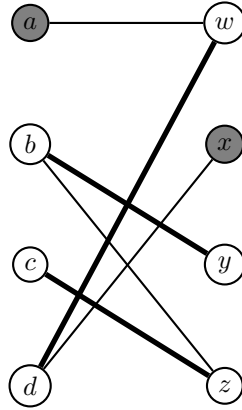


Fig. 10.1: This picture shows an example of an undirected bipartite graph  $G = (V, E)$ , where  $V = \{a, b, c, d, w, x, y, z\}$ ,  $E = \{\{a, w\}, \{b, y\}, \{b, z\}, \{c, z\}, \{d, w\}, \{d, x\}\}$ ,  $V_1 = \{a, b, c, d\}$  and  $V_2 = \{w, x, y, z\}$ . The thicker edges form a matching, vertices  $a$  and  $x$  (marked in grey) are exposed nodes and  $a, w, d, x$  is an alternating and augmenting path.  $C = \{b, w, x, z\}$  is a vertex cover.

## 10.2 Matching problems

Now that we know some basic definitions, we can concentrate on two of the most important problems about matchings:

1. **Maximum cardinality matching problem:** Find a matching  $M$  of maximum size.
2. **Maximum (or minimum) weight matching problem:** Given a graph  $G = (V, E)$  and a weight function  $w : E \rightarrow \mathbb{R}$ , find a matching  $M$  of maximum (or minimum) weight, where the weight of a matching is

$$w(M) = \sum_{e \in M} w(e).$$

From now we will focus on the case of an undirected bipartite graph  $G = (V, E)$ .

### 10.2.1 The maximum cardinality matching problem

Before giving a method to find a maximum cardinality matching in a graph, we want to show how we can prove its optimality. For this purpose we recall theorem 7.3, that gives an upper bound on the size of any matching in a given bipartite graph, and we prove another theorem, which gives us a method to see if a matching is of maximum size or not.

**Theorem 10.1 (König's theorem).** *In any bipartite graph, the number of edges in a maximum cardinality matching equals the number of vertices in a minimum vertex cover.*

**Theorem 10.2.** *A matching  $M$  of a graph  $G = (V, E)$  is of maximum cardinality if and only if there are no augmenting paths with respect to  $M$ .*

*Proof.* ( $\Rightarrow$ ) Suppose there exists an augmenting path in  $G$ , call it  $P$ . Consider the set  $M' = M \triangle P$ . By definition of augmenting path, we have that  $M'$  has exactly one edge more than  $M$ . Since  $M$  is a matching and  $P$  starts and ends in an exposed node, the only edges of  $M$  with an endpoint in  $P$  are the edges of  $P$ . This implies that  $M'$  is a matching.

( $\Leftarrow$ ) Let  $M$  be a maximum cardinality matching and let  $\tilde{M}$  be a strictly smaller matching. Consider  $X = M \triangle \tilde{M}$ .  $X$  is formed by cycles and alternating paths (with respect to  $M$  and  $\tilde{M}$ ). Since  $|M| > |\tilde{M}|$ ,  $X$  contains at least one alternating path  $P$  with more edges from  $M$  than from  $\tilde{M}$ . By definition of  $X$ , we see that  $P$  is an augmenting path.  $\square$

We are now ready to show an algorithm that allows us to find a maximum cardinality matching in any bipartite graph.

Let  $G = (V, E)$  be a bipartite graph with bipartition  $V = A \sqcup B$  and  $M$  a matching. Now turn  $G$  into a directed graph  $D = (V, A)$  by directing matching edges from  $A$  to  $B$  and non-matching edges from  $B$  to  $A$ . We are interested in a method to find augmenting paths or to assert that there aren't any (and thus, that  $M$  is of maximum size). For this purpose we illustrate the following claim and theorem.

*Claim.* Let  $G$  and  $D$  be as above. A path in  $D$  between two exposed nodes that starts in an exposed node in  $B$  (resp. in  $A$ ), ends in an exposed node in  $A$  (resp. in  $B$ ).

**Theorem 10.3.** *Let  $G$  and  $D$  be as above and  $M$  be a matching. There exists an augmenting path in  $G$  if and only if there exists a path from an exposed node in  $B$  to an exposed node in  $A$  in the directed graph  $D$ .*

*Proof.* ( $\Rightarrow$ ) This is a direct consequence of the choice of the direction of the arcs in  $D$  and of the previous claim.

( $\Leftarrow$ ) Trivial.

Let  $G$  and  $D$  be as illustrated before. How can we use the previous theorem for our purpose? Add a vertex  $s$  to our directed graph  $D$  and connect it to all exposed nodes in  $B$  by an arc whose tail is  $s$ . Now we have that finding a directed path, in  $D$ , from an exposed node  $u$  in  $B$  to an exposed node  $v$  in  $A$  is equivalent to finding a directed path from  $s$  to  $v$  passing by  $u$ . Thus, we deduce that there is an augmenting path in  $G$  if and only if there is at least one exposed node in  $A$  reachable from  $s$ . To find if there are such nodes reachable from  $s$  and the corresponding augmenting paths, we can use the *Breadth-First search* algorithm (see chapter 8.2.1) applied to vertex  $s$  in  $D$ . If there is any exposed node  $u$  in  $A$  with finite distance from  $s$ , then we can obtain an augmenting path by taking the shortest path  $s, a_1, a_2, \dots, u$  from  $s$  to  $u$  (which is given by Bread-First search) without  $s$  (i.e., the augmenting path would be  $a_1, a_2, \dots, u$ ).

To resume, we obtain the following algorithm.

Algorithm 10.1.

**Initialise**  $M = \emptyset$

**while** *There exists  $M$ -augmenting path*

*Update  $M$*

**return**  $M$

Finally, we are interested in the running time of our algorithm.

**Theorem 10.4.** *A maximum cardinality matching in a bipartite graph  $G = (V, E)$  can be computed in time  $O(|V| \cdot (|V| + |E|))$ . If we assume that  $G$  does not have any isolated vertex, then we can consider the running time to be  $O(|V| \cdot |E|)$ .*

*Proof.* The *while* loop runs at most  $|V|/2$  times and its execution requires  $O(|V| + |E|)$  ( $= O(|E|)$  if we have the assumption) operations.

### 10.2.2 The maximum weight matching problem

In this section we consider a graph  $G = (V, E)$  and a weight function  $w : E \rightarrow \mathbb{R}$ . First of all, notice that, by changing the sign of the weights of all edges, we have that finding a maximum weight matching or a minimum weight matching are equivalent problems.

We can also prove that the problem of finding a minimum weight matching can always be replaced by the problem of finding a minimum weight *perfect* matching. To see that, it is sufficient to create a copy  $G'$  of our graph  $G$  and add edges of weight 0 connecting all vertices  $v$  in  $G$  to their copy  $v'$  in  $G'$ . This will give us a new graph  $\tilde{G}$ . Notice that  $\tilde{G}$  has at least a perfect matching, thus there is one, call it  $\tilde{M}$ , of minimum weight. If we consider only the edges of  $\tilde{M}$  in  $G$ , we obtain a minimum weight matching in  $G$ .



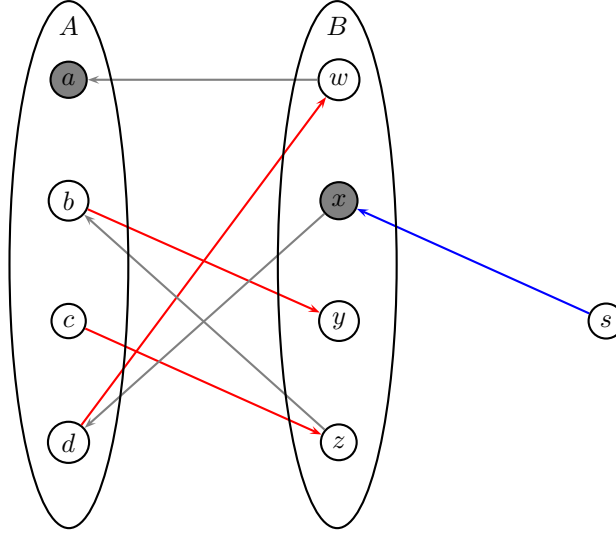


Fig. 10.2: This picture illustrate how transform a bipartite graph to apply the *maximum cardinality matching algorithm*, using the graph of figure 10.1 as an example.

Moreover, we have an important relation between the difficulty of this two problems, given by the following theorem.

**Theorem 10.5.** *If there exists a polynomial time algorithm for the minimum weight perfect matching problem, then there exists a polynomial time algorithm for the minimum weight matching problem.*

*Proof.* Let  $G = (V, E)$  be a graph. Using the process described above, we create the graph  $\tilde{G} = (\tilde{V}, \tilde{E})$ . Notice that  $|\tilde{V}| = 2|V|$  and  $|\tilde{E}| = 2|E| + |V|$ . By hypothesis, we can find a minimum weight perfect matching in  $\tilde{G}$  (and thus a minimum weight matching in  $G$ ) in time  $O((\tilde{V} + \tilde{E})^k) = O((2|V| + 2|E| + |V|)^k) = O((|V| + |E|)^{k'})$ , for some  $k, k' \in \mathbb{N}$  (that is, a polynomial time with respect to the size of  $G$ ).

Since we concentrate on bipartite graphs, we show how to reduce the problem of finding a minimum weight *bipartite* matching to the problem of finding a minimum weight *bipartite* perfect matching. Let  $G = (V, E)$  be a bipartite graph with bipartition  $V = A \sqcup B$ . We can suppose w.l.o.g. that  $|A| \geq |B|$ . If  $|A| > |B|$ , add a set  $C$  of cardinality  $|A| - |B|$  and add edges of weight zero connecting all nodes of  $A$  to all nodes of  $B \cup C$ . Call  $G'$  the bipartite graph obtained by this process. Obviously,  $G'$  admits at least one

bipartite perfect matching, and thus also one of minimum weight, call it  $M'$ . By taking out all edges of weight zero from  $M'$ , we obtain a minimum weight bipartite matching in  $G$ .

Also in this case, we have an important relation between the difficulty of these two matching problems.

**Theorem 10.6.** *If there exists a polynomial time algorithm for the minimum weight bipartite perfect matching problem, then there exists a polynomial time algorithm for the minimum weight bipartite matching problem.*

*Proof.* Let  $G = (V, E)$  be a bipartite graph with bipartition  $V = A \sqcup B$ . Construct  $G' = (V', E')$  as shown above. Notice that  $|V'| = 2|A| \leq 2|V|$  and  $|E'| = |E| + |A| \cdot (|A| - 1) \leq |E| + |V|^2$ . Hence, by hypothesis, finding a minimum weight bipartite matching in  $G$  (by finding a minimum weight bipartite perfect matching in  $G'$ ) can be done in time  $O((2|V| + |E| + |V|^2)^k) = O((|V| + |E|)^{k'})$ , for some  $k, k' \in \mathbb{N}$  (that is, a polynomial time with respect to the size of  $G$ ).

Let  $G = (V, E)$  be a bipartite graph. We have showed that the problem of finding a maximum weight (bipartite) matching can be reduced to the problem of finding a minimum weight (bipartite) perfect matching. Moreover, if there exists a polynomial time algorithm to solve the first problem, then there exists a polynomial time algorithm to solve the second one.

The following theorem is crucial for the introduction of the *Primal-Dual algorithm*.

**Theorem 10.7 (Complementary Slackness).** *Let*

$$\begin{aligned} \max \quad & c^T x \\ \text{subject to} \quad & Ax \leq b \end{aligned}$$

*be a primal LP and let*

$$\begin{aligned} \min \quad & b^T y \\ \text{subject to} \quad & A^T y = c \\ & y \geq 0 \end{aligned}$$

*be his dual. Let  $x^*$  and  $y^*$  be primal and dual feasible solutions respectively. Then they are both optimal if and only if*

$$(b - Ax^*)^T y^* = 0.$$

*Proof.*

$$\begin{aligned} (b - Ax^*)^T y^* &= 0 \\ \Leftrightarrow \\ y_i^* > 0 &\Rightarrow a_i^T x^* = b_i \\ \Leftrightarrow \\ b^T y^* &= (Ax^*)^T y^* = (x^*)^T A^T y^* = c^T x^* \end{aligned}$$

The result follows from strong duality theorem (theorem 5.2)

Let  $G = (V, E)$  be a bipartite graph. In order to prove the next theorem, we need to recall the *LP-relaxation* of the minimum weight perfect matching problem (we will call it the *dual LP*) and his dual (which we will call the *primal LP*).

$$\begin{array}{ll}
 \text{PRIMAL} & \text{DUAL} \\
 \max \sum_{v \in V} y_v & \min \sum_{e \in E} w_e \cdot x_e \\
 uv \in E : y_u + y_v \leq w_{uv} & v \in V : \sum_{e \in \delta(v)} x_e = 1 \\
 y \in \mathbb{R}^{|V|} & x \geq 0
 \end{array}$$

**Theorem 10.8.** *Suppose  $x^*$  and  $y^*$  are feasible solutions of the recalled primal and dual linear programs respectively. Then, they are optimal solutions if and only if*

$$\forall uv \in E : x_{uv}^* > 0 \Rightarrow y_u^* + y_v^* = w_{uv}$$

*Proof.* The proof follows immediately from theorem 10.7.

Another direct consequence of theorem 10.7 is the following criterion for the optimality of a minimum weight perfect matching.

Let  $y^* \in \mathbb{R}^{|V|}$  be feasible in the primal LP and let  $G_{y^*} = (V, E_{y^*})$  be the graph with edge set

$$E_{y^*} = \{uv \in E : y_u^* + y_v^* = w_{uv}\}$$

If  $G_{y^*}$  has a perfect matching  $M$ , then  $y^*$  and  $\chi^M$  are optimal solutions of the corresponding linear programs. In particular, since, by definition of  $\chi^M$ ,

$$\chi_{uv}^M > 0 \Rightarrow y_u^* + y_v^* = w_{uv},$$

$M$  is a *minimum weight perfect matching*.

### The Primal-Dual Algorithm

Let  $G = (V, E)$  be a complete bipartite graph with bipartition  $V = A \sqcup B$  and edge weights  $w : E \rightarrow \mathbb{R}$ . The primal-dual algorithm is a method to find a minimum weight matching in  $G$ . It starts from a feasible solution  $y^*$  of the dual LP recalled before theorem 10.8 (notice that a possible choice for an initial feasible solution can always be  $y^*$  such that  $y_v^*$  is equal to 0 or to the smallest edge-weight in the graph in the case it is smaller than 0) and it updates this solution until it is optimal. The last question we have to answer is: how do we update  $y^*$  in such a way that our algorithm terminates correctly?

Let  $G_{y^*}$  be as previously defined and let  $M$  be a maximum cardinality matching in  $G_{y^*}$ . Turn  $G_{y^*}$  into a directed graph by orientating the matching edges from  $A$  to  $B$  and the non-matching edges from  $B$  to  $A$  (call this new

graph  $\tilde{G}_{y^*}$ ). Let  $M$  be a maximum cardinality matching in  $G_{y^*}$ . Let  $L \subseteq A \sqcup B$  be the set of vertices reachable from the set of exposed nodes of  $B$  in the directed graph  $\tilde{G}_{y^*}$ .

*Claim.*  $G_{y^*}$  does not contain an edge  $uv$  with  $u \in A \setminus L$  and  $v \in B \cap L$ .

*Proof.* Let  $uv$  be such an edge. If  $uv$  is a non-matching edge, then, since  $v \in L$ , also  $u$  has to be in  $L$ . Suppose now that  $uv$  is a matching edge. Since  $v \in L$ , there must be a matching edge  $\tilde{u}v$  with  $\tilde{u} \in L$ . This implies that  $u = \tilde{u}$  and this leads to a contradiction since  $u \notin L$ .

We want to update our solution  $y^*$  in such a way to obtain an optimal solution. Let  $\delta = \min\{w_{uv} - y_u^* - y_v^* \mid uv \in E, u \in A \setminus L, v \in B \cap L\}$ . If  $y^*$  is not already optimal, then  $\delta$  is strictly positive. This allows us to construct a feasible solution  $\tilde{y}$  in this way:

- i)  $\tilde{y}_v = y_v^* + \delta$  if  $v \in A \setminus L$
- ii)  $\tilde{y}_v = y_v^* - \delta$  if  $v \in B \cap L$
- iii)  $\tilde{y}_v = y_v^*$  otherwise

This gives us a new graph  $G_{\tilde{y}}$ . Notice that the previous claim shows that all the edges of  $M$  are in  $G_{\tilde{y}}$ . Furthermore, by the definition of  $\delta$ , at least one edge connecting a node in  $A \setminus L$  to a vertex in  $B \cap L$  becomes tight. Thus,  $L$  is augmented. This implies that every time we update our solution there are two things that can happen:

- i)  $G_{\tilde{y}}$  contains an augmenting path with respect to  $M$ . This allows us to increase our matching.
- ii) The set  $L$  increases.

Since  $V$  is a finite set and by the criterion of optimality showed before, we can assert that the Primal-Dual algorithm terminates correctly.

We conclude by resuming the algorithm, giving his running time and an example, which can be useful to understand better how the Primal-Dual algorithm works.

Algorithm 10.2.

**Input:** A feasible solution  $y^*$

**Initialisation:** Turn  $G$  into a directed graph

**while**  $M$  is not optimal

Compute  $L$

Let  $\delta = \min\{w_{uv} - y_u^* - y_v^*\}$

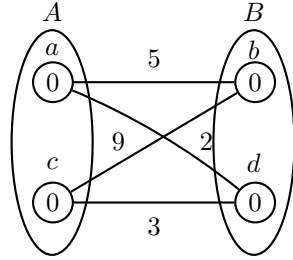
Update  $y^*$

**return**  $M$

**Theorem 10.9.** The Primal-Dual algorithm runs in  $O(|V|^2)$  time.

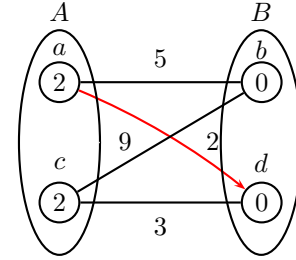
*Proof.* We can have at most  $|V|$  many augmentation of  $L$  without augmenting  $M$  and  $M$  can be augmented at most  $|V|/2$  times.

Figure 10.3 shows how to apply Primal-Dual algorithm on a complete bipartite graph.



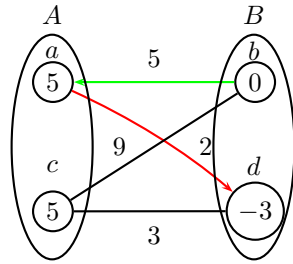
$$y^* = (0, 0, 0, 0), E_{y^*} = \emptyset$$

$$M = \emptyset, L = \{b, d\}, \delta = 2$$



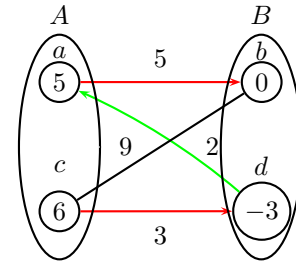
$$y^* = (2, 0, 2, 0), E_{y^*} = \{ad\}$$

$$M = \{ad\}, L = \{b\}, \delta = 3$$



$$y^* = (5, 0, 5, -3), E_{y^*} = \{ad, ba\}$$

$$M = \{ad\}, L = \{a, b, d\}, \delta = 1$$



$$y^* = (5, 0, 6, -3), E_{y^*} = \{ab, da, cd\}$$

$$M = \{ab, cd\}, L = \emptyset$$

Fig. 10.3: The edges of  $G_{y^*}$  are coloured, red for the matching edges, and green for the non-matching edges.

## Exercises

1. Let

$$\begin{aligned} \max \quad & c^T x \\ \text{subject to} \quad & Ax \leq b \end{aligned}$$

be a primal LP and let

$$\begin{aligned} \min \quad & b^T y \\ \text{subject to} \quad & A^T y = c \\ & y \geq 0 \end{aligned}$$

be his dual. Let  $x^*$  and  $y^*$  be primal and dual feasible solutions respectively. Suppose they are both optimal. Is the following true: if  $a_i^T x^* = b_i$ , then  $y_i^* > 0$ ?

2. Find  $G_{y^*}$  from the graph  $G = (V, E)$  shown in figure 10.4 as described in the criterion for the optimality of a minimum weight perfect graph.

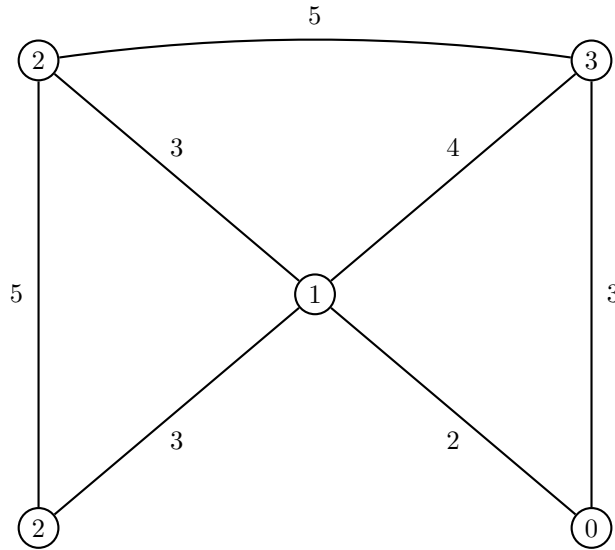


Fig. 10.4:  $G = (V, E)$

3. Are there rational numbers  $y_1, y_2, y_3, y_4$  such that

$$\begin{aligned} y_1 + y_3 &\leq 5 \\ y_1 + y_4 &= 2 \\ y_2 + y_3 &= 9 \\ y_2 + y_4 &\leq 3 \end{aligned}$$

4. Apply Primal-Dual Algorithm to the complete bipartite graph of figure 10.5 (we give an initial feasible solution  $y^* = (0, 0, 0, 0)$ ).

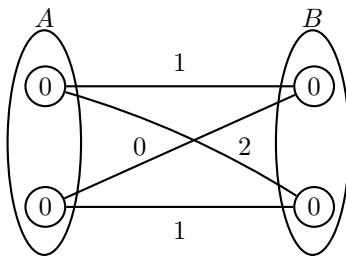


Fig. 10.5





# References

1. A. V. Aho, J. E. Hopcroft, and J. D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, Reading, 1974.
2. J. Edmonds. Maximum matching and a polyhedron with 0,1-vertices. *Journal of Research of the National Bureau of Standards*, 69:125–130, 1965.
3. J. Edmonds. Paths, trees and flowers. *Canadian Journal of Mathematics*, 17:449–467, 1965.
4. F. Eisenbrand, A. Karrenbauer, and C. Xu. Algorithms for longer oled lifetime. In *6th International Workshop on Experimental Algorithms, (WEA 07)*, pages 338–351, 2007.
5. M. Grötschel, L. Lovász, and A. Schrijver. *Geometric Algorithms and Combinatorial Optimization*, volume 2 of *Algorithms and Combinatorics*. Springer, 1988.
6. L. Khachiyan. A polynomial algorithm in linear programming. *Doklady Akademii Nauk SSSR*, 244:1093–1097, 1979.
7. V. Klee and G. J. Minty. How good is the simplex algorithm? In *Inequalities, III (Proc. Third Sympos., Univ. California, Los Angeles, Calif., 1969; dedicated to the memory of Theodore S. Motzkin)*, pages 159–175. Academic Press, New York, 1972.
8. T. Koch. *Rapid Mathematical Programming*. PhD thesis, Technische Universität Berlin, 2004. ZIB-Report 04-58.
9. B. Korte and J. Vygen. *Combinatorial optimization*, volume 21 of *Algorithms and Combinatorics*. Springer-Verlag, Berlin, second edition, 2002. Theory and algorithms.
10. E. L. Lawler. *Combinatorial optimization: networks and matroids*. Holt, Rinehart and Winston, New York, 1976.
11. L. Lovász. Graph theory and integer programming. *Annals of Discrete Mathematics*, 4:141–158, 1979.
12. J. E. Marsden and M. J. Hoffman. *Elementary Classical Analysis*. Freeman, 2 edition, 1993.
13. J. Matoušek and B. Gärtner. *Understanding and Using Linear Programming (Universitext)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.
14. N. Megiddo. Combinatorial optimization with rational objective functions. *Math. Oper. Res.*, 4(4):414–424, 1979.
15. A. S. Nemirovskiy and D. B. Yudin. Informational complexity of mathematical programming. *Izvestiya Akademii Nauk SSSR. Tekhnicheskaya Kibernetika*, (1):88–117, 1983.
16. A. Schrijver. *Theory of Linear and Integer Programming*. John Wiley, 1986.
17. N. Z. Shor. Cut-off method with space extension in convex programming problems. *Cybernetics and systems analysis*, 13(1):94–96, 1977.
18. V. Strassen. Gaussian elimination is not optimal. *Numerische mathematik*, 13(4):354–356, 1969.