
Discrete Optimization

lecture notes, summer term 2025

Yann Disser

June 30, 2025



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Contents

1	Introduction	1
1.1	Examples	2
1.2	Outline	4
2	Basic Polyhedral Theory	5
2.1	Orthogonal projections of polyhedra	5
2.2	Representations of polyhedra	9
2.3	The integer hull	14
3	Branch-and-Bound Method	19
3.1	Complexity of integer programs	19
3.2	The branch-and-bound method	20
4	Integral Polyhedra	29
4.1	Total unimodularity	29
4.2	The Hermite normal form	35
4.3	Total dual integrality	38
5	Cutting Planes	45
5.1	General cutting planes	45
5.2	Specialized cuts	53
6	Decomposition Methods	63
6.1	Lagrangian relaxation	63
6.2	Dantzig-Wolfe decomposition	74
6.3	Benders' decomposition	79
6.4	Connections between these approaches	82
7	Heuristics	85
7.1	The greedy algorithm	85
7.2	Local search	89
8	Approximation Algorithms	93
8.1	Approximation for TSP	94
8.2	Polynomial-time approximation schemes	96
8.3	LP Rounding	98

These lecture notes are based on the lecture notes of the german lecture “Diskrete Optimierung”, held by Marc Pfetsch at TU Darmstadt.

1 Introduction

This lecture is concerned with algorithmic problems of the following variety.

Definition 1.1. A *mixed-integer (linear) program (MIP or MILP)* is an optimization problem of the form

$$\begin{aligned} \max \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & \mathbf{Ax} \leq \mathbf{b}, \\ & \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}, \end{aligned}$$

with instances given by $\mathbf{c} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$, $n, m \in \mathbb{N}$ and $p \in \{1, \dots, n\}$.

Some special cases of a mixed-integer program are:

- linear program (LP): $p = 0$ (see *Introduction to Optimization*),
- integer (linear) program (IP or ILP): $p = n$,
- binary (linear) program (BP or BLP): $\mathbf{x} \in \{0, 1\}^n$.

As we will see, MIPs encompass all problems with linear objective function and a finite set of feasible solutions, such as the minimum spanning tree, the shortest path, the matching, and even the 3SAT problem (see *Algorithmic Discrete Mathematics*). We formalize this class of problems.

Definition 1.2. A *(linear) combinatorial optimization* problem is a problem of the form

$$\max_{F \in \mathcal{F}} \sum_{e \in F} c(e),$$

with instances given by an *objective function* $c(F) := \sum_{e \in F} c(e)$ with $c: E \rightarrow \mathbb{R}$ over a finite *ground set* E and a set $\mathcal{F} \subseteq 2^E$ of *feasible solutions*.

In particular, binary programs have a finite set of feasible solutions.

Observation 1.3. Every binary program can be formulated as a combinatorial optimization problem by setting $E := \{1, \dots, n\}$ and $c(i) := c_i$ for all $i \in \{1, \dots, n\}$, and $\mathcal{F} := \{F \subseteq E : \sum_{j \in F} \mathbf{A}_{.j} \leq \mathbf{b}\}$.

The following proof gives a preview of the typical reasoning we will employ in the lecture. It relies on a fundamental insight that will be shown in Chapter 2.

Proposition 1.4. Linear combinatorial optimization problems can be formulated as LPs.

Proof. Let a combinatorial optimization problem be given by $E = \{e_1, \dots, e_n\}$, $\mathcal{F} \subseteq 2^E$, and $c: E \rightarrow \mathbb{R}$. We define the characteristic vector $\chi^F \in \{0, 1\}^E$ of a set $F \subseteq E$ by

$$\chi_e^F := \begin{cases} 1, & \text{if } e \in F, \\ 0, & \text{otherwise.} \end{cases}$$

and we set $\mathcal{X} := \{\chi^F \in \{0, 1\}^E : F \in \mathcal{F}\}$ and $\bar{c} \in \mathbb{R}^n$ with $\bar{c}_i = c(e_i)$. Then,

$$\begin{aligned} \max_{F \in \mathcal{F}} c(F) &= \max\{\bar{c}^\top \mathbf{x} : \mathbf{x} \in \mathcal{X}\} \\ &\leq \max\{\bar{c}^\top \mathbf{x} : \mathbf{x} \in \text{conv}(\mathcal{X})\}. \end{aligned} \tag{1.1}$$

We will prove later that \mathcal{X} being finite implies that $\text{conv}(\mathcal{X})$ is a polytope (Corollary 2.9), which means that we can find $m \in \mathbb{N}$, $A \in \mathbb{R}^{m \times n}$, and $\mathbf{b} \in \mathbb{R}^m$ such that

$$\text{conv}(\mathcal{X}) = \mathcal{P}(A, \mathbf{b}) = \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}\}. \tag{1.2}$$

We already know that, since $\bar{c}^\top \mathbf{x}$ is linear, the maximum in (1.1) is attained in a vertex of this polytope, and its vertices are contained in \mathcal{X} , since they are extreme points and cannot be written as a convex combination of other points in $\mathcal{X} \subseteq \text{conv}(\mathcal{X})$ (see *Introduction to Optimization*). Hence (1.1) holds with equality. With (1.2), we obtain that

$$\max_{F \in \mathcal{F}} c(F) = \max\{\bar{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}, \mathbf{x} \in \mathbb{R}^n\}$$

can indeed be formulated as a linear program. □

Proposition 1.4 implies in particular that 3SAT (and thus every problem in NP) can be expressed as an LP. With Observation 1.3 it further follows that binary programs can be reduced to LPs. In Chapter 4, we will see that, under suitable assumptions, the same is true for MIPs in general.

On the other hand, we know that LPs can be solved in polynomial time with the ellipsoid method (see *Introduction to Optimization*). How does this fit together?

This apparent contradiction is resolved by the fact that, while it must exist, the LP representation of a combinatorial optimization problem may be difficult to find and it may be large. In particular, already the set of feasible solutions \mathcal{F} of a combinatorial optimization problem is often given implicitly and is often large (such as the one of Observation 1.3). Many of the solution methods we will see are based on finding or approximating the LP representation of a MIP.

1.1 Examples

We dealt with many integer and combinatorial optimization problems in the past (see *Algorithmic Discrete Mathematics* and *Introduction to Optimization*). Let us briefly revisit some important examples.

Example 1.5 (assignment problem). In a company, there are n employees and n jobs that need to be carried out. Each person can perform exactly one job. The cost incurred by the company when person i carries out job j is c_{ij} . How should the company deploy its employees as cost-efficiently as possible? In other words, we aim to find a perfect matching of minimum cost between employees and jobs (see *Algorithmic Discrete Mathematics*).

We use the variables

$$x_{ij} = \begin{cases} 1, & \text{if person } i \text{ is assigned to job } j, \\ 0, & \text{otherwise.} \end{cases}$$

This results in the following binary program:

$$\begin{aligned} \min \quad & \sum_{i=1}^n \sum_{j=1}^n c_{ij} x_{ij} \\ \text{s.t.} \quad & \sum_{j=1}^n x_{ij} = 1, && \text{for } i \in \{1, \dots, n\}, \\ & \sum_{i=1}^n x_{ij} = 1, && \text{for } j \in \{1, \dots, n\}, \\ & x_{ij} \in \{0, 1\}, && \text{for } i, j \in \{1, \dots, n\}. \end{aligned}$$

The problem can be cast as a combinatorial optimization problem by setting $E := \{1, \dots, n\}^2$, $c((i, j)) := c_{ij}$, and

$$\begin{aligned} \mathcal{F} := & \{F \subseteq E : |\{(i, j) \in F : j \in \{1, \dots, n\}\}| = 1 \forall i \in \{1, \dots, n\}\} \\ & \cap \{F \subseteq E : |\{(i, j) \in F : i \in \{1, \dots, n\}\}| = 1 \forall j \in \{1, \dots, n\}\} \end{aligned} \quad \Delta$$

Example 1.6 (knapsack problem). A thief wants to pack stolen goods of as large a value as possible in their knapsack. Each item $i \in \{1, \dots, n\}$ has weight $a_i > 0$ and value $c_i > 0$. The knapsack has capacity $\beta > 0$. Which items should the thief pack?

We use the variables

$$x_i = \begin{cases} 1, & \text{if item } i \text{ is selected,} \\ 0, & \text{otherwise.} \end{cases}$$

This results in the following binary program:

$$\begin{aligned} \max \quad & \sum_{i=1}^n c_i x_i \\ \text{s.t.} \quad & \sum_{i=1}^n a_i x_i \leq \beta, \\ & \mathbf{x} \in \{0, 1\}^n. \end{aligned}$$

The problem can be cast as a combinatorial optimization problem by setting $E := \{1, \dots, n\}$, $c(j) := c_j$, and

$$\mathcal{F} := \{F \subseteq E : \sum_{i \in F} a_i \leq \beta\}. \quad \Delta$$

Example 1.7 (set-packing/partitioning/covering problems). We are given a finite ground set $U = \{1, \dots, m\}$ and a family $\{S_j\}_{j \in \{1, \dots, n\}}$ of subsets of U . Each subset S_j is associated with a cost/utility c_j . The problem of finding a minimum cost / maximum utility selection of subsets such that every $e \in U$ is contained in at most, exactly, resp. at least one subset is called *set-packing*, *set-partitioning*, resp. *set-covering* problem. Examples are placing wind turbines without overlap in their areas of effect (set-packing), dividing a country into electoral districts (set-partitioning), and positioning schools in a city (set-covering).

We use the variables

$$x_j = \begin{cases} 1, & \text{if } S_j \text{ is selected,} \\ 0, & \text{otherwise.} \end{cases}$$

We also define the binary matrix $A \in \{0, 1\}^{m \times n}$, with $A_{ij} = 1$ if and only if $i \in S_j$. This yields the following binary programs:

$$\begin{aligned} \min / \max \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & \mathbf{Ax} \begin{cases} \leq \\ = \\ \geq \end{cases} \mathbf{1}, \\ & \mathbf{x} \in \{0, 1\}^n. \end{aligned}$$

The problem can be cast as a combinatorial optimization problem by setting $E := \{1, \dots, n\}$, $c(j) := c_j$, and

$$\mathcal{F} = \{F \subseteq E : |\{j \in F : i \in S_j\}| \begin{cases} \leq \\ = \\ \geq \end{cases} 1 \forall i \in U\}. \quad \triangle$$

Remark 1.8. The problem of Example 1.5 is polynomially time solvable, while the problems in Examples 1.6 and 1.7 are NP-hard (see *Algorithmic Discrete Mathematics*).

1.2 Outline

The main topic of this lecture are techniques for solving mixed-integer programs. We begin (Chapter 2) by extending our structural understanding of polyhedra from the lecture *Introduction to Optimization*. In particular, we will see that a set $P \subseteq \mathbb{R}^n$ is a polyhedron if and only if it can be written as the sum of a polytope and a polyhedral cone.

Subsequently, we present a general “Branch-and-Bound” approach, which reduces solving a MIP to repeatedly solving LPs (Chapter 3). The running time of this procedure is exponential in the worst case, and we will show that it is generally NP-complete just to decide whether a linear program has an integer solution. In contrast, we already know (see *Introduction to Optimization*) that *fractional* vertex solutions can be found in polynomial time. We can therefore efficiently solve integer programs when all vertex solutions of the LP relaxation are integral. We will show properties under which integer programs have this structure (Chapter 4).

We will then look at how we can use additional inequalities to exclude fractional solutions of the LP relaxation (Chapter 5), and how we can separately handle problematic inequalities/variables (Chapter 6).

Finally, we deal with alternative solution methods, such as heuristics (Chapter 7) and approximation algorithms (Chapter 8).

2 Basic Polyhedral Theory

We have already learned some simple facts about polyhedra (see *Introduction to Optimization*). In this chapter, we extend this understanding by techniques that we are going to need later.

Some of the arguments of this chapter hold for vector spaces over the field \mathbb{Q} in addition to over \mathbb{R} . To avoid repeating arguments, we work over $\mathbb{K} \in \{\mathbb{Q}, \mathbb{R}\}$ for the time being. In particular, for $\mathbb{K} = \mathbb{Q}$, we consider polyhedra as subsets of \mathbb{Q}^n for $n \in \mathbb{N}$.

2.1 Orthogonal projections of polyhedra

Projections are an important tool for establishing structural statements about polyhedra. We now discuss a method to compute orthogonal projections of a polyhedron, which amounts to the elimination of variables.

Definition 2.1. The *orthogonal projection* of a set $S \subseteq \mathbb{K}^n$ with respect to its k -th coordinate, $k \in \{1, \dots, n\}$ is

$$\text{Proj}_k(S) := \{\mathbf{x} \in \mathbb{K}^n : x_k = 0, \exists \lambda \in \mathbb{K} \text{ with } \mathbf{x} + \lambda \mathbf{e}^k \in S\}.$$

We are particularly interested in the case where $S = \mathcal{P}(A, \mathbf{b}) := \{\mathbf{x} \in \mathbb{K}^n : A\mathbf{x} \leq \mathbf{b}\}$ is a non-empty polyhedron with $A \in \mathbb{K}^{m \times n}$, $\mathbf{b} \in \mathbb{K}^m$ (see Figure 2.1). To construct the orthogonal projection with respect to the k -th coordinate of $\mathcal{P}(A, \mathbf{b})$, we divide the indices of the inequalities in $A\mathbf{x} \leq \mathbf{b}$ into the following three sets:

- $C_k^- := \{i \in \{1, \dots, m\} : A_{ik} < 0\}$ (negative coefficients),
- $C_k^0 := \{i \in \{1, \dots, m\} : A_{ik} = 0\}$ (zero coefficients),
- $C_k^+ := \{i \in \{1, \dots, m\} : A_{ik} > 0\}$ (positive coefficients).

Obviously, the inequalities $A_i \cdot \mathbf{x} \leq b_i$ for $i \in C_k^0$, do not contain the variable x_k and are therefore inherited by the projection.

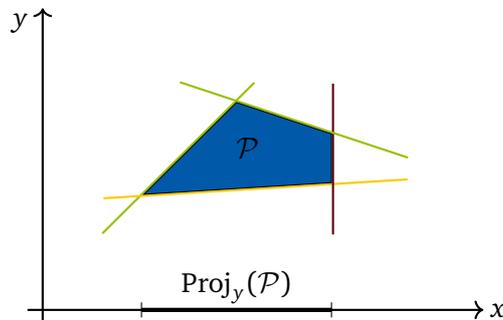


Figure 2.1: Example of an orthogonal projection with respect to the second coordinate. The colors of the hyperplanes correspond to the sets C_y^- (yellow), C_y^0 (red), and C_y^+ (green).

For $s \in C_k^-$ (i.e., $A_{sk} < 0$) we can isolate x_k in $A_s \mathbf{x} \leq b_s$ to obtain

$$x_k \geq \frac{b_s}{A_{sk}} - \sum_{\substack{j=1 \\ j \neq k}}^n \frac{A_{sj}}{A_{sk}} x_j. \quad (2.1)$$

Similarly, for $t \in C_k^+$ we obtain

$$x_k \leq \frac{b_t}{A_{tk}} - \sum_{\substack{j=1 \\ j \neq k}}^n \frac{A_{tj}}{A_{tk}} x_j. \quad (2.2)$$

Combining (2.1) and (2.2) yields

$$\frac{b_s}{A_{sk}} - \sum_{\substack{j=1 \\ j \neq k}}^n \frac{A_{sj}}{A_{sk}} x_j \leq \frac{b_t}{A_{tk}} - \sum_{\substack{j=1 \\ j \neq k}}^n \frac{A_{tj}}{A_{tk}} x_j \quad \Leftrightarrow \quad \sum_{\substack{j=1 \\ j \neq k}}^n (A_{tk} A_{sj} - A_{sk} A_{tj}) x_j \leq A_{tk} b_s - A_{sk} b_t. \quad (2.3)$$

Observe that we obtained a new inequality which no longer contains the variable x_k . In particular, the k -th column of the resulting coefficient matrix is $\mathbf{0}$.

Lemma 2.2. The vector $\tilde{\mathbf{x}} \in \mathbb{K}^n$ lies in $\text{Proj}_k(\mathcal{P}(\mathbf{A}, \mathbf{b}))$ if and only if $\tilde{x}_k = 0$ and $\tilde{\mathbf{x}}$ satisfies the inequalities (2.3) for all $(s, t) \in C_k^- \times C_k^+$ and the inequalities $\mathbf{A}_i \mathbf{x} \leq b_i$ for all $i \in C_k^0$.

Proof. Let $\tilde{\mathbf{x}} \in \text{Proj}_k(\mathcal{P}(\mathbf{A}, \mathbf{b}))$ and thus, in particular, $\tilde{x}_k = 0$. By assumption, there exists $\lambda \in \mathbb{K}$ such that $\mathbf{x} := \tilde{\mathbf{x}} + \lambda \mathbf{e}^k \in \mathcal{P}(\mathbf{A}, \mathbf{b})$, i.e., $\mathbf{A} \mathbf{x} \leq \mathbf{b}$. From the derivation above it follows that \mathbf{x} then also fulfills the inequalities (2.3). Since these do not contain x_k , they are also fulfilled by $\tilde{\mathbf{x}}$. Of course, this also applies to the inequalities belonging to $i \in C_k^0$.

Now assume that $\tilde{x}_k = 0$ and $\tilde{\mathbf{x}}$ satisfies the inequalities (2.3) for all $(s, t) \in C_k^- \times C_k^+$ and the inequalities $\mathbf{A}_i \mathbf{x} \leq b_i$ for all $i \in C_k^0$. Then, in particular, the maximum of the right-hand sides of (2.1) over all $s \in C_k^-$ cannot be greater than the minimum of the right-hand sides of (2.2) over all $t \in C_k^+$. We can therefore set $\lambda = x_k$ to an arbitrary value between these two bounds to satisfy all inequalities (2.1) and (2.2). The resulting vector $\mathbf{x} = \tilde{\mathbf{x}} + \lambda \mathbf{e}^k$ then lies in $\mathcal{P}(\mathbf{A}, \mathbf{b})$. Hence, $\tilde{\mathbf{x}} \in \text{Proj}_k(\mathcal{P}(\mathbf{A}, \mathbf{b}))$. \square

We can turn Lemma 2.2 into the following algorithm, the so-called *Fourier-Motzkin elimination*.

Algorithm: Fourier-Motzkin elimination

input: polyhedron $\mathcal{P}(\mathbf{A}, \mathbf{b})$, $\mathbf{A} \in \mathbb{K}^{m \times n}$, $\mathbf{b} \in \mathbb{K}^m$, index $k \in \{1, \dots, n\}$

output: $\mathcal{P}(\mathbf{D}, \mathbf{d})$ such that $\text{Proj}_k(\mathcal{P}(\mathbf{A}, \mathbf{b})) = \{\mathbf{x} \in \mathbb{K}^n : \mathbf{D} \mathbf{x} \leq \mathbf{d}, x_k = 0\}$

$C_k^- := \{i \in \{1, \dots, m\} : A_{ik} < 0\}$, $C_k^0 := \{i \in \{1, \dots, m\} : A_{ik} = 0\}$, $C_k^+ := \{i \in \{1, \dots, m\} : A_{ik} > 0\}$

$i \leftarrow 1$

for $z \in C_k^0$:

$\mathbf{D}_i \leftarrow \mathbf{A}_z$; $d_i \leftarrow b_z$; $i \leftarrow i + 1$

for $(s, t) \in C_k^- \times C_k^+$:

$\mathbf{D}_i \leftarrow A_{tk} \mathbf{A}_s - A_{sk} \mathbf{A}_t$

$d_i \leftarrow A_{tk} b_s - A_{sk} b_t$

$i \leftarrow i + 1$

(it follows that $D_{ik} = 0$)

return $\mathcal{P}(\mathbf{D}, \mathbf{d})$

We note the following immediate consequences.

Corollary 2.3. Let $\mathcal{P}(D, \mathbf{d})$ be the result of the Fourier-Motzkin elimination. We have

- (a) $\text{Proj}_k(\mathcal{P}(A, \mathbf{b})) = \{\mathbf{x} \in \mathbb{K}^n : x_k = 0, D\mathbf{x} \leq \mathbf{d}\};$
- (b) $D\mathbf{e}^k = \mathbf{0};$
- (c) The rows of $D\mathbf{x} \leq \mathbf{d}$ are non-negative linear combinations of the rows of $A\mathbf{x} \leq \mathbf{b}$, i.e., there exists a matrix $\bar{U} \geq 0$ with $\bar{U}A = D$ and $\bar{U}\mathbf{b} = \mathbf{d};$
- (d) $\mathcal{P}(A, \mathbf{b}) \neq \emptyset \Leftrightarrow \mathcal{P}(D, \mathbf{d}) \neq \emptyset \Leftrightarrow \text{Proj}_k(\mathcal{P}(A, \mathbf{b})) \neq \emptyset.$

Proof. Statement (a) holds because of Lemma 2.2 and the definition the Fourier-Motzkin elimination. Statement (b) follows from the fact that the algorithm ensures $D_{ik} = A_{tk}A_{sk} - A_{sk}A_{tk} = 0$. Statement (c) holds by definition of the algorithm with $A_{tk} > 0$ and $A_{sk} < 0$. It remains to prove (d).

If $\mathbf{x} \in \mathcal{P}(A, \mathbf{b})$ exists, then by (c) it holds that

$$D\mathbf{x} = \bar{U}A\mathbf{x} \stackrel{\bar{U} \geq 0}{\leq} \bar{U}\mathbf{b} = \mathbf{d}.$$

Thus, $\mathbf{x} \in \mathcal{P}(D, \mathbf{d})$ and therefore $\mathcal{P}(D, \mathbf{d}) \neq \emptyset$.

If $\mathbf{x} \in \mathcal{P}(D, \mathbf{d})$ exists, then, because of (b), also $\tilde{\mathbf{x}} := \sum_{i \neq k} x_i \mathbf{e}^i \in \mathcal{P}(D, \mathbf{d})$. Because of (a) and $\tilde{x}_k = 0$, it further follows that $\tilde{\mathbf{x}} \in \text{Proj}_k(\mathcal{P}(A, \mathbf{b}))$ and therefore $\text{Proj}_k(\mathcal{P}(A, \mathbf{b})) \neq \emptyset$.

If, in turn, $\mathbf{x} \in \text{Proj}_k(\mathcal{P}(A, \mathbf{b}))$ exists, then, by definition, there is $\lambda \in \mathbb{K}$ with $\mathbf{x} + \lambda \mathbf{e}^k \in \mathcal{P}(A, \mathbf{b})$ and therefore $\mathcal{P}(A, \mathbf{b}) \neq \emptyset$. \square

Corollary 2.3 (d) provides a certificate for the infeasibility of the system $A\mathbf{x} \leq \mathbf{b}$ via the infeasibility of the lower dimensional system $D\mathbf{x} \leq \mathbf{d}$. By iteratively projecting out all variables, we eventually obtain a zero-dimensional certificate in the form of a single point. It turns out that this point is exactly the certificate provided by the Farkas lemma (see *Introduction to Optimization*).

To see this, we start by applying Fourier-Motzkin elimination to the first variable, yielding $\mathcal{P}(D^{(1)}, \mathbf{d}^{(1)})$. Corollary 2.3 implies

- (i) $D^{(1)}\mathbf{e}^1 = \mathbf{0}.$
- (ii) There is a matrix $\bar{U}^{(1)} \geq 0$ with $\bar{U}^{(1)}A = D^{(1)}$, $\bar{U}^{(1)}\mathbf{b} = \mathbf{d}^{(1)}.$ (2.4)
- (iii) $\mathcal{P}(A, \mathbf{b}) = \emptyset \iff \mathcal{P}(D^{(1)}, \mathbf{d}^{(1)}) = \emptyset.$

In this manner, we can iteratively eliminate the first k variables. In the k -th step, we obtain a matrix $D^{(k)}$ and $\bar{U}^{(k)} \geq 0$ with $D^{(k)} = \bar{U}^{(k)}D^{(k-1)}$. We set $U^{(k)} = \bar{U}^{(k)}\bar{U}^{(k-1)} \dots \bar{U}^{(1)}$ such that $U^{(k)}A = D^{(k)}$ and $U^{(k)}\mathbf{b} = \mathbf{d}^{(k)}$, and obtain

- (i) $D^{(k)}\mathbf{e}^j = \mathbf{0}$, for all $j \in \{1, \dots, k\}.$
- (ii) There is a matrix $U^{(k)} \geq 0$ with $U^{(k)}A = D^{(k)}$ and $U^{(k)}\mathbf{b} = \mathbf{d}^{(k)}.$ (2.5)
- (iii) $\mathcal{P}(A, \mathbf{b}) = \emptyset \iff \mathcal{P}(D^{(k)}, \mathbf{d}^{(k)}) = \emptyset.$

Now let $\mathcal{P}(A, \mathbf{b}) = \emptyset$. According to (2.5) (iii), $\mathcal{P}(D^{(n)}, \mathbf{d}^{(n)}) = \emptyset$. Since $D^{(n)} = 0$ because of (2.5) (i), this implies that an index i must exist with $d_i^{(n)} < 0$. According to (2.5) (ii), it then holds for $\mathbf{u} := (U_i^{(n)})^\top$ that $\mathbf{u} \geq \mathbf{0}$,

$\mathbf{u}^\top A = \mathbf{D}_i^{(n)} = \mathbf{0}^\top$ and $\mathbf{u}^\top \mathbf{b} = d_i^n < 0$. Conversely, it immediately follows from the existence of such a $\mathbf{u} \geq \mathbf{0}$ that $\mathcal{P}(A, \mathbf{b}) = \emptyset$, since then

$$\mathbf{u}^\top A \mathbf{x} = 0 > d_i^n = \mathbf{u}^\top \mathbf{b} \quad \forall \mathbf{x} \in \mathbb{K}^n.$$

This means that exactly one of the following two systems has a solution (see *Introduction to Optimization*):

$$\exists \mathbf{x} : A \mathbf{x} \leq \mathbf{b} \quad \vee \quad \exists \mathbf{u} \geq \mathbf{0} : \mathbf{u}^\top A = \mathbf{0}, \mathbf{u}^\top \mathbf{b} < 0.$$

Moreover, Corollary 2.3 (a) immediately implies the following.

Corollary 2.4. The orthogonal projection $\text{Proj}_{\mathcal{K}}(\mathcal{P}(A, \mathbf{b}))$ is a polyhedron.

This allows to prove an important fact that opens a new perspective on polyhedra. We will later show that every polyhedron can be written in the following form.¹

Theorem 2.5. For every $A \in \mathbb{K}^{m \times n}$ and $B \in \mathbb{K}^{m \times n'}$, the set $\mathcal{P} = \text{conv}(A) + \text{cone}(B)$ is a polyhedron.

Proof. We have

$$\begin{aligned} \mathcal{P} &= \{\mathbf{x} \in \mathbb{K}^m : \exists \boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0} : \mathbf{1}^\top \boldsymbol{\lambda} = 1, \mathbf{x} = A\boldsymbol{\lambda} + B\boldsymbol{\mu}\} \\ &= \{\mathbf{x} \in \mathbb{K}^m : \exists \boldsymbol{\lambda}, \boldsymbol{\mu} : D \begin{pmatrix} \boldsymbol{\lambda} \\ \boldsymbol{\mu} \\ \mathbf{x} \end{pmatrix} \leq \mathbf{d}\}, \end{aligned}$$

with

$$D = \begin{pmatrix} -\mathbb{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & -\mathbb{I} & \mathbf{0} \\ \mathbf{1}^\top & \mathbf{0}^\top & \mathbf{0}^\top \\ -\mathbf{1}^\top & \mathbf{0}^\top & \mathbf{0}^\top \\ A & B & -\mathbb{I} \\ -A & -B & \mathbb{I} \end{pmatrix} \quad \text{and} \quad \mathbf{d} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ 1 \\ -1 \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix}.$$

Via iterative orthogonal projection of the variables in $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ in $D \begin{pmatrix} \boldsymbol{\lambda} \\ \boldsymbol{\mu} \\ \mathbf{x} \end{pmatrix} \leq \mathbf{d}$, by Corollary 2.4, we obtain a polyhedron $\text{Proj}_{\boldsymbol{\lambda}, \boldsymbol{\mu}}(\mathcal{P}(D, \mathbf{d}))$. This polyhedron is exactly \mathcal{P} embedded in $\mathbb{K}^{n+n'+m}$, i.e.,

$$\mathcal{P} = \{\mathbf{x} \in \mathbb{K}^m : \begin{pmatrix} \mathbf{0} \\ \mathbf{x} \end{pmatrix} \in \text{Proj}_{\boldsymbol{\lambda}, \boldsymbol{\mu}}(\mathcal{P}(D, \mathbf{d}))\},$$

or, conversely,

$$\text{Proj}_{\boldsymbol{\lambda}, \boldsymbol{\mu}}(\mathcal{P}(D, \mathbf{d})) = \{\begin{pmatrix} \mathbf{0} \\ \mathbf{x} \end{pmatrix} \in \mathbb{K}^{n+n'+m} : \mathbf{x} \in \mathcal{P}\}. \quad \square$$

Observe that, in the proof of Theorem 2.5, the only non-zero entries of \mathbf{d} are associated with $\boldsymbol{\lambda}$, i.e., with $\text{conv}(A)$. It follows that every conic hull is a polyhedral cone.

Corollary 2.6. For every $B \in \mathbb{K}^{m \times n}$, there exists a matrix $D \in \mathbb{K}^{k \times m}$ with $\text{cone}(B) = \mathcal{P}(D, \mathbf{0})$.

The argument of Theorem 2.5 essentially uses that any affine map of a polyhedron is a polyhedron. In particular, analogous arguments yield the following corollaries (exercises).

Corollary 2.7. For every $A \in \mathbb{K}^{m \times n}$, the sets $\text{lin}(A)$, $\text{aff}(A)$, $\text{cone}(A)$, and $\text{conv}(A)$ are polyhedra.

¹Throughout the lecture, for a matrix $A \in \mathbb{K}^{m \times n}$, we denote $\text{lin}(A) := \text{lin}(\{A_1, \dots, A_n\}) = \{A\boldsymbol{\lambda} : \boldsymbol{\lambda} \in \mathbb{K}^n\}$ and analogously for $\text{aff}(A) = \{A\boldsymbol{\lambda} : \mathbf{1}^\top \boldsymbol{\lambda} = 1\}$, $\text{cone}(A) = \{A\boldsymbol{\lambda} : \boldsymbol{\lambda} \geq \mathbf{0}\}$ and $\text{conv}(A) = \{A\boldsymbol{\lambda} : \mathbf{1}^\top \boldsymbol{\lambda} = 1, \boldsymbol{\lambda} \geq \mathbf{0}\}$.

Corollary 2.8. The Minkowski sum $\mathcal{P} + \mathcal{P}'$ of two polyhedra $\mathcal{P}, \mathcal{P}' \subseteq \mathbb{K}^n$ is a polyhedron.

We can already observe that not only does $\text{conv}(A)$ induce a polytope, in fact *every* polytope can be expressed this way. In the next section, we will complement Theorem 2.5 in same way for polyhedra.

Corollary 2.9. A set $\mathcal{P} \subseteq \mathbb{K}^n$ is a polytope if and only if there exists a finite set $\mathcal{V} \subseteq \mathbb{K}^n$ with $\mathcal{P} = \text{conv}(\mathcal{V})$.

Proof. Let $\mathcal{V} \subseteq \mathbb{K}^n$ be finite and $\mathcal{P} = \text{conv}(\mathcal{V})$. Then, \mathcal{P} is a polyhedron by Corollary 2.7, and \mathcal{P} is bounded by

$$\|\mathbf{x}\| \leq \sum_{\mathbf{v} \in \mathcal{V}} \|\mathbf{v}\| \quad \forall \mathbf{x} \in \mathcal{P}.$$

Hence, \mathcal{P} is a polytope.

Conversely, recall that every polytope $\mathcal{P} = \mathcal{P}^-(A, \mathbf{b})$ has a finite set of extreme points (i.e., vertices), each arising from a combination (i.e., a basis) of the columns of A – and as a compact and convex set, \mathcal{P} is given by the convex hull of its extreme points (see *Introduction to Optimization*). \square

2.2 Representations of polyhedra

Let's take another look at the Farkas lemma from a different point of view:

$$\exists \mathbf{x} \geq \mathbf{0}: A\mathbf{x} = \mathbf{b} \quad \vee \quad \exists \mathbf{y}: \mathbf{y}^\top A \leq \mathbf{0}^\top, \mathbf{y}^\top \mathbf{b} > 0,$$

or, expressed differently,

$$\exists \mathbf{x} \geq \mathbf{0}: A\mathbf{x} = \mathbf{b} \quad \iff \quad \forall \mathbf{y}: A^\top \mathbf{y} \leq \mathbf{0} \Rightarrow \mathbf{y}^\top \mathbf{b} \leq 0. \quad (2.6)$$

This characterizes all right-hand sides \mathbf{b} for which the system $A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}$ is feasible. By definition,

$$\text{cone}(A) = \{\mathbf{b} \in \mathbb{K}^m : \exists \mathbf{x} \geq \mathbf{0} \text{ with } A\mathbf{x} = \mathbf{b}\}.$$

Together with the Farkas lemma (2.6), this yields the following observation.

Proposition 2.10. For all matrices $A \in \mathbb{K}^{m \times n}$ it holds that

$$\text{cone}(A) = \{\mathbf{b} \in \mathbb{K}^m : \mathbf{y}^\top \mathbf{b} \leq 0 \forall \mathbf{y} \in \mathcal{P}(A^\top, \mathbf{0})\}.$$

The geometric interpretation of Proposition 2.10 could be written as

$$\left. \begin{array}{l} \text{feasible right-hand sides } \mathbf{b} \\ \text{of } A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0} \end{array} \right\} \triangleq \left. \begin{array}{l} \text{vectors that form an obtuse } (\geq \pi/2) \text{ angle} \\ \text{with all vectors from } \mathcal{P}(A^\top, \mathbf{0}) \end{array} \right\}.$$

The latter can be expressed more abstractly using the following definition.

Definition 2.11. The *polar cone* S° of $S \subseteq \mathbb{K}^n$ is the set of vectors that form an obtuse angle with all vectors in S , i.e.,

$$S^\circ := \{\mathbf{y} \in \mathbb{K}^n : \mathbf{y}^\top \mathbf{x} \leq 0 \forall \mathbf{x} \in S\}.$$

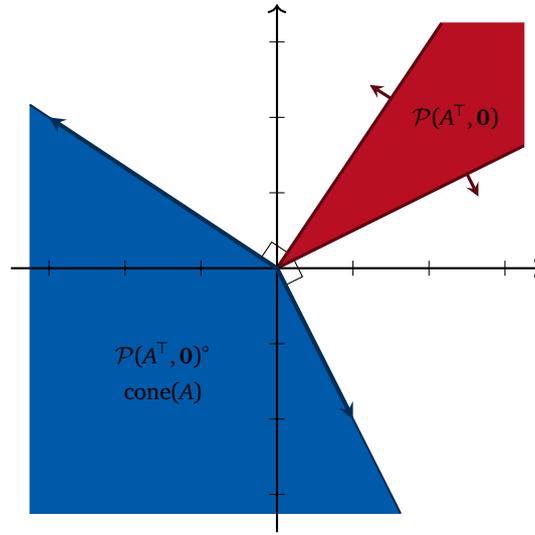


Figure 2.2: Illustration of Example 2.13 with $\mathcal{P}(A^\top, \mathbf{0})^\circ = \text{cone}(\{(-3, 1), (-2, -2)\})$.

With this, we can write Proposition 2.10 more concisely.²

Corollary 2.12. For all $A \in \mathbb{K}^{m \times n}$ it holds that $\text{cone}(A) = \mathcal{P}(A^\top, \mathbf{0})^\circ$.

Our reasoning so far can be summarized as

$$Ax = b, x \geq \mathbf{0} \text{ is feasible} \stackrel{(2.6)}{\iff} b \in \mathcal{P}(A^\top, \mathbf{0})^\circ \stackrel{\text{Cor. 2.12}}{\iff} b \in \text{cone}(A).$$

Example 2.13. For

$$A = \begin{pmatrix} -3 & 1 \\ 2 & -2 \end{pmatrix}$$

we have

$$\mathcal{P}(A^\top, \mathbf{0}) = \left\{ \begin{pmatrix} x \\ y \end{pmatrix} \in \mathbb{K}^2 : x/2 \leq y \leq 3x/2 \right\},$$

and indeed $\mathcal{P}(A^\top, \mathbf{0})^\circ = \text{cone}(A)$ (see Figure 2.2). Accordingly,

$$Ax = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, x \geq \mathbf{0} \text{ is infeasible, but } Ax = \begin{pmatrix} -1 \\ 0 \end{pmatrix}, x \geq \mathbf{0} \text{ is feasible.} \quad \triangle$$

To reverse Corollary 2.12, we observe that the polar cone already is a conic hull and vice-versa.

Lemma 2.14. For $S \subseteq \mathbb{K}^n$, it holds that

$$S^\circ = \text{cone}(S^\circ) = \text{cone}(S)^\circ.$$

Proof. exercise. □

We can now complement Corollary 2.12.

²For convenience, we write $\mathcal{P}(A, \mathbf{b})^\circ := (\mathcal{P}(A, \mathbf{b}))^\circ$, $\text{cone}(S)^\circ := (\text{cone}(S))^\circ$, and $S^{\circ\circ} := (S^\circ)^\circ$.

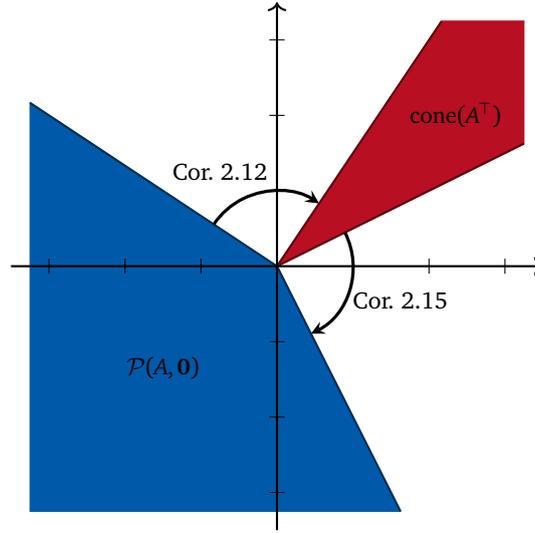


Figure 2.3: Correspondence between a polyhedral cone and a conic hull in “polar row space”, and vice versa.

Corollary 2.15. For all $A \in \mathbb{K}^{m \times n}$ it holds that $\mathcal{P}(A, \mathbf{0}) = \text{cone}(A^\top)^\circ$.

Proof. We have

$$(\text{cone}(A^\top))^\circ \stackrel{\text{Lem. 2.14}}{=} \{(A_1 \cdot)^\top, \dots, (A_m \cdot)^\top\}^\circ \stackrel{\text{Def. 2.11}}{=} \{y : y^\top A^\top \leq \mathbf{0}^\top\} = \{x : Ax \leq \mathbf{0}\} = \mathcal{P}(A, \mathbf{0}). \quad \square$$

We have all ingredients in place to prove that both polyhedral cones of the form $\mathcal{P}(A, \mathbf{0})$ and conic hulls of the form $\text{cone}(A)$ induce themselves as polar cones (see Figure 2.3).

Theorem 2.16. For every $A \in \mathbb{K}^{m \times n}$, it holds that

$$\begin{aligned} \mathcal{P}(A, \mathbf{0})^{\circ\circ} &= \mathcal{P}(A, \mathbf{0}), \\ \text{cone}(A)^{\circ\circ} &= \text{cone}(A). \end{aligned}$$

Proof. We use the relationship between polyhedral cones and conic hulls established in Corollaries 2.12 and 2.15:

$$\begin{aligned} \mathcal{P}(A, \mathbf{0}) &\stackrel{\text{Cor. 2.15}}{=} \text{cone}(A^\top)^\circ \stackrel{\text{Cor. 2.12}}{=} \mathcal{P}(A, \mathbf{0})^{\circ\circ}, \\ \text{cone}(A) &\stackrel{\text{Cor. 2.12}}{=} \mathcal{P}(A^\top, \mathbf{0})^\circ \stackrel{\text{Cor. 2.15}}{=} \text{cone}(A)^{\circ\circ}. \end{aligned} \quad \square$$

We now show that the objects in Theorem 2.16 are not only intimately related but actually different representations of the same sets.

Theorem 2.17 (Minkowski 1896). For every matrix $A \in \mathbb{K}^{m \times n}$ there is a matrix $B \in \mathbb{K}^{n \times k}$ with

$$\mathcal{P}(A, \mathbf{0}) = \text{cone}(B),$$

and vice versa.

Proof. For a given matrix B , the corresponding matrix A exists by Corollary 2.6. For given A , again by Corollary 2.6, there exists a matrix B with (see Figure 2.3)

$$\mathcal{P}(A, \mathbf{0}) \stackrel{\text{Cor. 2.15}}{=} \text{cone}(A^\top)^\circ \stackrel{\text{Cor. 2.6}}{=} \mathcal{P}(B^\top, \mathbf{0})^\circ \stackrel{\text{Cor. 2.12}}{=} \text{cone}(B). \quad \square$$

By reducing the general case to the conic case, we now prove the dual statement to Theorem 2.5, i.e., every polyhedron can be decomposed into a bounded and a conic part.

Theorem 2.18. For every $A \in \mathbb{K}^{m \times n}$, $\mathbf{b} \in \mathbb{K}^m$ there exist finite sets $\mathcal{V}, \mathcal{E} \subseteq \mathbb{K}^n$ with

$$\mathcal{P}(A, \mathbf{b}) = \text{conv}(\mathcal{V}) + \text{cone}(\mathcal{E}).$$

Proof. We embed $\mathcal{P}(A, \mathbf{b})$ into a polyhedral cone

$$\mathcal{H} = \mathcal{P}\left(\left(\begin{array}{cc} A & -\mathbf{b} \\ \mathbf{0}^\top & -1 \end{array}\right), \begin{pmatrix} \mathbf{0} \\ 0 \end{pmatrix}\right).$$

In particular, we add an additional dimension and place our polyhedron onto the plane $\left\{\begin{pmatrix} \lambda \\ 1 \end{pmatrix}\right\}_{\lambda \in \mathbb{K}^n}$, i.e., $\mathbf{x} \in \mathcal{P}(A, \mathbf{b}) \Leftrightarrow \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \in \mathcal{H}$.

According to Theorem 2.17, there is a matrix $B \in \mathbb{K}^{(n+1) \times d}$ with $\mathcal{H} = \text{cone}(B)$. Due to the last row in the description of \mathcal{H} , the last row of B has only non-negative entries. By scaling and swapping the columns of B we can transform B into a matrix \bar{B} with $\text{cone}(\bar{B}) = \text{cone}(B) = \mathcal{H}$, so that

$$\bar{B} = \left(\begin{array}{cc} \bar{V} & \bar{E} \\ \mathbf{1}^\top & \mathbf{0}^\top \end{array}\right),$$

where \bar{V}, \bar{E} may be empty. Let \mathcal{V} and \mathcal{E} be the sets of the columns of the matrices \bar{V} and \bar{E} , respectively. This means that

$$\begin{aligned} \mathbf{x} \in \mathcal{P}(A, \mathbf{b}) &\iff \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \in \mathcal{H} \\ &\iff \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \in \text{cone}(\bar{B}) \\ &\iff \exists \boldsymbol{\lambda}, \boldsymbol{\mu} \geq 0: \begin{array}{l} \mathbf{x} = \bar{V}\boldsymbol{\lambda} + \bar{E}\boldsymbol{\mu} \\ 1 = \mathbf{1}^\top \boldsymbol{\lambda} \end{array} \\ &\iff \mathbf{x} \in \text{conv}(\mathcal{V}) + \text{cone}(\mathcal{E}). \end{aligned} \quad \square$$

Together, Theorems 2.5 and 2.18 yield the central structure theorem of polyhedral theory.

Theorem 2.19 (representation theorem). A subset $\mathcal{P} \subseteq \mathbb{K}^n$ is a polyhedron if, and only if there are finite sets $\mathcal{V}, \mathcal{E} \subseteq \mathbb{K}^n$ with

$$\mathcal{P} = \text{conv}(\mathcal{V}) + \text{cone}(\mathcal{E}).$$

Intuitively, the convex hull in the representation of a polyhedron corresponds to its bounded part, while the conic hull corresponds to the unbounded part. Accordingly, a polytope, i.e., a bounded polyhedron reduces to a convex hull. We can phrase Theorem 2.19 by saying that every polyhedron is the sum of a polytope (Corollary 2.9) and a polyhedral cone (Theorem 2.17).

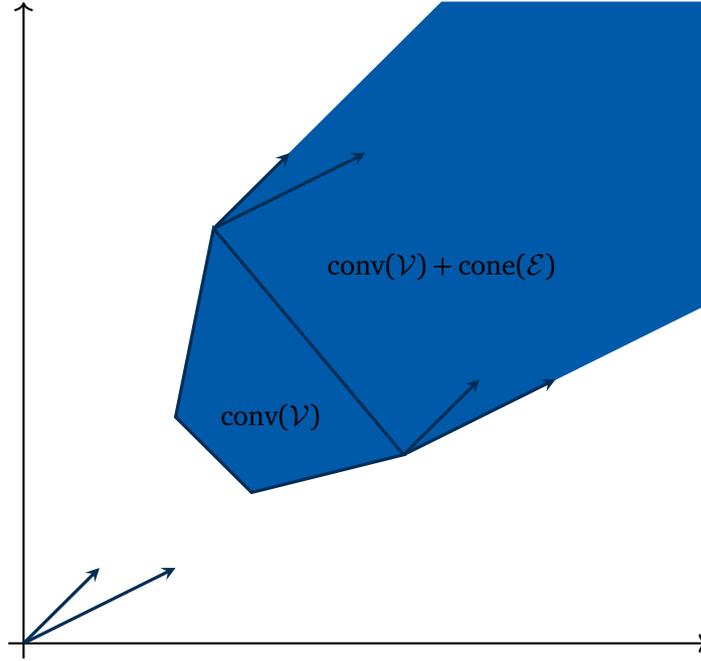


Figure 2.4: Internal structure of a polyhedron.

Remark 2.20. Note that $\text{conv}(\emptyset) = \emptyset$ (\emptyset is convex) and $\text{cone}(\emptyset) = \{\mathbf{0}\}$. However, because $\emptyset + S = \emptyset$ for every set S , we have that \mathcal{V} is non-empty if, and only if \mathcal{P} is non-empty. Furthermore, the representation is not unique. For example, the vectors in \mathcal{E} can be scaled arbitrarily, and if \mathcal{P} is a linear subspace, any $\mathbf{v} \in \mathcal{P}$ can be chosen and any basis $\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(k)}$ (with $k = \dim(\mathcal{P})$) in $\mathcal{P} = \text{conv}(\{\mathbf{v}\}) + \text{cone}(\{\mathbf{b}^{(1)}, -\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(k)}, -\mathbf{b}^{(k)}\})$.

We now know two representations for polyhedra (see Figure 2.4):

- (a) The *exterior representation* or *halfspace representation* (*H-representation*):

$$\mathcal{P} = \mathcal{P}(A, \mathbf{b}) = \bigcap_{i=1}^m \{\mathbf{x} \in \mathbb{K}^n : A_i \cdot \mathbf{x} \leq b_i\},$$

i.e., \mathcal{P} is regarded as the intersection of larger objects (halfspaces).

- (b) The *interior representation* or *vertex representation* (*V-representation*):

$$\mathcal{P} = \text{conv}(\mathcal{V}) + \text{cone}(\mathcal{E}),$$

i.e., \mathcal{P} is described by its vertices and extreme rays.

We will generally consider polyhedra as subsets of \mathbb{R}^n . The representation theorem still yields a correspondence between rational interior and exterior descriptions of polyhedra.

Corollary 2.21. Let $\mathcal{P} \subseteq \mathbb{R}^n$ and $\mathbb{K} = \mathbb{R}$. There exist $A \in \mathbb{Q}^{m \times n}$, $\mathbf{b} \in \mathbb{Q}^m$, and $m \in \mathbb{N}$ with $\mathcal{P} = \mathcal{P}(A, \mathbf{b})$ if and only if there exist finite sets $\mathcal{V}, \mathcal{E} \subseteq \mathbb{Q}^n$ with $\mathcal{P} = \text{conv}(\mathcal{V}) + \text{cone}(\mathcal{E})$.

Proof. Theorem 2.19 for $\mathbb{K} = \mathbb{Q}$ yields the existence of A and \mathbf{b} , respectively of $\mathcal{V} = \{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}\}$ and $\mathcal{E} = \{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(\ell)}\}$, such that

$$\{\mathbf{x} \in \mathbb{Q}^n : A\mathbf{x} \leq \mathbf{b}\} = \left\{ \sum_{i=1}^k \lambda_i \mathbf{v}^{(i)} + \sum_{j=1}^{\ell} \mu_j \mathbf{r}^{(j)} : \lambda \in \mathbb{Q}_{\geq 0}^k, \mu \in \mathbb{Q}_{\geq 0}^{\ell}, \mathbf{1}^{\top} \lambda = 1 \right\}.$$

We claim that the set on the left-hand-side is dense in $\mathcal{P}(A, \mathbf{b}) = \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}\}$ and the set on the right-hand-side is dense in $\text{conv}(\mathcal{V}) + \text{cone}(\mathcal{E})$. This completes the proof since the latter sets are closed in \mathbb{R}^n . The right-hand-side is dense in $\text{conv}(\mathcal{V}) + \text{cone}(\mathcal{E})$, since \mathbb{Q} is dense in \mathbb{R} and since, by linearity of the involved sums, every real point in $\text{conv}(\mathcal{V}) + \text{cone}(\mathcal{E})$ can be approached by the sum of a rational convex and a rational conic combination of the (rational) points in \mathcal{V} and \mathcal{E} . For the left-hand-side, we observe that every point $\mathbf{x} \in \mathcal{P}(A, \mathbf{b})$ lies in some polytope $\mathcal{P}' := \mathcal{P}(A, \mathbf{b}) \cap \{\mathbf{y} \in \mathbb{R}^n : \mathbb{1}\mathbf{y} \leq \beta \mathbf{1}, \mathbb{1}\mathbf{y} \geq -\beta \mathbf{1}\}$ with $\beta \in \mathbb{Q}$. Thus, \mathbf{x} can again be approached by a rational convex combination of the extreme points of \mathcal{P}' . These extreme points are solutions of systems of equations with rational coefficients (see *Introduction to Optimization*) and thus lie in \mathbb{Q}^n . Hence, $\{\mathbf{x} \in \mathbb{Q}^n : A\mathbf{x} \leq \mathbf{b}\}$ is dense in $\mathcal{P}(A, \mathbf{b})$. \square

Remark 2.22. In order to convert from the exterior to the interior representation, we can proceed as in Theorem 2.18 and obtain the matrix B as in the proofs of Theorems 2.17 and 2.5. Conversely, we can go from the interior to the exterior description directly as in Theorem 2.5. Observe that both directions use the Fourier-Motzkin elimination.

Remark 2.23. Unfortunately, when computing the vertices of a polytope with Fourier-Motzkin elimination, it can happen that in intermediate iterations, the number of generated vertices can be exponential in the dimension, even though the final result is small.

2.3 The integer hull

We now transition back to working with vector spaces over \mathbb{R} . In particular, we deal with the set of feasible solutions of a mixed-integer program of the form

$$\begin{aligned} \max \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} \leq \mathbf{b}, \\ & \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}, \end{aligned}$$

which is a subset of the polyhedron $\mathcal{P}(A, \mathbf{b})$.

Definition 2.24. The *integer hull* of a polyhedron \mathcal{P} is the set $\mathcal{P}_{1,p} := \text{conv}(\{\mathbf{x} \in \mathcal{P} : \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}\})$ for $p \in \{0, 1, \dots, n\}$. We write $\mathcal{P}_1 := \mathcal{P}_{1,n}$.

We are particularly interested in the fully integral case $p = n$.

Definition 2.25. A polyhedron P is *integral* if $P = \mathcal{P}_1$.

Intuitively, one might think that the integer hull, i.e., in particular, the convex hull of all feasible solutions of a MIP, is a polyhedron. However, this is not the case, as the following example shows.

Example 2.26. Consider the MIP

$$\sup \{x_1 - \sqrt{2}x_2 : \mathbf{x} \in \mathcal{P}_1\}$$

with $\mathcal{P} = \{\mathbf{x} \in \mathbb{R}^2 : x_1 - \sqrt{2}x_2 \leq 0, x_2 \geq 1\}$ (see Figure 2.5).

The problem is feasible since, e.g., $\begin{pmatrix} 1 \\ 1 \end{pmatrix} \in \mathcal{P} \cap \mathbb{Z}^2 \subset \text{conv}(\mathcal{P} \cap \mathbb{Z}^2) = \mathcal{P}_1$, and bounded by $x_1 - \sqrt{2}x_2 \leq 0$, but does not have an optimum as we will see. By strong duality (see *Introduction to Optimization*), this implies that the MIP above cannot be an LP, i.e., that \mathcal{P}_1 cannot be a polyhedron.

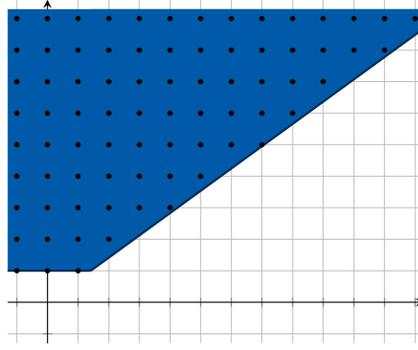


Figure 2.5: Illustration of Example 2.26.

To see that the supremum is not attained, first observe that, for points in \mathcal{P}_1 , the inequality $x_1 - \sqrt{2}x_2 \leq 0$ is equivalent to $\frac{x_1}{x_2} \leq \sqrt{2}$, because $x_2 > 0$. Since $\sqrt{2}$ is irrational, $x_1 - \sqrt{2}x_2 < 0$ holds for all points in $\mathcal{P} \cap \mathbb{Z}^2$ and thus for all points in $\mathcal{P}_1 = \text{conv}(\mathcal{P} \cap \mathbb{Z}^2)$. We construct a sequence $(\mathbf{x}^{(i)})_{i \in \mathbb{N}}$ with $\mathbf{x}^{(i)} \in \mathbb{N}^2$ and $0 < \sqrt{2}x_2^{(i)} - x_1^{(i)} < \frac{1}{i}$ to prove that the supremum is 0, which is not attained in \mathcal{P}_1 .³ Let $y_j := j\sqrt{2} - \lfloor j \cdot \sqrt{2} \rfloor \in [0, 1)$ for $j \in \{0, \dots, i\}$, and let $\{\tilde{y}_0, \dots, \tilde{y}_i\} = \{y_0, \dots, y_i\}$ be sorted such that $\tilde{y}_0 \leq \tilde{y}_1 \leq \dots$. By pigeonhole principle, there are $k > \ell$ with $0 \leq \tilde{y}_k - \tilde{y}_\ell < \frac{1}{i}$. We can now let $x_1^{(i)} := \lfloor k\sqrt{2} \rfloor - \lfloor \ell\sqrt{2} \rfloor$ and $x_2^{(i)} := k - \ell \leq 1$ to obtain $x_1^{(i)} - \sqrt{2}x_2^{(i)} = \tilde{y}_\ell - \tilde{y}_k \in (-\frac{1}{i}, 0]$, as desired. \triangle

The root (literally) of the issue above lies in the irrationality of the data. We now establish that this is the only obstruction to \mathcal{P}_1 being a polyhedron.

Definition 2.27. A polyhedron $\mathcal{P} \subseteq \mathbb{R}^n$ is called *rational* if there are $A \in \mathbb{Q}^{m \times n}$ and $\mathbf{b} \in \mathbb{Q}^m$ with $\mathcal{P} = \mathcal{P}(A, \mathbf{b})$.

We begin with the following observation.

Proposition 2.28. If $\mathcal{P}_1 \subseteq \mathbb{R}^n$ is a polyhedron, then \mathcal{P}_1 is rational.

Proof. If \mathcal{P}_1 is a polyhedron, then, according to Theorem 2.19, there are finite sets \mathcal{V}, \mathcal{E} with $\mathcal{P}_1 = \text{conv}(\mathcal{V}) + \text{cone}(\mathcal{E})$. Since \mathcal{P}_1 is the convex hull of integer points, all vertices of \mathcal{P}_1 are integral and every extreme ray contains integer points. So we can choose the sets \mathcal{V}, \mathcal{E} to consist of integer vectors. We can use Corollary 2.21 to obtain a representation by rational inequalities. \square

We collect sufficient conditions for \mathcal{P}_1 to be a polyhedron.

Proposition 2.29. Let $\mathcal{P} \subseteq \mathbb{R}^n$ be a polyhedron.

- (a) If \mathcal{P} is a polytope, then \mathcal{P}_1 is also a polytope.
- (b) If \mathcal{P} is a rational polyhedral cone, then $\mathcal{P} = \mathcal{P}_1$.

Proof.

- (a) Since \mathcal{P} is bounded, $\mathcal{X} = \{\mathbf{x} \in \mathbb{Z}^n : \mathbf{x} \in \mathcal{P}\}$ is finite. By Corollary 2.9, $\text{conv}(\mathcal{X}) = \mathcal{P}_1$ is a polytope.

³It is insufficient to require $x_1^{(i)}/x_2^{(i)} \rightarrow \sqrt{2}$, since this does not imply $x_1^{(i)} - \sqrt{2}x_2^{(i)} \rightarrow 0$, e.g., for $x_1^{(i)} = \sqrt{i} + i\sqrt{2}$ and $x_2^{(i)} = i$.

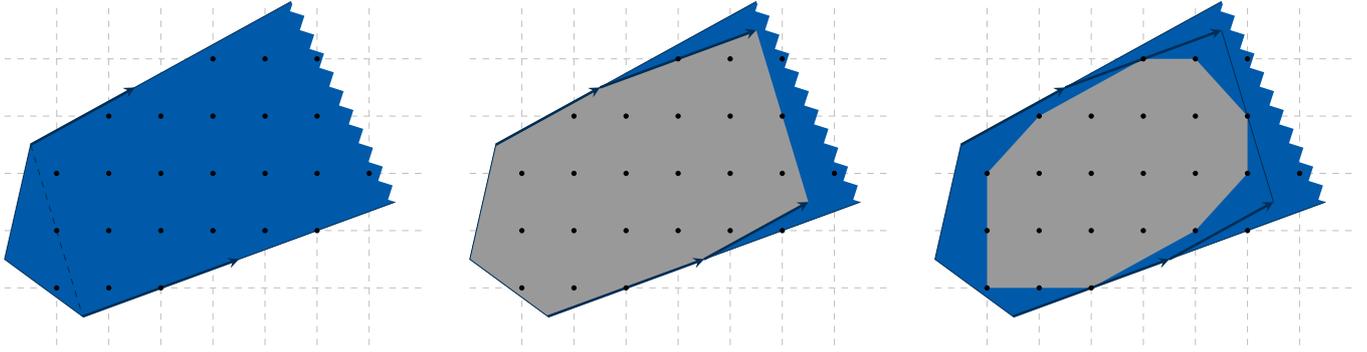


Figure 2.6: Illustration of the proof of Theorem 2.30: *left*: \mathcal{P} with $\mathbf{y}^{(1)} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$ and $\mathbf{y}^{(2)} = \begin{pmatrix} 3 \\ 1 \end{pmatrix}$, *center*: $\mathcal{Q} + \mathcal{B}$ in gray, *right*: $(\mathcal{Q} + \mathcal{B})_I$ in gray.

- (b) By Theorem 2.17 and analogously to Corollary 2.21, there exist rational $\mathbf{r}^{(i)}$ with $\mathcal{P} = \text{cone}(\{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(k)}\})$. By scaling, we can assume that $\mathbf{r}^{(i)}$ is integral for all i .

It is clear that $\mathcal{P} \supseteq \mathcal{P}_I$, since \mathcal{P} is convex and \mathcal{P}_I is the convex hull of points in \mathcal{P} . For the opposite inclusion, let $\mathbf{x} \in \mathcal{P}$. For $\mathbf{x} = \mathbf{0}$ we have $\mathbf{x} \in \mathcal{P}_I$. Otherwise, there exists a vector $\boldsymbol{\lambda} \geq \mathbf{0}$ with

$$\mathbf{x} = \sum_{i=1}^k \lambda_i \mathbf{r}^{(i)}.$$

To express \mathbf{x} as a convex combination of integral vectors, we rescale each $\mathbf{r}^{(i)}$ by the same integral factor to decrease the sum of the coefficients to below 1 and add a suitable multiple of $\mathbf{0}$. Specifically, we set $\mu := \mathbf{1}^\top \boldsymbol{\lambda} > 0$ (since $\mathbf{x} \neq \mathbf{0}$) and define $\bar{\boldsymbol{\lambda}} := \frac{1}{\lceil \mu \rceil} \boldsymbol{\lambda}$ to write

$$\mathbf{x} = \sum_{i=1}^k \bar{\lambda}_i (\lceil \mu \rceil \mathbf{r}^{(i)}) + (1 - \mathbf{1}^\top \bar{\boldsymbol{\lambda}}) \mathbf{0}$$

with $\sum_{i=1}^k \bar{\lambda}_i + (1 - \mathbf{1}^\top \bar{\boldsymbol{\lambda}}) = 1$. Because the vectors $\mathbf{0}, \lceil \mu \rceil \mathbf{r}^{(1)}, \dots, \lceil \mu \rceil \mathbf{r}^{(k)}$ are integral and contained in \mathcal{P} , we have that \mathbf{x} is a convex combination of integer points in \mathcal{P} , i.e., $\mathbf{x} \in \mathcal{P}_I$. \square

We finally conclude that it is sufficient to require the system $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ to have rational coefficients in order to ensure that the convex hull of all integral solutions to form a polyhedron.

Theorem 2.30. If $\mathcal{P} \subseteq \mathbb{R}^n$ is a rational polyhedron, then \mathcal{P}_I is also a (rational) polyhedron.

Proof. Once we establish that \mathcal{P}_I is a polyhedron, rationality follows by Proposition 2.28.

By Corollary 2.21, we know that \mathcal{P} has a representation of the form $\mathcal{P} = \mathcal{Q} + \mathcal{C}$, where \mathcal{Q} is a convex hull of rational vectors and \mathcal{C} is the conic hull of rational vectors. By scaling, we can thus determine $\mathbf{y}^{(i)} \in \mathbb{Z}^n$, $i \in \{1, \dots, k\}$, with $\mathcal{C} = \text{cone}(\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(k)}\})$. We set

$$\mathcal{B} := \left\{ \sum_{i=1}^k \mu_i \mathbf{y}^{(i)} : \mathbf{0} \leq \boldsymbol{\mu} \leq \mathbf{1} \right\}.$$

Note that the (bounded) set \mathcal{B} is a polytope by Corollary 2.8. By Corollaries 2.6 and 2.8, $\mathcal{Q} + \mathcal{B}$ is a polytope, since it is bounded. According to Proposition 2.29 (a), $(\mathcal{Q} + \mathcal{B})_1$ is therefore a polytope as well. Finally, by Theorem 2.17 and Corollary 2.8, it follows that $(\mathcal{Q} + \mathcal{B})_1 + \mathcal{C}$ is a polyhedron.

With this, it suffices to show

$$\mathcal{P}_1 = (\mathcal{Q} + \mathcal{B})_1 + \mathcal{C}.$$

\subseteq : We argued that $(\mathcal{Q} + \mathcal{B})_1 + \mathcal{C}$ is a polyhedron, i.e., in particular, it is convex. Since \mathcal{P}_1 is the convex hull of integer points, it suffices to show that all integer points in \mathcal{P}_1 also lie in $(\mathcal{Q} + \mathcal{B})_1 + \mathcal{C}$. So let $\mathbf{p} \in \mathcal{P}_1$ be integral. Then there exist $\mathbf{q} \in \mathcal{Q}$ and $\mathbf{c} \in \mathcal{C}$ with $\mathbf{p} = \mathbf{q} + \mathbf{c}$. We can find $\boldsymbol{\mu} \geq \mathbf{0}$ with

$$\mathbf{c} = \sum_{i=1}^k \mu_i \mathbf{y}^{(i)}.$$

Now let

$$\mathbf{c}' := \sum_{i=1}^k \lfloor \mu_i \rfloor \mathbf{y}^{(i)}$$

and $\mathbf{b} := \mathbf{c} - \mathbf{c}'$. It follows that $\mathbf{c}' \in \mathcal{C} \cap \mathbb{Z}^n$ and $\mathbf{q} + \mathbf{b} = \mathbf{p} - \mathbf{c}' \in \mathbb{Z}^n$, since \mathbf{p} and \mathbf{c}' are integral. Furthermore, $\mathbf{b} \in \mathcal{B}$ by definition and therefore $\mathbf{q} + \mathbf{b} \in (\mathcal{Q} + \mathcal{B})_1$. Hence, $\mathbf{p} = \mathbf{q} + \mathbf{b} + \mathbf{c}' \in (\mathcal{Q} + \mathcal{B})_1 + \mathcal{C}$.

\supseteq : Because of $\mathcal{B} \subseteq \mathcal{C}$ we have

$$(\mathcal{Q} + \mathcal{B})_1 + \mathcal{C} \subseteq (\mathcal{Q} + \mathcal{C})_1 + \mathcal{C} = \mathcal{P}_1 + \mathcal{C},$$

and, by Proposition 2.29 (b),

$$\begin{aligned} \mathcal{P}_1 + \mathcal{C} &= \mathcal{P}_1 + \mathcal{C}_1 \\ &\subseteq (\mathcal{P} + \mathcal{C})_1 = \mathcal{P}_1, \end{aligned}$$

which establishes $(\mathcal{Q} + \mathcal{B})_1 + \mathcal{C} \subseteq \mathcal{P}_1$. □

Corollary 2.31. Let $\mathcal{P} \subseteq \mathbb{R}^n$ be a rational polyhedron such that $\mathcal{P}_1 \neq \emptyset$. Then, for all $\mathbf{c} \in \mathbb{R}^n$, $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \mathcal{P}\}$ is bounded if and only if $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \mathcal{P}_1\}$ is bounded.

Proof. In the proof of Theorem 2.30, we have $\mathcal{P} = \mathcal{Q} + \mathcal{C}$ and $\mathcal{P}_1 = (\mathcal{Q} + \mathcal{B})_1 + \mathcal{C}$, where \mathcal{C} is the only unbounded term in either equality. In addition, $\mathcal{P}_1 \neq \emptyset$ (and thus $\mathcal{P} \neq \emptyset$) implies that $(\mathcal{B} + \mathcal{Q})_1 \neq \emptyset$ and $\mathcal{Q} \neq \emptyset$. Hence, both $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \mathcal{P}\}$ and $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \mathcal{P}_1\}$ are unbounded if and only if there exists $\mathbf{r} \in \mathcal{C}$ with $\mathbf{c}^\top \mathbf{r} > 0$. □

With analogous arguments, Theorem 2.30 can be extended to the mixed-integer case.

Theorem 2.32. If \mathcal{P} is a rational polyhedron, then $\mathcal{P}_{1,p}$ is a rational polyhedron for all $p \in \{0, \dots, n\}$.

We conclude the chapter with a useful characterization.

Lemma 2.33. Let \mathcal{P} be a rational polyhedron. Then the following are equivalent:

- (a) $\mathcal{P} = \mathcal{P}_1$.
- (b) Every non-empty face of \mathcal{P} contains an integer point.
- (c) Every non-empty minimal face of \mathcal{P} , i.e., every face containing no other faces, contains an integer point.
- (d) $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \mathcal{P}\}$ is assumed at an integer point, for all $\mathbf{c} \in \mathbb{R}^n$ (or $\mathbf{c} \in \mathbb{Z}^n$), for which the maximum is finite.

Proof. exercise. □

In particular, a MIP reduces to an LP if the underlying polyhedron is integral.

Corollary 2.34. If $\mathcal{P}(A, \mathbf{b})$ is integral, then the optimum vertex solutions of the following problems coincide

$$\begin{array}{ll} \max & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} & A\mathbf{x} \leq \mathbf{b}, \\ & \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}, \end{array} \qquad \begin{array}{ll} \max & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} & A\mathbf{x} \leq \mathbf{b}, \\ & \mathbf{x} \in \mathbb{R}^n. \end{array}$$

3 Branch-and-Bound Method

In this chapter we deal with the branch-and-bound method, an exact method for solving mixed-integer programs. The method exhaustively searches the solution space and thus guarantees that an optimum solution is found if one exists. In case the underlying problem is bounded, the method terminates, but with exponential running time. In practice, the branch-and-bound method is very often efficient and is the basis of all modern integer optimization solvers. We begin by showing that, from a theoretical perspective, solving mixed-integer programs is computationally difficult.

3.1 Complexity of integer programs

We already know (Proposition 2.29 (a) and Theorem 2.30) that every bounded or rational MIP can be reduced to a linear program, i.e., there exist $D \in \mathbb{Q}^{m \times n}$ and $\mathbf{d} \in \mathbb{Q}^m$, with $\mathcal{P}_1 = \mathcal{P}(D, \mathbf{d})$. Here, m can become very large. We now show that, in general, we cannot hope to solve a rational MIP efficiently. More precisely, we show the following.

Theorem 3.1. The decision problem “Has a given rational system of inequalities $A\mathbf{x} \leq \mathbf{b}$ an integral solution?” is NP-complete (even if \mathbf{x} is binary).

We begin by recalling central definitions (see *Algorithmic Discrete Mathematics*).

Definition 3.2. A decision problem is a tuple $(\mathcal{I}, (S_I)_{I \in \mathcal{I}})$ with $S_I \in \{\{\text{‘yes’}\}, \{\text{‘no’}\}\}$ for all $I \in \mathcal{I}$. The input size $|I|$ of an instance $I \in \mathcal{I}$ is defined to be equal to the number of bits in the binary representation of \mathcal{I} . An algorithm has polynomial running time if it computes $A(I) = S_I$ in time $|I|^{O(1)}$ for all $I \in \mathcal{I}$.

Roughly speaking, NP is the class of decision problems for which ‘yes’-instances can be verified in polynomial time. For a rigorous definition we refer to *Algorithmic Discrete Mathematics*. Based on this class, we can capture the computational intractability as follows.

Definition 3.3. A decision problem $\Pi = (\mathcal{I}, (S_I)_{I \in \mathcal{I}})$ is NP-hard if, for all $\Pi' = (\mathcal{I}', (S'_{I'})_{I' \in \mathcal{I}'}) \in \text{NP}$, there is a reduction $R: \mathcal{I}' \rightarrow \mathcal{I}$ that satisfies $S'_{I'} = S_{R(I')}$ and can be computed in polynomial time. If, additionally, $\Pi \in \text{NP}$, then Π is NP-complete.

One of the most important NP-complete problems is the satisfiability problem (SAT).

Satisfiability (SAT) Problem

input: CNF formula given by clauses $\mathcal{C} = \{C_i\}_{i=1, \dots, m}$ over variables $\mathcal{X} = \{x_i\}_{i=1, \dots, n}$
problem: Is there a satisfying assignment $\alpha: \mathcal{X} \rightarrow \{0, 1\}$?

An example of a CNF formula is $(x_1 \vee \bar{x}_2 \vee x_4) \wedge (x_2 \vee \bar{x}_3) \wedge (x_3 \vee \bar{x}_4)$ with satisfying assignment, e.g., $\alpha(x_1) = \alpha(x_2) = 1$ and $\alpha(x_3) = \alpha(x_4) = 0$.

We need the following observation about the *coding length* of a solution, i.e., the length of a binary representation (for a proof see [39, Chapter 17]).

Lemma 3.4. Let $\mathcal{P} = \mathcal{P}(A, \mathbf{b})$ be a rational polyhedron such that the coding length of each inequality is at most $\varphi \in \mathbb{N}$. If $\mathcal{P}_1 \neq \emptyset$, then there is an integer point in \mathcal{P} , whose coding length is at most $6n^3\varphi$.

Proof of Theorem 3.1. From Lemma 3.4 we obtain a polynomial certificate for $\mathcal{P} \neq \emptyset$. It follows that the decision problem is in NP.

We can model a given SAT instance with variables x_1, \dots, x_n and clauses C_1, \dots, C_m as an integer system $A\hat{\mathbf{x}} \leq \mathbf{b}$ consisting of the (negated) inequalities

$$\sum_{j: x_j \in C_i} \hat{x}_j + \sum_{j: \bar{x}_j \in C_i} (1 - \hat{x}_j) \geq 1 \quad \forall i \in \{1, \dots, m\},$$

with $\hat{\mathbf{x}} \in \{0, 1\}^n$. An integral solution of $A\hat{\mathbf{x}} \leq \mathbf{b}$ can then be interpreted as

$$\hat{x}_j = \begin{cases} 1, & \text{if variable } x_j \text{ is set to "true",} \\ 0, & \text{if variable } x_j \text{ is set to "false".} \end{cases}$$

Now, $A\hat{\mathbf{x}} \leq \mathbf{b}$ has an integral solution (i.e. $\{\hat{\mathbf{x}} \in \{0, 1\}^n : A\hat{\mathbf{x}} \leq \mathbf{b}\} \neq \emptyset$) if and only if the given SAT instance has a solution. The transformation can be carried out in polynomial time and $\hat{\mathbf{x}}$ is obviously bounded. \square

The following is an immediate consequence of Theorem 3.1.

Corollary 3.5. Every problem in NP can be modeled as a (binary) integer program.

Since finding a feasible solution to a MIP is NP-hard and thus solving a MIP is difficult, but LPs can be solved in polynomial time (see *Introduction to Optimization*), we cannot expect to efficiently find a complete description of \mathcal{P}_1 or to solve the associated separation problem (unless $P = NP$). With this in mind, we now introduce the branch-and-bound method for solving general mixed-integer programs. In Chapter 4, we will deal with cases that allow to find solutions efficiently.

3.2 The branch-and-bound method

The *branch-and-bound* method revolves around the two operations of *branching* and *bounding*. We want to solve the mixed-integer optimization problem

$$\begin{aligned} \max \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} \leq \mathbf{b} \\ & \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}, \end{aligned} \tag{3.1}$$

where $A \in \mathbb{Q}^{m \times n}$, $\mathbf{c} \in \mathbb{Q}^n$, $\mathbf{b} \in \mathbb{Q}^m$, and $p \in \{1, \dots, n\}$. Note that we assume that all data are rational, because this is the case in practice and because the irrational case causes theoretical problems, as we saw in Section 2.3.

We use a “relaxation” that yields an upper bound on the optimum value of the problem (3.1) by forgoing the integer conditions. We obtain the *LP relaxation*

$$\begin{aligned} \max \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} \leq \mathbf{b} \\ & \mathbf{x} \in \mathbb{R}^n. \end{aligned} \tag{3.2}$$

Such linear programs can be solved efficiently (see *Introduction to Optimization*). If the computed optimum solution of the LP relaxation is integral, we have solved (3.1). Otherwise, we branch, i.e., we generate subproblems, which we solve recursively. In this way, we generate a collection of subproblems with shrinking feasible regions. Every optimum solution to our integer problem appears in one of the subproblems, so that no optimum solutions are lost on the way.

The LP relaxation can be further enhanced by inequalities (see Chapter 5) and other relaxations can be used (see Chapter 6). Furthermore, the procedure can be accelerated by heuristics (see Chapter 7).

3.2.1 Divide and conquer

As a first step, we solve the LP relaxation (3.2) to obtain a fractional optimum solution $\hat{\mathbf{x}}$. If we are lucky,

$$\hat{\mathbf{x}} \in \mathcal{X} := \{\mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p} : A\mathbf{x} \leq \mathbf{b}\},$$

i.e., $\hat{\mathbf{x}}$ is feasible for (3.1), and we have already found an optimum solution to the mixed-integer problem. Sufficient conditions that guarantee that this case occurs will be discussed in Chapter 4.

In general, the solution $\hat{\mathbf{x}}$ of the LP relaxation is not integral (i.e., $\hat{\mathbf{x}} \notin \mathbb{Z}^p \times \mathbb{R}^{n-p}$). In this case, there exists an index $i \in \{1, \dots, p\}$ such that $\hat{x}_i \notin \mathbb{Z}$. We now create two subproblems

$$\mathcal{X}^\leq := \{\mathbf{x} \in \mathcal{X} : x_i \leq \lfloor \hat{x}_i \rfloor\} \subset \mathcal{X} \quad \text{and} \quad \mathcal{X}^\geq := \{\mathbf{x} \in \mathcal{X} : x_i \geq \lceil \hat{x}_i \rceil\} \subset \mathcal{X}.$$

Evidently, $\mathcal{X} = \mathcal{X}^\leq \cup \mathcal{X}^\geq$ and $\mathcal{X}^\leq \cap \mathcal{X}^\geq = \emptyset$, so that every (optimum) solution of the MIP (3.1) must either lie in \mathcal{X}^\leq or in \mathcal{X}^\geq . These sets are associated with the two polyhedra

$$\mathcal{P}^\leq := \{\mathbf{x} \in \mathcal{P} : x_i \leq \lfloor \hat{x}_i \rfloor\} \quad \text{and} \quad \mathcal{P}^\geq := \{\mathbf{x} \in \mathcal{P} : x_i \geq \lceil \hat{x}_i \rceil\},$$

where $\mathcal{P} := \mathcal{P}(A, \mathbf{b}) = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$. Because $\hat{\mathbf{x}} \notin \mathcal{P}^\leq \cup \mathcal{P}^\geq$, the LP relaxations of both subproblems yield optimum solutions different from $\hat{\mathbf{x}}$. We can now apply the procedure recursively to the two subproblems.

The recursive subdivision into subproblems creates a *branch-and-bound tree* whose nodes correspond to the subproblems, see Figure 3.1. In each step, a subproblem that has not yet been processed is selected and an optimum solution $\hat{\mathbf{x}}$ of the corresponding LP relaxation is computed; if it is infeasible, we continue with another subproblem. If $\hat{\mathbf{x}} \notin \mathcal{X}$, we create two new subproblems. If $\hat{\mathbf{x}} \in \mathcal{X}$, we have found a feasible solution for the MIP. In this case we keep track of the best solution \mathbf{x}^* we encountered so far and continue with another subproblem. The procedure terminates when all subproblems have been processed. The optimum solution then is \mathbf{x}^* if any solution was found, otherwise the problem is infeasible.

In practice, the procedure can be accelerated using the following simple insight. Let us assume that we already have found some feasible solution $\mathbf{x}^* \in \mathcal{X}$. Then, obviously,

$$\max \{ \mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \mathcal{X} \} \geq \mathbf{c}^\top \mathbf{x}^*.$$

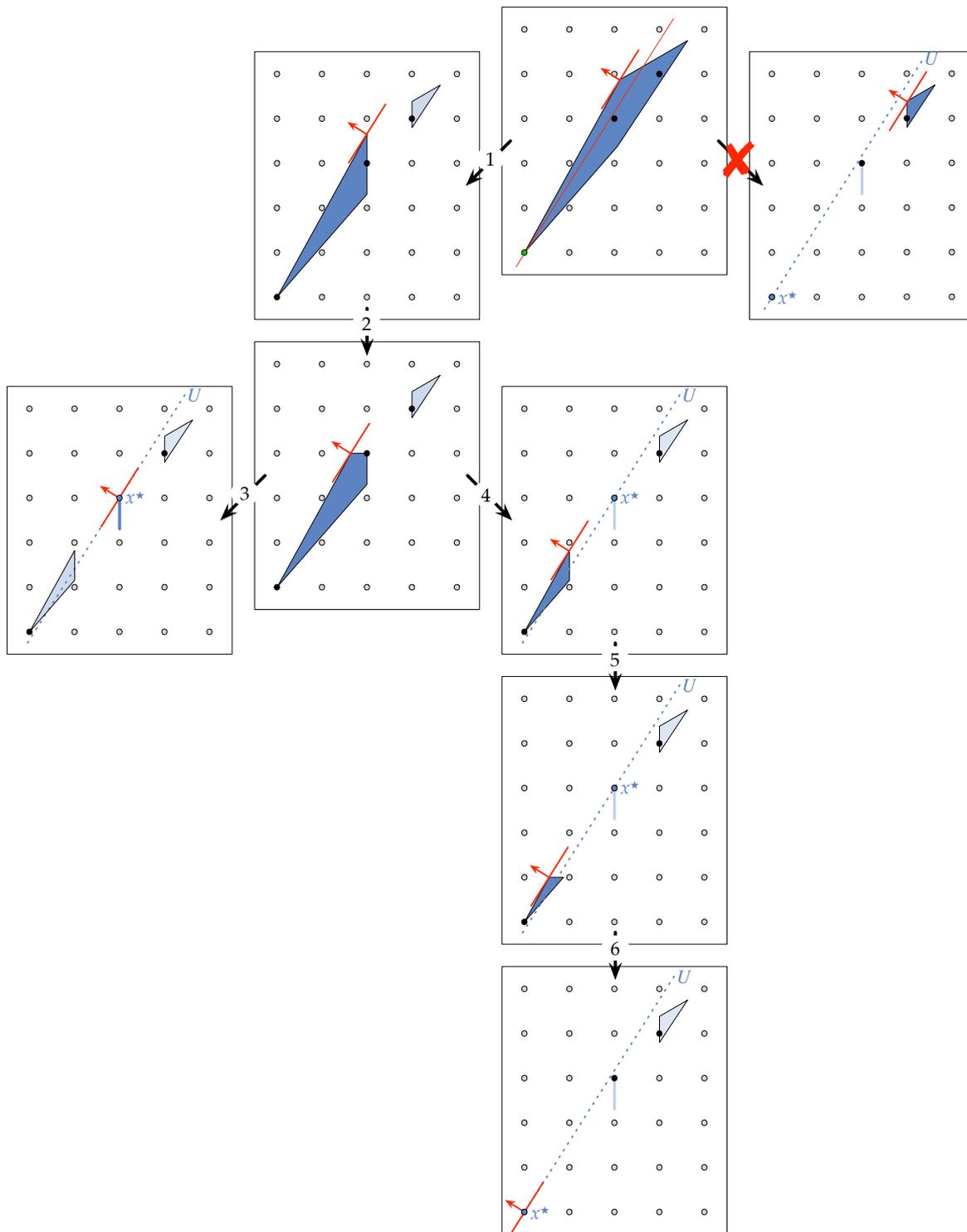


Figure 3.1: Example of an application of the branch-and-bound method: The tree is searched from left to right. In the lower parts, the respective other subproblems are infeasible and omitted.

The solution \mathbf{x}^* therefore provides us with a lower bound on the optimum objective function value of our MIP. If we now solve a subproblem with associated polyhedron $\hat{\mathcal{P}}$ satisfying

$$\max \{ \mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \hat{\mathcal{P}} \} \leq \mathbf{c}^\top \mathbf{x}^*,$$

then there cannot be any better feasible solution than \mathbf{x}^* in this subproblem, and this subproblem need not be further considered. Accordingly, we may prune (i.e., “cut off”) the corresponding subtree.

3.2.2 Formal definition

The above procedure results in the *branch-and-bound method*, which is formally stated below.

Algorithm: branch-and-bound method

input: MIP $\max \{ \mathbf{c}^\top \mathbf{x} : \mathbf{A}\mathbf{x} \leq \mathbf{b}, \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p} \}$ with $\mathcal{P}(\mathbf{A}, \mathbf{b})$ bounded

output: optimum solution \mathbf{x}^* or “infeasible” if $\underline{z} = -\infty$

```

A ← {P(A, b)}                                     (active nodes)
z̄ ← -∞, x̄ ← ∅                                     (lower bound, best solution)
while A ≠ ∅:
  choose P̂ ∈ A                                     (subproblem selection)
  if P̂ ≠ ∅:
    compute x̃ ∈ P̂ ∩ (Zp × Rn-p) via primal heuristics (see Chapter 7)
    if c⊤x̃ > z̄:
      z̄ ← c⊤x̃, x̄ ← x̃
    x̂ ← arg max{c⊤x : x ∈ P̂}                       (or use Chapter 6)
    if c⊤x̂ > z̄:
      if x̂ ∈ Zp × Rn-p:
        z̄ ← c⊤x̂, x̄ ← x̂                             (see Chapter 4)
      else
        choose i ∈ {1, ..., p} with x̂i ∉ Z           (variable selection)
        P̂≤ ← {x ∈ P̂ : xi ≤ ⌊x̂i⌋}, P̂≥ ← {x ∈ P̂ : xi ≥ ⌈x̂i⌉} (or use Chapter 5)
        A ← A ∪ {P̂≤, P̂≥}
    A ← A \ {P̂}
return (x̄, z̄)

```

Remark 3.6. The grayed out part of the algorithm is optional and intended for acceleration – it is not necessary for correct execution. Techniques from Chapter 7 can be used here.

The key observation for correctness of branch-and-bound is that the algorithm maintains the following invariants.

Lemma 3.7. Branch-and-bound maintains the invariants that

- (a) if $\underline{z} > -\infty$, then $\bar{\mathbf{x}}$ is feasible for (3.1) with $\mathbf{c}^\top \bar{\mathbf{x}} = \underline{z}$, and
- (b) for every optimum solution \mathbf{x}^* of (3.1), either $\underline{z} = \mathbf{c}^\top \mathbf{x}^*$, or there is an active node $\hat{\mathcal{P}} \in A$ with $\mathbf{x}^* \in \hat{\mathcal{P}}$.

Proof. The first part of the invariant follows since \underline{z} and \bar{x} are only updated in unison and only when \bar{x} is a feasible solution of value \underline{z} .

For the second part, observe that only fractional solutions are eliminated when branching. This means that every optimum solution \mathbf{x}^* remains in some active node $\hat{\mathcal{P}}$ until this node can be solved without branching, which only happens when the corresponding LP relaxation yields an integral optimum solution. This solution cannot be better than \mathbf{x}^* since the integral solutions in $\hat{\mathcal{P}}$ are a subset of \mathcal{X} , and it cannot be worse than \mathbf{x}^* , since $\mathbf{x}^* \in \hat{\mathcal{P}}$. By the first part of the invariant and optimality of \mathbf{x}^* , we have $\underline{z} \leq \mathbf{c}^\top \mathbf{x}^*$. Hence, after considering $\hat{\mathcal{P}}$, we have $\underline{z} = \mathbf{c}^\top \mathbf{x}^*$ (with or without update). \square

This immediately yields the following bounds.

Corollary 3.8. The branch-and-bound method maintains the invariant that \underline{z} is a lower bound and

$$\bar{z} := \max \{ \underline{z}, \max \{ \mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \hat{\mathcal{P}}, \hat{\mathcal{P}} \in \mathcal{A} \} \}$$

is an upper bound on the optimum value of (3.1).

We show that the branch-and-bound method computes the correct result in finite time.

Theorem 3.9. If \mathcal{P} is bounded, then the branch-and-bound method terminates after a finite number of steps and the final value of \underline{z} is equal to the optimum value of the MIP (3.1).

Proof. Assume that branch-and-bound terminates. For the case that (3.1) is infeasible, Lemma 3.7 (a) implies that the correct value $\underline{z} = -\infty$ is maintained. If (3.1) is feasible, Lemma 3.7 (b) implies that branch-and-bound can only terminate once $\underline{z} = \mathbf{c}^\top \mathbf{x}^*$ for every optimum solution \mathbf{x}^* and that, at this point, \bar{x} must be optimal. We now show that the algorithm terminates. Since $\mathcal{P} = \mathcal{P}(\mathbf{A}, \mathbf{b})$ is bounded, there is a constant $\varphi \in \mathbb{N}$ so that

$$\mathcal{P} \subseteq \{ \mathbf{x} \in \mathbb{R}^n : -\varphi \mathbf{1} \leq \mathbf{x} \leq \varphi \mathbf{1} \}.$$

Every LP relaxation that occurs during branch-and-bound is defined by $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ and constraints on the variables of the form

$$x_i \geq \xi \quad \text{with } \xi \in \{-\varphi + 1, \dots, \varphi\} \quad \text{or} \quad (3.3)$$

$$x_i \leq \xi \quad \text{with } \xi \in \{-\varphi, \dots, \varphi - 1\} \quad (3.4)$$

for $i \in \{1, \dots, p\}$. Because $\hat{\mathbf{x}}$ is not feasible for either of the two subproblems generated when branching, and since all active subproblems have disjoint feasible regions, all LP relaxations encountered by the algorithm are different. This means that the LP relaxations use different subsets of inequalities of the type in (3.3) and (3.4). However, there are only $4p\varphi$ different inequalities, i.e., $2^{4p\varphi}$ different combinations, and thus subproblems, to consider. \square

Remark 3.10. For binary programs, the height of the branch-and-bound tree is bounded by the dimension n , hence it has at most 2^n leaves and thus there are at most $2^{n+1} - 1$ subproblems to consider.

Remark 3.11. If \mathcal{P} is unbounded, the method may not terminate (*exercise*). However, an estimation of the absolute values of the components of an optimum solution, as in Lemma 3.4, can be used to constrain the problem without losing all optimum solutions.

Thus, we can say that the branch-and-bound method terminates. However, it may encounter exponentially many subproblems, even for infeasible binary programs (*exercise*). This performance is not surprising, since it is NP-hard to find a feasible solution (Theorem 3.1).

3.2.3 Practical considerations

Dual simplex

A key insight for an efficient implementation is that the dual simplex algorithm is ideally suited to solve the LP relaxations. This is because when a subproblem is subdivided and an additional constraint on the variable x_i is imposed, the current basis remains dually feasible (not primally, of course). The same is true for the addition of further inequalities in the branch-and-cut method. The dual simplex can resume directly with the optimum basis of the parent node. This is particularly efficient when moving on to a child in subproblem selection, but is still useful otherwise if the information about the final basis every parent node is stored in memory.

Selection rules

So far, we have discussed a generic version of the branch-and-bound method. Two degrees of freedom remain: The selection of the next active subproblem and the selection of the variable $\hat{x}_i \notin \mathbb{Z}$ to branch on, i.e., the variable whose domain is subdivided.

Subproblem selection When selecting the next subproblem, there are, e.g., the following options.

- With *depth-first search*, the active node created most recently is selected next, i.e., in particular, one of the children of the previously considered node if possible. Depth-first search has the advantage that consecutive subproblems differ only slightly from one another and thus the required modification is small. Furthermore, the chance of finding an integral solution (and possibly improving the bound \underline{z}) increases with increasing depth, as long as the problems are feasible and are not pruned.
- With *breadth-first search*, the active node created least recently is processed next, i.e., a node of smallest depth in the tree. This method systematically works through the tree level by level. Under the assumption that the local upper bounds $\bar{z}(\hat{P})$ of the nodes \hat{P} decrease with increasing depths, this method primarily reduces the upper bound \bar{z} of Lemma 3.8.
- With *best-bound search*, the active node \hat{P} of maximum local upper bound $\bar{z}(\hat{P}) := \max\{\mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \hat{P}\}$ is selected next. The idea is to attempt to reduce \bar{z} as quickly as possible. The disadvantage of this variant (and of breadth-first search) lies in the fact that feasible integral solutions are found relatively late, i.e., that the lower bound \underline{z} has no effect in the beginning.

In practice, best-bound search is typically used, but a few depth-first steps are carried out from time to time to find feasible integral solutions a bit sooner.

Variable selection – branching rule For the selection of the variable $\hat{x}_i \notin \mathbb{Z}$, there are various options, including the following.

- Choose $i \in \arg \max \{ \min\{\hat{x}_j - \lfloor \hat{x}_j \rfloor, \lceil \hat{x}_j \rceil - \hat{x}_j\}, j \in \{1, \dots, p\} \}$, i.e., the variable of largest distance to the next integer, or a maximally infeasible variable with respect to integrality. The hope is that the values of such variables must change considerably in both subproblems and therefore cause large changes in the objective function value.
- Choose $i = \arg \min \{ \min\{\hat{x}_j - \lfloor \hat{x}_j \rfloor, \lceil \hat{x}_j \rceil - \hat{x}_j\}, \hat{x}_j \notin \mathbb{Z}, j \in \{1, \dots, p\} \}$, i.e., the variable of smallest distance to the next integer, or a minimally infeasible variable with respect to integrality. The hope is that these variables only need little change to become integral.

- *Pseudo-cost branching* uses historical information about the change in the objective function when a certain variable is selected as the branching variable. If, in the past, branches with variable j changed the objective function value by Δ_j^- or Δ_j^+ in the children whose subproblems were rounded down, respectively up, we set

$$\zeta_j^- := \frac{\Delta_j^-}{\hat{x}_j - \lfloor \hat{x}_j \rfloor} \quad \text{and} \quad \zeta_j^+ := \frac{\Delta_j^+}{\lceil \hat{x}_j \rceil - \hat{x}_j}.$$

Then, we calculate the arithmetic means $\bar{\zeta}_j^-$ and $\bar{\zeta}_j^+$ of these values, which provide an estimator for the change of the objective function value per unit, similarly to reduced costs in the simplex method (see *Introduction to Optimization*). In the current node, we select a variable j of largest valuation $B((\hat{x}_j - \lfloor \hat{x}_j \rfloor) \cdot \bar{\zeta}_j^-, (\lceil \hat{x}_j \rceil - \hat{x}_j) \cdot \bar{\zeta}_j^+)$. The valuation function could, for example, be $B(q^-, q^+) = (1 - \mu) \cdot \min\{q^-, q^+\} + \mu \cdot \max\{q^-, q^+\}$ (where $\mu \in [0, 1]$). An alternative is $B(q^-, q^+) = \max\{q^-, \varepsilon\} \cdot \max\{q^+, \varepsilon\}$ (with $\varepsilon > 0$).

This branching rule works particularly well if a lot of information about the variables is already available, especially if there has been a lot of branching. In the beginning, the selection is essentially random.

- *Strong branching* generates the two subproblems for all possible candidates $i \in \{1, \dots, p\}$ and calculates the corresponding changes in objective function value. It then selects the variable that causes the largest change. This procedure is very time-consuming, but almost always yields small trees. In practice, it is often accelerated by limiting the number of simplex iterations and the number of candidates.
- *Reliability branching* combines pseudo-cost and strong branching. In the beginning, the pseudo-costs are initialized via strong-branching. Later, when the values are more reliable, the pseudo costs are used.

For details on the last three rules, we refer to the literature (see, e.g., [1]). It turns out that the first two rules – although intuitively good – perform very poorly in practice: They are about as good as a random selection rule. The last rule, on the other hand, yields significantly better results.

Presolving

Presolving is essential for the practical performance of the branch-and-bound method. It tries to simplify the MIP with the aim of making the resulting problem easier. Typically, presolving focuses on reducing the number of variables and constraints. We list examples of presolving steps:

- Eliminate trivial constraints, e.g., $\mathbf{0}^\top \mathbf{x} \leq 1$.
- Eliminate duplicate constraints.
- Eliminate redundant constraints, e.g., $2x_1 + x_2 \leq 3$; $x_1, x_2 \in \{0, 1\}$.
- Simplify constraints, e.g., $2x_1 + 2x_2 = 2 \rightarrow x_1 + x_2 = 1$.
- Strengthen coefficients, e.g., $2x_1 + x_2 \geq 1$; $x_1, x_2 \in \{0, 1\} \rightarrow x_1 + x_2 \geq 1$.
- Aggregate variables: Consider $\alpha_1 x_1 + \dots + \alpha_k x_k = \beta$ with $\alpha_i \neq 0$. Then the equation can be solved for x_i and be eliminated together with x_i . The objective function must be adapted accordingly. Note that aggregation can make the coefficient matrix of a MIP denser, i.e., introduce more non-zero elements. This can have a negative impact on the the solving speed for the LP relaxations.
- Recognize special types of constraints, e.g., knapsack inequalities, set-packing constraints etc. For these, we can then, for example, generate the inequalities of Section 5.
- Presolve dual: E.g., if the only constraints involving x_2 are $2x_1 + x_2 \geq 3$ and $x_2 \leq 5$, the objective function coefficient of x_2 is non-negative, and $x_1 \geq 0$, we can fix $x_2 = 5$ and eliminate the first inequality.

- Use probing: Fix binary variables to 0 or 1 and test for consequences, e.g., for $x_1 + x_2 \leq 1$, $x_1 - x_2 \leq 0$, and $x_1, x_2 \in \{0, 1\}$, trying $x_1 = 1$ yields $x_2 = 0$, which yields $x_1 = 0$ – a contradiction. So we can safely set $x_1 = 0$. Similarly, implications such as, e.g., $x_i = 1 \Rightarrow x_j = 0$ can be recognized. These implications then lead to set-packing conditions (e.g., $x_i + x_j \leq 1$).

Such presolving steps can dramatically decrease the size of MIPs, especially those generated by modeling languages.

Fixing reduced costs

During the branch-and-bound procedure, it is possible to use the reduced costs of the solution to the LP relaxation. To simplify notation, let us consider the problem

$$\max \{ \mathbf{c}^\top \mathbf{x} : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}, \mathbf{x} \in \mathbb{Z}^n \}.$$

Let $\hat{\mathbf{x}}$ be an optimum basis solution of the relaxation $\max \{ \mathbf{c}^\top \mathbf{x} : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0} \}$ with associated basis B and non-basis $N = \{1, \dots, n\} \setminus B$, i.e., $\hat{\mathbf{x}}_B = A_{.B}^{-1} \mathbf{b}$, $\hat{\mathbf{x}}_N = \mathbf{0}$. Let \underline{z} be a lower bound on the optimum value. For every feasible integral solution \mathbf{x} with objective function value at least \underline{z} , we have

$$\begin{aligned} \underline{z} - \mathbf{c}^\top \hat{\mathbf{x}} &\leq \mathbf{c}^\top \mathbf{x} - \mathbf{c}^\top \hat{\mathbf{x}} = \mathbf{c}_N^\top (\mathbf{x}_N - \hat{\mathbf{x}}_N) + \mathbf{c}_B^\top (\mathbf{x}_B - \hat{\mathbf{x}}_B) \\ &= \mathbf{c}_N^\top \mathbf{x}_N + \mathbf{c}_B^\top (A_{.B}^{-1} \mathbf{b} - A_{.B}^{-1} A_{.N} \mathbf{x}_N - A_{.B}^{-1} \mathbf{b}) \\ &= \mathbf{c}_N^\top \mathbf{x}_N - \mathbf{c}_B^\top A_{.B}^{-1} A_{.N} \mathbf{x}_N \\ &= (\mathbf{c}_N^\top - \mathbf{c}_B^\top A_{.B}^{-1} A_{.N}) \mathbf{x}_N = \mathbf{z}_N^\top \mathbf{x}_N \leq z_j x_j, \end{aligned}$$

for all $j \in N$. For the reduced costs, we have $\mathbf{z}_N \leq \mathbf{0}$, because the basis B is optimal. If now $z_j < 0$, then

$$x_j \leq \frac{\underline{z} - \mathbf{c}^\top \hat{\mathbf{x}}}{z_j} = \frac{\mathbf{c}^\top \hat{\mathbf{x}} - \underline{z}}{|z_j|}.$$

This results in an upper bound on the value of x_j . If, for example, $(\mathbf{c}^\top \hat{\mathbf{x}} - \underline{z})/|z_j| < 1$, the integer variable x_j can be fixed to 0. The same procedure can be used for lower bounds.

Branch-and-cut

The combination of the branch-and-bound method with the generation of cutting planes to eliminate half-spaces is called branch-and-cut (see [38]). Here, the inequalities from Chapter 5 can be used. There are many possibilities for variation: The exact configuration must be empirically evaluated for each application. We are still far away from a true theoretical understanding of the advantages and disadvantages of the different variants. In general, the branch-and-cut method works very well in practice, but requires careful tuning of the parameters.

4 Integral Polyhedra

In this chapter we deal with polyhedra whose faces always contain integer points. Vertex solutions of the corresponding linear programs are therefore integral. For these special cases, the integer optimization problem can be solved via its LP relaxation, and thus in polynomial time (see *Introduction to Optimization*).

4.1 Total unimodularity

We first examine sufficient conditions on the constraint matrix $A \in \mathbb{Z}^{m \times n}$ for the polyhedron $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ to be integral for every integral right-hand side $\mathbf{b} \in \mathbb{Z}^m$. In particular, this allows to solve the corresponding MIP in polynomial time, since it is sufficient to solve its LP relaxation.

Definition 4.1. A matrix $A \in \mathbb{Z}^{m \times n}$ is called

- (a) *unimodular* if A has full row rank, and the determinant of every $(m \times m)$ -submatrix of A is in $\{-1, 0, 1\}$.
- (b) *totally unimodular (TU)* if the determinant of every square submatrix of A is in $\{-1, 0, 1\}$.

Note that all entries (the square submatrices of size size 1) of a totally unimodular matrix are either -1 , 0 , or 1 . Furthermore, the property of a matrix being TU is closed under taking submatrices.

Remark 4.2. Every totally unimodular matrix of full row rank is also unimodular. The converse does not hold: The integer matrix

$$\begin{pmatrix} 2 & 3 & 2 \\ 4 & 2 & 3 \\ 9 & 6 & 7 \end{pmatrix}$$

has determinant 1 (and therefore also full rank), so it is unimodular, but not totally unimodular, since its entries are not in $\{-1, 0, 1\}$:

Total unimodularity can be characterized in various ways that will prove useful later on.

Proposition 4.3. The following are equivalent:

- A is totally unimodular
- $[A, \mathbb{I}]$ is unimodular
- $\begin{pmatrix} A \\ -A \\ \mathbb{I} \\ -\mathbb{I} \end{pmatrix}$ is totally unimodular.
- $\begin{pmatrix} A & 0 \\ \mathbb{I} & \mathbb{I} \end{pmatrix}$ is totally unimodular.
- A^\top is totally unimodular.

Proof. exercise. □

Proposition 4.4. The vertex-arc incidence matrix of a directed graph $G = (V, E)$ (i.e., $A = (a_{v,e})_{v \in V, e \in E}$ with $a_{v,e} = 1$, if $e \in \delta^+(v)$, $a_{v,e} = -1$ if $e \in \delta^-(v)$, and $a_{v,e} = 0$ otherwise) is totally unimodular.

Proof. exercise. □

The following three results exhibit that a linear program with a (totally) unimodular constraint matrix always has an integral optimum solution, provided one exists at all.

Theorem 4.5. An integer matrix $A \in \mathbb{Z}^{m \times n}$ of full row rank is unimodular if and only if the polyhedron $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ is integral for all $\mathbf{b} \in \mathbb{Z}^m$.

Proof. Assume that A is unimodular and let $\mathbf{b} \in \mathbb{Z}^m$ be an arbitrary integer vector. First observe that, since every minimal face \mathcal{F} defines an affine space $\text{aff}(\mathcal{F}) = \mathcal{F}$ and since $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \geq \mathbf{0}\}$ does not contain any affine spaces that consist of more than a single point, every face of $\mathcal{P}^=(A, \mathbf{b}) = \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ contains some vertex. By Lemma 2.33, it therefore suffices to show that every vertex $\bar{\mathbf{x}}$ of $\mathcal{P}^=(A, \mathbf{b})$ is integral. Since A has full row rank, there is a basis $B \subseteq \{1, \dots, n\}$, $|B| = m$, such that $\bar{\mathbf{x}}_B = (A_{.B})^{-1}\mathbf{b}$ and $\bar{\mathbf{x}}_N = \mathbf{0}$, where $N := \{1, \dots, n\} \setminus B$ (see *Introduction to Optimization*). By Cramer's rule, for $i \in B$ we have

$$\bar{x}_i = \frac{\det(\mathbf{A}_{.d_1}, \dots, \mathbf{A}_{.d_{i-1}}, \mathbf{b}, \mathbf{A}_{.d_{i+1}}, \dots, \mathbf{A}_{.d_m})}{\det A_{.B}},$$

where $\{d_1, \dots, d_m\} := B$. Since A is unimodular, $\bar{\mathbf{x}}_B$ is integral (the numerator is an integer and the denominator is ± 1) and therefore also $\bar{\mathbf{x}}$.

Conversely, assume that $\mathcal{P}^=(A, \mathbf{b})$ is integral for all integer vectors $\mathbf{b} \in \mathbb{Z}^m$. Let $B \subseteq \{1, \dots, n\}$, $|B| = m$ be chosen arbitrarily such that $A_{.B} \in \mathbb{Z}^{m \times m}$ is invertible, i.e., $\det A_{.B} \neq 0$. We need to show that $\det A_{.B} \in \{-1, 1\}$. We claim that $(A_{.B})^{-1}\mathbf{z} \in \mathbb{Z}^m$ for all $\mathbf{z} \in \mathbb{Z}^m$. Note that this immediately follows from integrality of $\mathcal{P}^=(A, \mathbf{z})$ in case $(A_{.B})^{-1}\mathbf{z} \geq \mathbf{0}$.

Applying the claim for $\mathbf{z} = \mathbf{e}^j$, $j \in \{1, \dots, m\}$, yields that $(A_{.B})^{-1}\mathbf{e}^j \in \mathbb{Z}^m$, i.e., the j -th column of $(A_{.B})^{-1}$ is integral, and thus the entire matrix $(A_{.B})^{-1}$ is integral. Because then $\det(A_{.B})$ and $\det((A_{.B})^{-1})$ are integers, $\det(A_{.B}) \cdot \det((A_{.B})^{-1}) = 1$ implies that $\det(A_{.B}) \in \{-1, 1\}$. Since we have chosen $A_{.B}$ as an arbitrary invertible $(m \times m)$ -submatrix, it follows that A is unimodular.

To show the claim, take $\mathbf{y} \in \mathbb{Z}^m$ such that $\mathbf{t} := \mathbf{y} + (A_{.B})^{-1}\mathbf{z} \geq \mathbf{0}$ and set $\mathbf{b} := A_{.B}\mathbf{t} = A_{.B}\mathbf{y} + \mathbf{z} \in \mathbb{Z}^m$. Then, $\mathbf{x}_B := (A_{.B})^{-1}\mathbf{b} = \mathbf{t} \geq \mathbf{0}$ and $\mathbf{x}_N := \mathbf{0}$ for $N := \{1, \dots, n\} \setminus B$ defines a vertex \mathbf{x} of $\mathcal{P}^=(A, \mathbf{b})$ (see *Introduction to Optimization*). Since $\mathcal{P}^=(A, \mathbf{b})$ is integral, \mathbf{x} and thus \mathbf{t} are also integral, implying $(A_{.B})^{-1}\mathbf{z} \in \mathbb{Z}^m$. □

Corollary 4.6. For every square, unimodular matrix $U \in \mathbb{Z}^{m \times m}$, we have $U^{-1} \in \mathbb{Z}^{m \times m}$, and, in particular, U^{-1} is unimodular as well.

Proof. First note that U is invertible, since it has full (row) rank. The proof that $U^{-1}\mathbf{b} \in \mathbb{Z}^m$ for every $\mathbf{b} \in \mathbb{Z}^m$ is analogous to the proof of Theorem 4.5 via Cramer's rule applied to U (instead of $A_{.B}$). Since this also applies for $\mathbf{b} = \mathbf{e}^i$ and all $i \in \{1, \dots, m\}$, we have $U^{-1} \in \mathbb{Z}^{m \times m}$. From $\det(U) \cdot \det(U^{-1}) = 1$ and $\det(U) \in \{-1, 1\}$ (since U has full row rank) it follows that $\det(U^{-1}) \in \{-1, 1\}$, hence U^{-1} is unimodular. □

Example 4.7. The inverse of the matrix of Remark 4.2 is

$$\begin{pmatrix} -4 & -9 & 5 \\ -1 & -4 & 2 \\ 6 & 15 & -8 \end{pmatrix}$$

with determinant 1. △

Corollary 4.8 ([25]). An integer matrix $A \in \mathbb{Z}^{m \times n}$ is totally unimodular if and only if the polyhedron $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ is integral for all integer vectors $\mathbf{b} \in \mathbb{Z}^m$.

Proof. Proposition 4.3 (a) states that A is totally unimodular if and only if $[A, \mathbb{I}]$ is unimodular. Furthermore, for an integer vector \mathbf{b} it holds that $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ is integral if and only if $\{\mathbf{z} \in \mathbb{R}^{n+m} : [A, \mathbb{I}]\mathbf{z} = \mathbf{b}, \mathbf{z} \geq \mathbf{0}\}$ is integral, since the slack variables $\mathbf{s} = \mathbf{b} - A\mathbf{x} \geq \mathbf{0}$ can be chosen integral if \mathbf{x} , \mathbf{b} and A are integral. The application of Theorem 4.5 for the matrix $[A, \mathbb{I}]$ completes the proof. □

By simple transformations that preserve integrality and total unimodularity (by Proposition 4.3), we can extend Corollary 4.8 as follows.

Corollary 4.9. An integer matrix $A \in \mathbb{Z}^{m \times n}$ is totally unimodular if and only if the polyhedron $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{a} \leq A\mathbf{x} \leq \mathbf{b}, \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}\}$ is integral for all integer vectors $\mathbf{a}, \mathbf{b} \in \mathbb{Z}^m$ and $\mathbf{l}, \mathbf{u} \in \mathbb{Z}^n$.

We mention a similar extension of Theorem 4.5 without proof.

Corollary 4.10. An integer matrix $A \in \mathbb{Z}^{m \times n}$ of full row rank is unimodular, if and only if the polyhedron $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}, \mathbf{0} \leq \mathbf{x} \leq \mathbf{u}\}$ is integral for all integer vectors $\mathbf{b} \in \mathbb{Z}^m$ and $\mathbf{u} \in \mathbb{Z}^n$.

Note that we need bounds on the variables to obtain a characterization.

Proposition 4.11. If an integer matrix $A \in \mathbb{Z}^{m \times n}$ is totally unimodular, then the polyhedron $\mathcal{P}(A, \mathbf{b})$ is integral for all integer vectors $\mathbf{b} \in \mathbb{Z}^m$. The converse does not hold.

Proof. Consider any minimal face¹ $\mathcal{F} := \{\mathbf{x} \in \mathbb{R}^n : A_R \mathbf{x} = \mathbf{b}_R\}$ of $\mathcal{P}(A, \mathbf{b})$, where $R \subseteq \{1, \dots, m\}$ can be chosen such that A_R has full row rank. By Lemma 2.33, it is sufficient to show that \mathcal{F} contains an integer point. By rearranging the variables, we can assume that $A_R = [U, V]$ with U invertible. Because A is totally unimodular, so is its submatrix U , and by Corollary 4.6, U^{-1} is integral as well. Thus $\begin{pmatrix} U^{-1}\mathbf{b}_R \\ \mathbf{0} \end{pmatrix}$ is an integer point of \mathcal{F} .

For the second part of the statement, consider the matrix

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

This matrix is unimodular, since it has determinant -1 , but not totally unimodular, since the determinant of the upper left (2×2) -submatrix is 2. For every $\mathbf{b} \in \mathbb{Z}^3$, $\{\mathbf{x} \in \mathbb{R}^3 : A\mathbf{x} = \mathbf{b}\} = \{A^{-1}\mathbf{b}\} \subset \mathbb{Z}^3$ is contained in every face of $\mathcal{P}(A, \mathbf{b})$. Thus, by Lemma 2.33, the rational polyhedron $\mathcal{P}(A, \mathbf{b})$ is integral for every $\mathbf{b} \in \mathbb{Z}^3$, but A is not totally unimodular. □

¹Note that a minimal face is a affine space (see exercise).

Importantly, total unimodularity allows to extend strong duality to integer optimization problems.

Corollary 4.12 (integral duality). If $A \in \mathbb{Z}^{m \times n}$ is totally unimodular, then, for all integer vectors $\mathbf{b} \in \mathbb{Z}^m$ and $\mathbf{c} \in \mathbb{Z}^n$, both sides of

$$\begin{aligned} \max \{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\} &= \min \{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} \geq \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}, \quad \text{and} \\ \max \{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}\} &= \min \{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\} \end{aligned}$$

have integral optimum solutions, provided the optima exist.

Proof. For the first version, by Corollary 4.8 and Lemma 2.33, it follows that A is totally unimodular if, for all integer vectors $\mathbf{c} \in \mathbb{Z}^n$ and $\mathbf{b} \in \mathbb{Z}^m$, the optimum of the linear program $\max \{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ is attained by an integer point, if at all. Since A is totally unimodular if and only if $-A^\top$ is totally unimodular (Proposition 4.3), the same applies for the dual linear program $\min \{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} \geq \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$.

For the second version, by Proposition 4.11 and Lemma 2.33, if A is totally unimodular, then, for all integer vectors $\mathbf{c} \in \mathbb{Z}^n$ and $\mathbf{b} \in \mathbb{Z}^m$, the optimum of the linear program $\max \{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$ is attained by an integer point, if at all. For the dual LP, integrality follows from Theorem 4.5 with Lemma 2.33, because A^\top is totally unimodular (Proposition 4.3). Here, by possibly omitting redundant rows of A^\top , we may assume that the matrix has full row rank; the resulting matrix is then unimodular. \square

A very useful characterization of total unimodularity is mentioned next.

Theorem 4.13 ([18]). A matrix $A \in \mathbb{Z}^{m \times n}$ is totally unimodular if and only if for every subset of rows $I \subseteq \{1, \dots, m\}$ a vector $\mathbf{r} \in \{-1, 1\}^I$ exists, with $\sum_{i \in I} r_i \mathbf{A}_i \in \{-1, 0, 1\}^n$.

Proof. For a proof see [39]. \square

4.1.1 Applications in combinatorial optimization

Total unimodularity yields an different perspective on the max-flow min-cut theorem (see *Algorithmic Discrete Mathematics*).

Theorem 4.14. Let $((V, E), \mu, s, t)$ be a flow network with integer capacities μ . Then, an integral maximum s - t -flow f^* exists and its value is equal to the minimum capacity over all s - t -cuts, i.e.,

$$|f^*| = \min_{S \subseteq V \setminus \{s, t\}} \sum_{e \in \delta^+(S \cup \{s\})} \mu(e).$$

Proof. Let A be the vertex-arc incidence matrix of the directed graph $G = (V, E \cup \{(t, s)\})$ with $\mu((t, s)) := \infty$, and define $\boldsymbol{\mu} \in \mathbb{Z}_{\geq 0}^E$ by $\mu_e := \mu(e)$. Every s - t -flow in (V, E) corresponds to a circulation in G , i.e., a flow satisfying flow conservation also at s and t (see *Combinatorial Optimization*). Accordingly, a maximum s - t -flow is defined by the solution of the linear program

$$\max \{x_{(t,s)} : A\mathbf{x} = \mathbf{0}, \mathbf{0} \leq \mathbf{x} \leq \boldsymbol{\mu}\} = \max \{x_{(t,s)} : \begin{pmatrix} A \\ -\mathbb{I} \end{pmatrix} \mathbf{x} \leq \begin{pmatrix} \mathbf{0} \\ \boldsymbol{\mu} \end{pmatrix}, \mathbf{x} \geq \mathbf{0}\}.$$

This LP is feasible (e.g., $\mathbf{x} = \mathbf{0}$) and bounded by $x_{(t,s)} = \sum_{e \in \delta^-(t)} x_e \leq \sum_{e \in \delta^-(t)} \mu_e$. By Propositions 4.4 and 4.3 and Corollary 4.12, the LP and its dual

$$\min \{ \boldsymbol{\mu}^\top \mathbf{y} : A^\top \mathbf{z}^+ - A^\top \mathbf{z}^- + \mathbf{y} \geq \mathbf{e}^{(t,s)}; \mathbf{y}, \mathbf{z}^+, \mathbf{z}^- \geq \mathbf{0} \} = \min \{ \boldsymbol{\mu}^\top \mathbf{y} : A^\top \mathbf{z} + \mathbf{y} \geq \mathbf{e}^{(t,s)}, \mathbf{y} \geq \mathbf{0} \}$$

therefore have integral optimum solutions $\bar{\mathbf{x}}$, $\bar{\mathbf{y}}$ and $\bar{\mathbf{z}}$. In other words, the value $\bar{x}_{(t,s)}$ of a maximum s - t -flow is equal to the minimum value $\boldsymbol{\mu}^\top \bar{\mathbf{y}}$, where $\bar{\mathbf{y}}$ fulfills

$$\begin{aligned} \bar{z}_v - \bar{z}_w + \bar{y}_{(v,w)} &\geq 0, \quad \text{for } (v,w) \in E \setminus \{(t,s)\} \\ \bar{z}_t - \bar{z}_s + \bar{y}_{(t,s)} &\geq 1, \\ \bar{\mathbf{y}} &\geq \mathbf{0}. \end{aligned}$$

By weak complementary slackness (see *Introduction to Optimization*), and since $\mu_{(t,s)} = \infty$, we conclude that $\bar{y}_{(t,s)} = 0$, and therefore $\bar{z}_t \geq 1 + \bar{z}_s$. Since $s \notin W := \{w \in V : \bar{z}_w \geq \bar{z}_t\}$ and $t \in W$, it follows that $\delta^-(W)$ is an s - t -cut. Furthermore, $(v,w) \in \delta^-(W)$ implies $\bar{z}_w \geq \bar{z}_t > \bar{z}_v$, and thus

$$\bar{y}_{(v,w)} > \bar{z}_v - \bar{z}_w + \bar{y}_{(v,w)} \geq 0 \quad \forall (v,w) \in \delta^-(W)$$

yields, together with integrality of $\bar{\mathbf{y}}$, that $\bar{y}_{\delta^-(W)} \geq 1$. From this, and $\bar{\mathbf{y}}, \boldsymbol{\mu} \geq \mathbf{0}$, we conclude

$$\bar{x}_{(t,s)} = \boldsymbol{\mu}^\top \bar{\mathbf{y}} \geq \sum_{e \in \delta^-(W)} \mu_e.$$

In other words, there is an s - t -cut $\delta^-(W)$ whose capacity is not larger than the maximum flow value. Since there cannot be a cut of capacity smaller than the maximum flow value (see *Algorithmic Discrete Mathematics*), the theorem follows. \square

Totally unimodular matrices also naturally occur for undirected graphs. Let $G = (V, E)$ be an undirected graph and let $A \in \{0, 1\}^{E \times V}$ be its edge-vertex incidence matrix. In the following, we investigate sufficient conditions for the polyhedron $\mathcal{P} = \{ \mathbf{x} : A\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0} \}$ to be integral. This is interesting, e.g., because for the case $\mathbf{b} = \mathbf{1}$ the integral points of \mathcal{P} correspond to stable sets in G (i.e., a subset S of the vertices, with $e \not\subseteq S$ for all $e \in E$) or, more generally, feasible solutions of the set packing problem from Example 1.7 (with ground set $U = E$ and subsets $S_v = \delta(v)$ for all $v \in V$).

Theorem 4.15. The edge-vertex incidence matrix A of an undirected graph G is totally unimodular if and only if G is bipartite.

Proof. First assume that G is not bipartite. Then, G contains a cycle $C \subseteq E$ of odd length. We may assume that C does not contain any chords, i.e., edges between vertices not adjacent along the cycle, since every odd cycle with a chord contains a smaller odd cycle. The submatrix of A that belongs to the edges $E(C)$ and vertices $V(C)$ of C is (after suitable permutation of rows and columns)

$$M := \begin{pmatrix} 1 & 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & 1 & \dots & 0 & 0 \\ \vdots & & & \ddots & 0 & 0 \\ 0 & 0 & 0 & \dots & 1 & 1 \\ 1 & 0 & 0 & \dots & 0 & 1 \end{pmatrix} \in \{0, 1\}^{|C| \times |C|}.$$

To see that this matrix is not totally unimodular, consider the Leibniz formula for the determinant

$$\det(M) = \sum_{\sigma \in S_{|C|}} \left(\operatorname{sgn}(\sigma) \cdot \prod_{i=1}^{|C|} M_{i\sigma(i)} \right),$$

where $S_{|C|}$ denotes the symmetric group of all permutations of $\{1, \dots, |C|\}$. There are exactly two permutations with non-vanishing contributions to the sum, namely the identity (diagonal) and the cyclic map $\sigma(1) = 2, \sigma(2) = 3, \dots, \sigma(|C|) = 1$. The number of inversions of the latter permutation is $|C| - 1$, and therefore it has a positive sign since $|C|$ is odd. Hence, the determinant of M is 2 and, consequently, A is not totally unimodular.

Now, conversely, assume that G is bipartite. We let $\{V_1, V_2\} = V$ denote the corresponding partition of the vertex set. We arbitrarily fix $\mathbf{b} \in \mathbb{Z}^m$ and $\mathbf{c} \in \mathbb{R}^n$ and show that the linear program $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{A}\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ has an integral solution, provided it is bounded. This is sufficient because of Lemma 2.33 (d) and Corollary 4.8.

Let \mathbf{x}^* be an optimum solution of the linear problem $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{A}\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$. Let $\tau \sim \mathcal{U}(0, 1)$ be drawn uniformly at random from the interval $(0, 1)$. For $v \in V$, define $\mathbf{z} \in \mathbb{Z}_{\geq 0}^V$ via

$$z_v := \begin{cases} \lfloor x_v^* \rfloor + 1 & \text{if } v \in V_1 \text{ and } x_v^* - \lfloor x_v^* \rfloor \geq \tau, \\ \lfloor x_v^* \rfloor + 1 & \text{if } v \in V_2 \text{ and } x_v^* - \lfloor x_v^* \rfloor > 1 - \tau, \\ \lfloor x_v^* \rfloor & \text{otherwise.} \end{cases}$$

We claim that \mathbf{z} is feasible, i.e., $\mathbf{A}\mathbf{z} \leq \mathbf{b}$. To see this, consider $e = \{u, v\} \in E$ with $u \in V_1$ and $v \in V_2$ and the corresponding row $\mathbf{A}_e = (\mathbf{e}^u + \mathbf{e}^v)^\top$ of A . We need to show that \mathbf{z} satisfies $\mathbf{A}_e \mathbf{z} \leq b_e$, i.e., that $z_u + z_v \leq b_e$. We distinguish three cases.

- If $\lfloor x_u^* \rfloor + \lfloor x_v^* \rfloor = b_e$, then $x_u^*, x_v^* \in \mathbb{Z}$ and $z_u = \lfloor x_u^* \rfloor$ and $z_v = \lfloor x_v^* \rfloor$ (because $0 < \tau < 1$), hence $z_u + z_v = b_e$.
- If $\lfloor x_u^* \rfloor + \lfloor x_v^* \rfloor \leq b_e - 2$, then $\lfloor x_u^* \rfloor + 1 + \lfloor x_v^* \rfloor + 1 \leq b_e$ and therefore $z_u + z_v \leq b_e$.
- Otherwise, $\lfloor x_u^* \rfloor + \lfloor x_v^* \rfloor = b_e - 1$ holds. Then, $x_u^* - \lfloor x_u^* \rfloor + x_v^* - \lfloor x_v^* \rfloor = x_u^* + x_v^* - (b_e - 1) \leq 1$. If $x_u^* - \lfloor x_u^* \rfloor \geq \tau$, then $x_v^* - \lfloor x_v^* \rfloor \leq 1 - \tau$ and vice versa. So not both values x_u^* and x_v^* are rounded up, resulting in $z_u + z_v \leq b_e$.

Now, for $u \in V_1$ we obtain $\Pr_\tau[z_u = \lfloor x_u^* \rfloor + 1] = \Pr_\tau[\tau \leq x_u^* - \lfloor x_u^* \rfloor] = x_u^* - \lfloor x_u^* \rfloor$ and for $v \in V_2$ that $\Pr_\tau[z_v = \lfloor x_v^* \rfloor + 1] = \Pr_\tau[1 - \tau < x_v^* - \lfloor x_v^* \rfloor] = 1 - \Pr_\tau[1 - \tau \geq x_v^* - \lfloor x_v^* \rfloor] = 1 - \Pr_\tau[U \leq 1 - (x_v^* - \lfloor x_v^* \rfloor)] = 1 - (1 - (x_v^* - \lfloor x_v^* \rfloor)) = x_v^* - \lfloor x_v^* \rfloor$. Overall, we have shown

$$\Pr_\tau[z_v = \lfloor x_v^* \rfloor + 1] = x_v^* - \lfloor x_v^* \rfloor \quad \forall v \in V.$$

Let c^{IP} denote the largest objective function value of a feasible integer vector with respect to \mathbf{c} . By feasibility of $\mathbf{z} \in \mathbb{Z}^V$ we obtain for the expectation that

$$\begin{aligned} \mathbf{c}^\top \mathbf{x}^* &\geq c^{\text{IP}} \geq \mathbb{E}_\tau[\mathbf{c}^\top \mathbf{z}] = \sum_{v \in V} c_v \lfloor x_v^* \rfloor + \sum_{v \in V} c_v \Pr_\tau[z_v = \lfloor x_v^* \rfloor + 1] \\ &= \sum_{v \in V} c_v \lfloor x_v^* \rfloor + \sum_{v \in V} c_v (x_v^* - \lfloor x_v^* \rfloor) = \mathbf{c}^\top \mathbf{x}^*. \end{aligned}$$

It follows that $\mathbf{c}^\top \mathbf{x}^* = c^{\text{IP}}$. Thus, by Lemma 2.33 (d), the polyhedron $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ is integral. Corollary 4.8 implies that A is totally unimodular. \square

Remark 4.16. The technique used in the second half of the proof of Theorem 4.15 is often called *randomized rounding*.

Corollary 4.17. Consider the assignment problem from Example 1.5. The LP-relaxation of the binary program (obtained by replacing $x_{ij} \in \{0, 1\}$ by $x_{ij} \in [0, 1]$) always has an integral optimum solution.

Proof. Let A be the edge-vertex incidence matrix of the bipartite graph in Example 1.5. The LP relaxation of the assignment problem can be formulated as

$$\min \{ \mathbf{c}^\top \mathbf{x} : A^\top \mathbf{x} = \mathbf{1}, \mathbf{x} \geq \mathbf{0} \},$$

where the constraint $\mathbf{x} \leq \mathbf{1}$ is implied. According to Theorem 4.15, Proposition 4.3 (third and fifth item), and Corollary 4.8, this LP has an integral optimum solution. \square

4.2 The Hermite normal form

In this section we show that every rational matrix can be brought into so-called Hermite normal form. This form will help us to prove an integer analogue of the Farkas lemma. We will later need this lemma to show the integrality of polyhedra.

Definition 4.18. A matrix with full row rank is in *Hermite normal form* if it has the form $[B, 0]$, where B is a non-negative lower triangular matrix and each row has a unique maximum entry that is on the main diagonal.

Note that $[B, 0]$ being in Hermite normal form implies that B is invertible, since it is lower triangular and its determinant thus is the product of its strictly positive diagonal entries. Crucially, we can bring matrices into this form without disturbing too much the structural properties of matrices that we aim to establish. It suffices to allow the following operations.

Definition 4.19. The (*elementary*) *unimodular column operations* on a matrix are

- swapping two columns,
- multiplication of a column by -1 ,
- addition of an integer multiple of a column to another column.

The name of these operations is justified because they can be encoded as multiplications with unimodular matrices.

Observation 4.20. The unimodular column operations of Definition 4.19 can be performed by multiplying the following unimodular matrices from the right:

$$\begin{pmatrix} 1 & & & \\ & 0 & 1 & \\ & 1 & 0 & \\ & & & 1 \end{pmatrix}, \begin{pmatrix} 1 & & & \\ & -1 & & \\ & & & \\ & & & 1 \end{pmatrix}, \begin{pmatrix} 1 & & & \\ & 1 & 0 & \\ & k & 1 & \\ & & & 1 \end{pmatrix}.$$

With this, we can show how to bring matrices into the desired form.

Theorem 4.21. Every matrix $A \in \mathbb{Q}^{m \times n}$ of full row rank can be converted to Hermite normal form using unimodular column operations.

Proof. Let A be a rational matrix of full row rank. We may assume that $A \in \mathbb{Z}^{m \times n}$ is integral: Otherwise, we scale the matrix by a suitable positive factor before carrying out the following operations and by its reciprocal afterwards. This is possible, since all unimodular column operations commute with scalar multiplication.

We first show that A can be transformed into a matrix $[B, 0]$, where B is a lower triangular matrix and $B_{ii} > 0$ for all $i \in \{1, \dots, m\}$. Suppose we have already transformed A into the form

$$\begin{pmatrix} B & 0 \\ C & D \end{pmatrix},$$

where B is a lower triangular matrix with a positive diagonal. Using unimodular column operations, we establish $D_{11} \geq D_{12} \geq \dots \geq D_{1k} \geq 0$ for the first row of $D \in \mathbb{Z}^{\ell \times k}$ by multiplying columns by -1 and then reordering them. Iterative subtraction of the j -th column from the i -th column for $i < j$ and reordering can be used to bring the first row of D into non-decreasing order in such a way that, additionally, $\sum_{i=1}^k D_{1i}$ cannot further be reduced.

Since A has full row rank, we have $D_{11} > 0$. We claim that $D_{1i} = 0$ for $i \in \{2, \dots, k\}$. Otherwise, $D_{12} > 0$. If we subtract the second column from the first one and rearrange these columns if needed, we maintain that the first row of D is in non-increasing order, but we decreased $\sum_{i=1}^k D_{1i}$, in contradiction its minimality. Overall, we have brought one more row of the matrix into the desired form.

After we have performed this procedure m times, we have transformed A into a matrix $[B, 0]$, where B is a lower triangular matrix with a positive main diagonal.

In order to reach Hermite normal form, we need to establish $0 \leq B_{ij} < B_{ii}$ for every $j < i$. This can be achieved by considering for $i \leftarrow 2, \dots, m$ every $j \in \{1, \dots, i-1\}$ and adding an integer multiple (positive if $B_{ij} < 0$, negative if $B_{ij} \geq B_{ii}$) of the i -th column to the j -th column, such that $B_{ij} \in \{1, \dots, B_{ii} - 1\}$. In this way, the entries above row i are not affected and we achieve Hermite normal form. \square

Example 4.22. Consider the matrix

$$A = \begin{pmatrix} 4 & 8 & 0 & -4 & 4 \\ 3 & -6 & 6 & 9 & 0 \\ 2 & -2 & 0 & 4 & -4 \end{pmatrix}.$$

Multiplication of the first column by -2 , 1 , -1 and addition to the columns 2, 4, 5 yields

$$\begin{pmatrix} 4 & 0 & 0 & 0 & 0 \\ 3 & -12 & 6 & 12 & -3 \\ 2 & -6 & 0 & 6 & -6 \end{pmatrix}.$$

Negating and swapping the last column forward yields

$$\begin{pmatrix} 4 & 0 & 0 & 0 & 0 \\ 3 & 3 & -12 & 6 & 12 \\ 2 & 6 & -6 & 0 & 6 \end{pmatrix}.$$

Addition of the third column to column 5 and multiplication of the second column by 4 and -2 and addition to columns 3 and 4 yields

$$\begin{pmatrix} 4 & 0 & 0 & 0 & 0 \\ 3 & 3 & 0 & 0 & 0 \\ 2 & 6 & 18 & -12 & 0 \end{pmatrix}.$$

Addition of the fourth column to column 3 yields

$$\begin{pmatrix} 4 & 0 & 0 & 0 & 0 \\ 3 & 3 & 0 & 0 & 0 \\ 2 & 6 & 6 & -12 & 0 \end{pmatrix}.$$

Addition of twice the third column to column 4 yields

$$\begin{pmatrix} 4 & 0 & 0 & 0 & 0 \\ 3 & 3 & 0 & 0 & 0 \\ 2 & 6 & 6 & 0 & 0 \end{pmatrix}.$$

Subtraction of the second column from first column yields

$$\begin{pmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ -4 & 6 & 6 & 0 & 0 \end{pmatrix}.$$

Finally, subtraction of the third column from column 2 and addition to column 1 yields a matrix in Hermite normal form:

$$\begin{pmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 2 & 0 & 6 & 0 & 0 \end{pmatrix}.$$

△

The product of the unimodular matrices of Observation 4.20 involved in the proof of Theorem 4.21 results in a unimodular matrix U , because of the multiplicative property $\det(AB) = \det(A) \cdot \det(B)$ of the determinant.

Corollary 4.23. For every matrix $A \in \mathbb{Q}^{m \times n}$ of full row rank, there is a unimodular matrix $U \in \mathbb{Z}^{n \times n}$ such that

$$[B, 0] = AU,$$

is in Hermite normal form and is integral if A is.

Remark 4.24. For every rational matrix of full row rank, the Hermite normal form is unique. (U is unique if A is also regular). Furthermore, the Hermite normal form can be computed in polynomial time (see [39]).

We are now ready to prove an integer analogue of the Farkas lemma.

Theorem 4.25. For $A \in \mathbb{Q}^{m \times n}$ and $\mathbf{b} \in \mathbb{Q}^m$ exactly one of the systems

$$\begin{array}{l} Ax = \mathbf{b}, \\ \mathbf{x} \in \mathbb{Z}^n \end{array} \quad \text{or} \quad \begin{array}{l} \mathbf{y}^\top A \in \mathbb{Z}^m \\ \mathbf{y}^\top \mathbf{b} \notin \mathbb{Z} \\ \mathbf{y} \in \mathbb{Q}^m \end{array}$$

has a solution.

Proof. Both systems of equations cannot have a solution at the same time, because otherwise integrality of \mathbf{x} and $\mathbf{y}^\top A$ would imply integrality of $\mathbf{y}^\top \mathbf{b} = \mathbf{y}^\top A \mathbf{x}$.

Assume that the second system has no solution, i.e., for all $\mathbf{y} \in \mathbb{Q}^m$ for which $\mathbf{y}^\top \mathbf{A}$ is integral, $\mathbf{y}^\top \mathbf{b}$ is also integral. We have to show that $\mathbf{A}\mathbf{x} = \mathbf{b}$, $\mathbf{x} \in \mathbb{Z}^n$ has a solution.

First, observe that $\mathbf{A}\mathbf{x} = \mathbf{b}$ is (fractionally) feasible, since otherwise a rational vector $\mathbf{y} \in \mathbb{Q}^m$ with $\mathbf{y}^\top \mathbf{A} = \mathbf{0}$ and $\mathbf{y}^\top \mathbf{b} \neq 0$ would exist by Gaussian elimination. After appropriate rescaling, there would then be a \mathbf{y} with $\mathbf{y}^\top \mathbf{A} = \mathbf{0}$ and $\mathbf{y}^\top \mathbf{b} = \frac{1}{2}$ – a contradiction to our assumption.

We may therefore assume that \mathbf{A} has full row rank (after eliminating redundant rows from $\mathbf{A}\mathbf{x} = \mathbf{b}$). According to Corollary 4.23, there is a unimodular matrix $\mathbf{U} \in \mathbb{Z}^{n \times n}$ with $\mathbf{A}\mathbf{U} = [\mathbf{B}, \mathbf{0}]$ in Hermite normal form. Since $\mathbf{B}^{-1}\mathbf{A}\mathbf{U} = \mathbf{B}^{-1}[\mathbf{B}, \mathbf{0}] = [\mathbf{I}, \mathbf{0}] \in \mathbb{Z}^{m \times n}$ and, by Corollary 4.6, $\mathbf{U}^{-1} \in \mathbb{Z}^{n \times n}$, we have $\mathbf{B}^{-1}\mathbf{A} \in \mathbb{Z}^{m \times n}$. It follows from our assumption (with $\mathbf{y} = (\mathbf{B}^{-1})_i^\top$ for $i \in \{1, \dots, m\}$) that $\mathbf{B}^{-1}\mathbf{b} \in \mathbb{Z}^m$. Hence, $\hat{\mathbf{x}} := \mathbf{U} \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0} \end{pmatrix} \in \mathbb{Z}^n$ is integral. Because of

$$\mathbf{A}\hat{\mathbf{x}} = \mathbf{A}\mathbf{U} \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0} \end{pmatrix} = [\mathbf{B}, \mathbf{0}] \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0} \end{pmatrix} = \mathbf{b},$$

we obtain that $\hat{\mathbf{x}}$ is an integral solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$. □

Corollary 4.26. We can in polynomial time find an integer solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$ with $\mathbf{A} \in \mathbb{Q}^{m \times n}$, $\mathbf{b} \in \mathbb{Q}^m$ or a certificate that none exists.

Proof. We first solve $\mathbf{A}\mathbf{x} = \mathbf{b}$ fractionally via Gaussian elimination. If no solution exists, we obtain a vector \mathbf{y} with $\mathbf{y}^\top \mathbf{A} = \mathbf{0}$ and $\mathbf{y}^\top \mathbf{b} = \frac{1}{2}$, that certifies that no integral solution exists according to Theorem 4.25. If a solution exists, we eliminate redundant rows and compute the Hermite normal form $[\mathbf{B}, \mathbf{0}]$ in polynomial time via Remark 4.24, and invert \mathbf{B} via Gaussian elimination. Now, either $\hat{\mathbf{x}} = \mathbf{U} \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0} \end{pmatrix}$ is an integer solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$ or a certificate $\mathbf{y} = (\mathbf{B}^{-1})_i^\top$ with $(\mathbf{B}^{-1}\mathbf{b})_i \notin \mathbb{Z}$ for infeasibility has been found according to Theorem 4.25. □

Remark 4.27. Corollary 4.26 only holds for systems of equations of the form $\mathbf{A}\mathbf{x} = \mathbf{b}$. If, additionally, \mathbf{x} is constrained, i.e., via $\mathbf{x} \in \{0, 1\}^n$, the hardness result of Theorem 3.1 carries over.

4.3 Total dual integrality

In Section 4.1, we saw that total unimodularity of an integer matrix \mathbf{A} ensures that, for every integral right-hand side \mathbf{b} , the polyhedron $\mathcal{P}_1 = \text{conv}(\{\mathbf{x} \in \mathbb{Z}^n : \mathbf{A}\mathbf{x} \leq \mathbf{b}\})$ is completely described by the original inequalities $\mathbf{A}\mathbf{x} \leq \mathbf{b}$. If, on the other hand, we decide on a particular right-hand side \mathbf{b} , then TDI (total dual integrality) is the right concept to guarantee integrality of a polyhedron.

Definition 4.28. The system of inequalities $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ with $\mathbf{A} \in \mathbb{Q}^{m \times n}$, $\mathbf{b} \in \mathbb{Q}^m$ is called *totally dual integral (TDI)* if, for every integer vector $\mathbf{c} \in \mathbb{Z}^m$, there is an integral optimum solution of the LP $\min \{\mathbf{b}^\top \mathbf{y} : \mathbf{A}^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$ if an optimum exists.

The TDI property of a system $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ has a geometric interpretation: Let \mathbf{c} be an integer vector in the cone spanned by the rows of \mathbf{A} , i.e., there is $\mathbf{y} \geq \mathbf{0}$ with $\mathbf{A}^\top \mathbf{y} = \mathbf{c}$. Among all the possibilities to express \mathbf{c} as a conic combination of the rows of \mathbf{A} , let \mathcal{S} be the set minimum conic combinations with respect to \mathbf{b} , i.e.,

$$\mathcal{S} := \arg \min \{\mathbf{b}^\top \mathbf{y} : \mathbf{A}^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}.$$

Then, $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ is TDI if and only if \mathcal{S} contains an integer vector. In other words, this means that among all shortest ways to express \mathbf{c} as a conic combination of the rows of \mathbf{A} , there is one that is an integer. We will later express this interpretation in terms of Hilbert bases, which play an important role in integer programming. Total dual integrality is directly connected to total unimodularity as in Corollary 4.12, but even for $\mathbf{b} \in \mathbb{Q}^m$.

Observation 4.29. If $A \in \mathbb{Z}^{m \times n}$ is totally unimodular, then $Ax \leq \mathbf{b}$ is TDI for all $\mathbf{b} \in \mathbb{Q}^m$.

Proof. If A is totally unimodular, by Corollary 4.8, we have that

$$\{\mathbf{y} \in \mathbb{R}^m : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\} = \{\mathbf{y} \in \mathbb{R}^m : \begin{pmatrix} A^\top \\ -A^\top \end{pmatrix} \mathbf{y} \leq \begin{pmatrix} \mathbf{c} \\ -\mathbf{c} \end{pmatrix}, \mathbf{y} \geq \mathbf{0}\}$$

is integral for all $\mathbf{c} \in \mathbb{Z}^n$. Lemma 2.33 then implies that

$$\min\{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\} = -\max\{-\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$$

has an integral optimum solution for all $\mathbf{b} \in \mathbb{R}^m$ for which an optimum exists. By definition, $Ax \leq \mathbf{b}$ is TDI. \square

Note also that TDI is a property of a system of inequalities and not a property of the polyhedron $\mathcal{P} = \{\mathbf{x} \in \mathbb{R}^n : Ax \leq \mathbf{b}\}$. We will see that there are many ways to describe \mathcal{P} , but there are only a few that also possess the TDI property. To obtain such a description, it may be necessary to add many redundant inequalities to the initial formulation.

We now establish a connection between $Ax \leq \mathbf{b}$ being TDI and the integrality of the corresponding polyhedron $\mathcal{P}(A, \mathbf{b})$. We first strengthen Lemma 2.33.

Lemma 4.30. If $\max\{\mathbf{c}^\top \mathbf{x} : Ax \leq \mathbf{b}\}$ is an integer for all $\mathbf{c} \in \mathbb{Z}^n$ for which it is finite, then $\mathcal{P}(A, \mathbf{b})$ is integral.

Proof. Let $\mathcal{F} \neq \emptyset$ be a minimal face of $\mathcal{P}(A, \mathbf{b})$. By Lemma 2.33, it suffices to show that \mathcal{F} contains integral points. Since \mathcal{F} is minimal, it follows that $\mathcal{F} = \{\mathbf{x} \in \mathcal{P}(A, \mathbf{b}) : A_{\text{eq}(\mathcal{F})} \mathbf{x} = \mathbf{b}_{\text{eq}(\mathcal{F})}\} = \{\mathbf{x} \in \mathbb{R}^n : A_{\text{eq}(\mathcal{F})} \mathbf{x} = \mathbf{b}_{\text{eq}(\mathcal{F})}\}$.

If \mathcal{F} does not contain any integral points, then, by Theorem 4.25 there exists $\mathbf{y} \in \mathbb{Q}^{\text{eq}(\mathcal{F})}$ with $\mathbf{c}^\top := \mathbf{y}^\top A_{\text{eq}(\mathcal{F})} \in \mathbb{Z}^n$ and $\gamma := \mathbf{y}^\top \mathbf{b}_{\text{eq}(\mathcal{F})} \notin \mathbb{Z}$. We may assume that \mathbf{y} is non-negative, otherwise choose $\mathbf{s} \in \mathbb{Z}_{\geq 0}^{\text{eq}(\mathcal{F})}$ large enough so that $\mathbf{y} + \mathbf{s} \geq \mathbf{0}$, $(\mathbf{y} + \mathbf{s})^\top A_{\text{eq}(\mathcal{F})} \in \mathbb{Z}^n$ and $(\mathbf{y} + \mathbf{s})^\top \mathbf{b}_{\text{eq}(\mathcal{F})} \notin \mathbb{Z}$, and use $\mathbf{y}' := \mathbf{y} + \mathbf{s}$ instead of \mathbf{y} . We claim that $\max\{\mathbf{c}^\top \mathbf{x} : Ax \leq \mathbf{b}\}$ exists and is assumed by every $\hat{\mathbf{x}} \in \mathcal{F}$. This is the case since $\mathcal{F} \neq \emptyset$ and, for all $\mathbf{x} \in \mathcal{P}(A, \mathbf{b})$ it holds that

$$\mathbf{c}^\top \mathbf{x} = \mathbf{y}^\top A_{\text{eq}(\mathcal{F})} \mathbf{x} \stackrel{\mathbf{y} \geq \mathbf{0}}{\leq} \mathbf{y}^\top \mathbf{b}_{\text{eq}(\mathcal{F})} = \mathbf{y}^\top A_{\text{eq}(\mathcal{F})} \hat{\mathbf{x}} = \mathbf{c}^\top \hat{\mathbf{x}}.$$

It also follows that $\gamma = \mathbf{y}^\top \mathbf{b}_{\text{eq}(\mathcal{F})} = \max\{\mathbf{c}^\top \mathbf{x} : Ax \leq \mathbf{b}\}$. By assumption, the right-hand side is an integer, which contradicts our choice of γ . Therefore, \mathcal{F} contains integral points. \square

We can now weaken the condition of Prop 4.11 for the case where we have a specific right-hand side $\mathbf{b} \in \mathbb{Z}^m$.

Theorem 4.31. If $Ax \leq \mathbf{b}$ is TDI and $\mathbf{b} \in \mathbb{Z}^m$, then $\mathcal{P}(A, \mathbf{b})$ is integral.

Proof. By definition of $Ax \leq \mathbf{b}$ being TDI, and since \mathbf{b} is an integer vector, the value $\min\{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$ is an integer for all $\mathbf{c} \in \mathbb{Z}^n$ for which it is finite. Due to strong duality (see *Introduction to Optimization*), this number coincides with $\max\{\mathbf{c}^\top \mathbf{x} : Ax \leq \mathbf{b}\}$. The statement thus follows from Lemma 4.30. \square

We have already hinted at the connection between the geometric interpretation of TDI and Hilbert bases. We now make this precise.

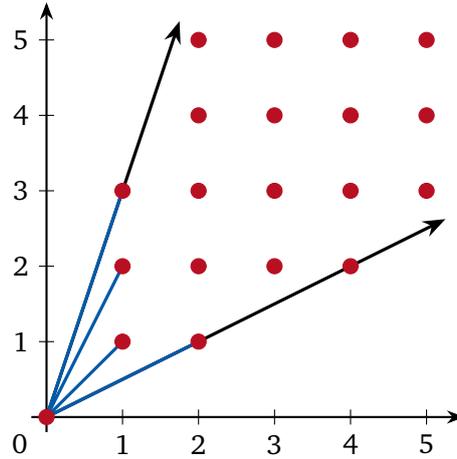


Figure 4.1: Hilbert basis of a polyhedral cone.

Definition 4.32. A Hilbert basis of a rational polyhedral cone $\mathcal{C} = \mathcal{P}(A, \mathbf{0}) \subseteq \mathbb{R}^n$ is a finite subset $\mathcal{H} \subseteq \mathcal{C}$ such that every $\mathbf{z} \in \mathcal{C} \cap \mathbb{Z}^n$ is a non-negative integral linear combination of elements from \mathcal{H} , i.e.,

$$\mathbf{z} = \sum_{\mathbf{h} \in \mathcal{H}} \lambda(\mathbf{h}) \mathbf{h} \quad \text{with } \lambda: \mathcal{H} \rightarrow \mathbb{Z}_{\geq 0}.$$

The Hilbert basis \mathcal{H} is integral if $\mathcal{H} \subseteq \mathbb{Z}^n$.

Example 4.33. Consider $\mathcal{C} = \text{cone}(\{(\frac{1}{3}), (\frac{2}{1})\})$ (see Figure 4.1). An integer Hilbert basis of \mathcal{C} is given by $\mathcal{H} = \{(\frac{1}{3}), (\frac{1}{2}), (\frac{1}{1}), (\frac{2}{1})\}$. △

For proofs of the following facts see [39, Chapter 16.4].

Remark 4.34.

- (a) For every rational polyhedral cone there exists an integral Hilbert basis.
- (b) If the cone has a vertex, the minimal integral Hilbert basis is uniquely determined.

We derive a characterization for $A\mathbf{x} \leq \mathbf{b}$ to be TDI.

Theorem 4.35. Let $A \in \mathbb{Q}^{m \times n}$, $\mathbf{b} \in \mathbb{Q}^m$ and $\mathcal{P}(A, \mathbf{b}) \neq \emptyset$. Then, the system of inequalities $A\mathbf{x} \leq \mathbf{b}$ is TDI, if and only if the rows of $A_{\text{eq}(\mathcal{F})}$ form a Hilbert basis of $\text{cone}((A_{\text{eq}(\mathcal{F})})^\top)$ for every (minimal) face \mathcal{F} of $\mathcal{P}(A, \mathbf{b})$.

To understand the proof of Theorem 4.35, the following geometric interpretation may be helpful. Let $\mathbf{c} \in \text{cone}(A^\top)$ be an integer vector such that $\min\{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$ is finite. Then, we know that $\mathbf{c}^\top \mathbf{x}$ assumes its optimum on a face of \mathcal{P} . By complementary slackness (see *Introduction to Optimization*), it follows that if $\mathbf{c}^\top \mathbf{x}$ assumes the optimum on face \mathcal{F} , then \mathbf{c} lies in the cone spanned by the rows of $A_{\text{eq}(\mathcal{F})}$. The condition that the rows induced by \mathcal{F} form a Hilbert basis is equivalent to the fact that \mathbf{c} is a non-negative integer combination of the rows induced by \mathcal{F} . Formally, there is an integer vector $\mathbf{y}^* \in \mathbb{Z}^m$, $A^\top \mathbf{y}^* = \mathbf{c}$, $\mathbf{y}^* \geq \mathbf{0}$ with $\mathbf{b}^\top \mathbf{y}^* = \min\{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$. This corresponds exactly to the TDI property of the system $A\mathbf{x} \leq \mathbf{b}$.

Proof. Let $A\mathbf{x} \leq \mathbf{b}$ be TDI. Further, let $\mathcal{F} \neq \emptyset$ be a face of $\mathcal{P} := \mathcal{P}(A, \mathbf{b})$ and $\mathbf{c} \in \text{cone}((A_{\text{eq}(\mathcal{F})})^\top) \cap \mathbb{Z}^n$. We need to show that \mathbf{c} is a non-negative integer combination of the row vectors of $A_{\text{eq}(\mathcal{F})}$. By assumption, \mathbf{c} is a

non-negative linear combination of the row vectors of $A_{\text{eq}(\mathcal{F})}$, i.e., there is $\hat{\mathbf{y}} \geq \mathbf{0}$ with $\hat{y}_i = 0$ for all $i \notin \text{eq}(\mathcal{F})$ and $\mathbf{c}^\top = \hat{\mathbf{y}}^\top A$.

We now claim that $\max\{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$ exists and is assumed for every $\hat{\mathbf{x}} \in F$. This follows since $\mathcal{F} \neq \emptyset$ and, for all $\mathbf{x} \in \mathcal{P}$ and $\hat{\mathbf{x}} \in \mathcal{F}$, it holds that

$$\mathbf{c}^\top \mathbf{x} = \hat{\mathbf{y}}^\top A\mathbf{x} \leq \hat{\mathbf{y}}^\top \mathbf{b} = \hat{\mathbf{y}}_{\text{eq}(\mathcal{F})}^\top \mathbf{b}_{\text{eq}(\mathcal{F})} = \hat{\mathbf{y}}_{\text{eq}(\mathcal{F})}^\top A_{\text{eq}(\mathcal{F})} \hat{\mathbf{x}} = \hat{\mathbf{y}}^\top A\hat{\mathbf{x}} = \mathbf{c}^\top \hat{\mathbf{x}},$$

because $\hat{y}_i = 0$ for all $i \notin \text{eq}(\mathcal{F})$. By strong duality (see *Introduction to Optimization*), we conclude that the value $\min\{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$ is finite. Using that $A\mathbf{x} \leq \mathbf{b}$ is TDI, we obtain

$$\min\{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \in \mathbb{Z}^m, \mathbf{y} \geq \mathbf{0}\} = \max\{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}.$$

Let \mathbf{y}^* be an integer optimum solution of $\min\{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$. For every $\hat{\mathbf{x}} \in \mathcal{F}$ we have $A_i \hat{\mathbf{x}} = b_i$ for all $i \in \text{eq}(\mathcal{F})$ and, by minimality of \mathcal{F} , $A_i \hat{\mathbf{x}} < b_i$ for all $i \notin \text{eq}(\mathcal{F})$. Since $\hat{\mathbf{x}} \in \mathcal{F}$ is an optimal solution of $\max\{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$, it follows from weak complementary slackness that $y_i^* = 0$ for all $i \notin \text{eq}(\mathcal{F})$. This completes the first part of the proof, since \mathbf{y}^* is integral, $y_i^* = 0$ for all $i \notin \text{eq}(\mathcal{F})$, $y_i^* \geq 0$ for all $i \in \text{eq}(\mathcal{F})$, and \mathbf{c} is a non-negative integer combination of the rows of $A_{\text{eq}(\mathcal{F})}$.

Now, let $\mathbf{c} \in \mathbb{Z}^n$ such that $\gamma := \min\{\mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$ exists. By strong duality, we know that $\gamma = \max\{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$ exists. Let \mathcal{F} be a minimal face of \mathcal{P} on which the optimum value of $\mathbf{c}^\top \mathbf{x}$ is attained.

Note that $\mathbf{c} \in \text{cone}((A_{\text{eq}(\mathcal{F})})^\top)$: Due to weak complementary slackness, there is a vector $\mathbf{y}^* \geq \mathbf{0}$ with $y_i^* = 0$ for all $i \notin \text{eq}(\mathcal{F})$ and $\mathbf{c}^\top = \sum_{i \in \text{eq}(\mathcal{F})} A_i \cdot y_i^*$.

Since the rows of $A_{\text{eq}(\mathcal{F})}$ are a Hilbert basis of $\text{cone}((A_{\text{eq}(\mathcal{F})})^\top)$, there is a non-negative integer vector $\tilde{\mathbf{y}}$ with $\mathbf{c}^\top = \sum_{i \in \text{eq}(\mathcal{F})} A_i \cdot \tilde{y}_i$. We can set $y_i := 0$ for all $i \notin \text{eq}(\mathcal{F})$ and $y_i := \tilde{y}_i$ for all $i \in \text{eq}(\mathcal{F})$ and obtain an integer vector \mathbf{y} with $A^\top \mathbf{y} = \mathbf{c}$ and $\mathbf{y} \geq \mathbf{0}$. Furthermore, for $\hat{\mathbf{x}} \in \mathcal{F}$, we have

$$\mathbf{b}^\top \mathbf{y} = \sum_{i \in \text{eq}(\mathcal{F})} y_i b_i = \sum_{i \in \text{eq}(\mathcal{F})} y_i A_i \hat{\mathbf{x}} = \mathbf{c}^\top \hat{\mathbf{x}} = \gamma.$$

Hence, $A\mathbf{x} \leq \mathbf{b}$ is TDI. □

We apply Theorem 4.35 to show that the TDI property of a system $A\mathbf{x} \leq \mathbf{b}$ of inequalities is preserved if we restrict ourselves to subsystems corresponding to faces of the polyhedron $\mathcal{P}(A, \mathbf{b})$.

Corollary 4.36. If the system $A\mathbf{x} \leq \mathbf{b}$, $\mathbf{c}^\top \mathbf{x} \leq d$ with $A \in \mathbb{Q}^{m \times n}$, $\mathbf{b} \in \mathbb{Q}^m$, $\mathbf{c} \in \mathbb{Q}^n$ and $d \in \mathbb{Q}$ is TDI, then the system $A\mathbf{x} \leq \mathbf{b}$, $\mathbf{c}^\top \mathbf{x} \leq d$, $-\mathbf{c}^\top \mathbf{x} \leq -d$ is TDI as well.

Proof. Let $\mathcal{F} \neq \emptyset$ be a face of the polyhedron $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}, \mathbf{c}^\top \mathbf{x} \leq d, -\mathbf{c}^\top \mathbf{x} \leq -d\}$. Then, \mathcal{F} is also a face of the polyhedron $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}, \mathbf{c}^\top \mathbf{x} \leq d\}$. Since the system $A\mathbf{x} \leq \mathbf{b}$, $\mathbf{c}^\top \mathbf{x} \leq d$ is TDI, by Theorem 4.35 that $((A_{\text{eq}(\mathcal{F})})^\top, \mathbf{c})$ is a Hilbert basis of $\text{cone}((A_{\text{eq}(\mathcal{F})})^\top, \mathbf{c})$. Let $\mathbf{z} \in \text{cone}((A_{\text{eq}(\mathcal{F})})^\top, \mathbf{c}, -\mathbf{c}) \cap \mathbb{Z}^n$. Then, \mathbf{z} is of the form

$$\mathbf{z} = \sum_{i \in \text{eq}(\mathcal{F})} \lambda_i A_i^\top + \mu \mathbf{c} - \sigma \mathbf{c}$$

with $\lambda_i \geq 0$, $i \in \text{eq}(\mathcal{F})$, $\mu \geq 0$, $\sigma \geq 0$. This is equivalent to $\mathbf{z} + \sigma \mathbf{c} \in \text{cone}((A_{\text{eq}(\mathcal{F})})^\top, \mathbf{c})$. Let $\sigma' \in \mathbb{Z}_{\geq 0}$, $\sigma' \geq \sigma$ be chosen such that $\sigma' \mathbf{c} \in \mathbb{Z}^n$. It follows that $\mathbf{z} + \sigma' \mathbf{c} \in \text{cone}((A_{\text{eq}(\mathcal{F})})^\top, \mathbf{c}) \cap \mathbb{Z}^n$. Since $((A_{\text{eq}(\mathcal{F})})^\top, \mathbf{c})$ is a Hilbert basis of $\text{cone}((A_{\text{eq}(\mathcal{F})})^\top, \mathbf{c})$, we can write $\mathbf{z} + \sigma' \mathbf{c}$ as

$$\mathbf{z} + \sigma' \mathbf{c} = \sum_{i \in \text{eq}(\mathcal{F})} \lambda'_i A_i^\top + \mu' \mathbf{c},$$

where $\lambda'_i \geq 0$ is an integer for all $i \in \text{eq}(\mathcal{F})$, and $\mu' \geq 0$ is also an integer. It follows that

$$\mathbf{z} = \sum_{i \in \text{eq}(\mathcal{F})} \lambda'_i \mathbf{A}_i^\top + \mu' \mathbf{c} - \sigma' \mathbf{c}.$$

That is, \mathbf{z} can be defined as a non-negative integer combination of the generators of $\text{cone}((\mathbf{A}_{\text{eq}(\mathcal{F})})^\top, \mathbf{c}, -\mathbf{c})$. Hence, $((\mathbf{A}_{\text{eq}(\mathcal{F})})^\top, \mathbf{c}, -\mathbf{c})$ is a Hilbert basis of $\text{cone}((\mathbf{A}_{\text{eq}(\mathcal{F})})^\top, \mathbf{c}, -\mathbf{c})$. Applying Theorem 4.35 completes the proof. \square

We mention without proof a characterization of TDI via Hilbert bases that is helpful for algorithmically checking whether an inequality system is TDI.

Theorem 4.37. Let $A \in \mathbb{Z}^{m \times n}$ and $\mathbf{b} \in \mathbb{Q}^m$ with $\mathcal{P}(A, \mathbf{b}) \neq \emptyset$. The system $A\mathbf{x} \leq \mathbf{b}$ is TDI if and only if

- the rows of A form a Hilbert basis of the cone spanned by the rows of A , and
- for every $S \subseteq \{1, \dots, m\}$ the linear program

$$\min \{ \mathbf{b}^\top \mathbf{y} : A^\top \mathbf{y} = \sum_{i \in S} (\mathbf{A}_i)^\top, \mathbf{y} \geq \mathbf{0} \}$$

has an integral optimum solution.

The following theorem shows that every rational polyhedron \mathcal{P} can be described by a TDI system. The key idea of the proof is to construct, for every minimal face \mathcal{F} of \mathcal{P} , a Hilbert basis of the cone spanned by inequalities induced by \mathcal{F} .

Theorem 4.38. For every rational polyhedron \mathcal{P} there is a TDI system $A\mathbf{x} \leq \mathbf{b}$ with integral matrix A such that $\mathcal{P} = \mathcal{P}(A, \mathbf{b})$.

Proof. Let $D \in \mathbb{Q}^{m \times n}$, $\mathbf{d} \in \mathbb{Q}^m$ be such that $\mathcal{P} = \mathcal{P}(D, \mathbf{d})$ and the rows in $D\mathbf{x} \leq \mathbf{d}$ are not redundant. Let $\mathcal{F}_1, \dots, \mathcal{F}_t$ be all (minimal) faces of \mathcal{P} and let $\mathcal{H}_i \subseteq \mathbb{Z}^n$ be a minimal integer Hilbert basis of $\text{cone}((D_{\text{eq}(\mathcal{F}_i)})^\top)$ for every $i \in \{1, \dots, t\}$ (exists by Remark 4.34). Define A as the matrix whose rows are exactly the vectors in $\bigcup_{i=1}^t \mathcal{H}_i$ and consider the k -th row vector \mathbf{A}_k of A . Then, \mathbf{A}_k is an element of some integer Hilbert basis \mathcal{H}_i . Consequently, A is an integral matrix. Let \mathbf{b} be defined via $b_k := \max \{ \mathbf{A}_k \cdot \mathbf{x} : \mathbf{x} \in \mathcal{P} \}$. Note that this maximum exists because $\mathbf{A}_k^\top = \sum_{j \in \text{eq}(\mathcal{F}_i)} \lambda_j (D_j)^\top \in \text{cone}((D_{\text{eq}(\mathcal{F}_i)})^\top)$, and thus

$$b_k = \max \{ \mathbf{A}_k \cdot \mathbf{x} : \mathbf{x} \in \mathcal{P} \} = \mathbf{A}_k \cdot \hat{\mathbf{x}} \text{ for all } \hat{\mathbf{x}} \in \mathcal{F}_i. \quad (4.1)$$

We claim that $\mathcal{P} = \mathcal{P}(A, \mathbf{b})$.

Clearly, $\mathcal{P} \subseteq \mathcal{P}(A, \mathbf{b})$, since, for every $\mathbf{y} \in \mathcal{P}$ and row \mathbf{A}_k , it holds that $\mathbf{A}_k \cdot \mathbf{y} \leq \max \{ \mathbf{A}_k \cdot \mathbf{x} : \mathbf{x} \in \mathcal{P} \} = b_k$.

Conversely, if $\mathbf{y} \notin \mathcal{P}$, then there is a row index $\ell \in \{1, \dots, m\}$ such that $D_\ell \cdot \mathbf{y} > d_\ell$. Let $i \in \{1, \dots, t\}$ such that $\ell \in \text{eq}(\mathcal{F}_i)$, which exists since $D\mathbf{x} \leq \mathbf{d}$ is not redundant.

Now, let $\delta \geq 0$ be such that $\delta(D_\ell)^\top \in \mathbb{Z}^n$. Since $\{\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(s)}\} := \mathcal{H}_i$ is a Hilbert basis of $\text{cone}((D_{\text{eq}(\mathcal{F}_i)})^\top)$, there are non-negative integer multiples $\delta_1, \dots, \delta_s$ of δ with

$$\delta D_\ell^\top = \sum_{r=1}^s \delta_r \mathbf{a}^{(r)}.$$

Denote by k_r the row of A with $A_{k_r}^\top = \mathbf{a}^{(r)}$ for all $r \in \{1, \dots, s\}$. Taking into account (4.1), we obtain for $\hat{\mathbf{x}} \in \mathcal{F}_i$ that

$$\sum_{r=1}^s \delta_r A_{k_r} \cdot \mathbf{y} = \sum_{r=1}^s \delta_r (\mathbf{a}^{(r)})^\top \mathbf{y} = \delta \mathbf{D}_\ell \cdot \mathbf{y} > \delta d_\ell = \delta \mathbf{D}_\ell \cdot \hat{\mathbf{x}} = \sum_{r=1}^s \delta_r (\mathbf{a}^{(r)})^\top \hat{\mathbf{x}} = \sum_{r=1}^s \delta_r b_{k_r}.$$

So there is a $\hat{r} \in \{1, \dots, s\}$ with $A_{k_{\hat{r}}} \cdot \mathbf{y} > b_{k_{\hat{r}}}$. Thus $\mathbf{y} \notin \mathcal{P}(A, \mathbf{b})$. This completes the proof that $\mathcal{P} = \mathcal{P}(A, \mathbf{b})$.

Finally, the system $A\mathbf{x} \leq \mathbf{b}$ is TDI as required: If \mathcal{F}_i is a face of \mathcal{P} and $\{\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(s)}\} := \mathcal{H}_i$, then $A_{k_r} \cdot \hat{\mathbf{x}} = (\mathbf{a}^{(r)})^\top \hat{\mathbf{x}} = b_{k_r}$ for all $\hat{\mathbf{x}} \in \mathcal{F}_i$ and $r \in \{1, \dots, s\}$, see (4.1). Thus $(\mathbf{a}^{(r)})^\top$ is a row of $A_{\text{eq}(\mathcal{F}_i)}$. In addition, $\{\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(s)}\}$ is an (integer) Hilbert basis of $\text{cone}((D_{\text{eq}(\mathcal{F}_i)})^\top) = \text{cone}((A_{\text{eq}(\mathcal{F}_i)})^\top)$ because of $\mathcal{P}(D, \mathbf{d}) = \mathcal{P}(A, \mathbf{b})$. Applying Theorem 4.35 completes the proof. \square

In fact, it is even true that \mathcal{P} is integral if and only if \mathbf{b} can be chosen as an integer vector. One direction is obtained from the proof of Theorem 4.38, because \mathcal{P} is integral if and only if every face contains an integer point (Lemma 2.33). For the TDI system in the proof of Theorem 4.38, every integer vector \mathbf{z} on a face satisfies some of the inequalities of $A\mathbf{x} \leq \mathbf{b}$ with equality. Due to (4.1), the associated components of \mathbf{b} must be integers. Conversely, Theorem 4.31 implies that integrality of \mathbf{b} implies the integrality of \mathcal{P} . We obtain the following characterization.

Theorem 4.39. A rational polyhedron $\mathcal{P} \subseteq \mathbb{R}^n$ is integral if and only if there is a TDI system $A\mathbf{x} \leq \mathbf{b}$ with $A \in \mathbb{Q}^{m \times n}$ and $\mathbf{b} \in \mathbb{Z}^m$ and $\mathcal{P} = \mathcal{P}(A, \mathbf{b})$.

4.3.1 Applications in combinatorial optimization

There are many examples where TDI systems naturally occur, two of which are presented below.

Example 4.40. A matrix $A \in \{0, 1\}^{m \times n}$ is called *balanced* if it does not have a $(k \times k)$ -submatrix for odd k and with exactly two ones per row and column. For example, the matrix

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

is balanced, but not totally unimodular (its determinant is -2). It can be shown (see [16]) that A being balanced implies that the system $A\mathbf{x} \leq \mathbf{1}, \mathbf{x} \geq \mathbf{0}$ is TDI. \triangle

Often, TDI systems are used to prove discrete min-max results, such as Theorem 4.14.

Example 4.41. Let $G = (V, E)$ be a directed graph and $r \in V$. An *r-arborescence* is a subset $A \subseteq E$ with $|A| = |V| - 1$ such that every vertex is reachable from r in (V, A) , i.e., in particular, $|\delta^-(v) \cap A| = 1$ for all $v \in V \setminus \{r\}$ and $|\delta^-(r) \cap A| = 0$. An *r-cut* is a directed cut $\delta^-(U)$ in G with $\emptyset \neq U \subseteq V \setminus \{r\}$.

Let M be the 0/1-matrix whose rows are the incidence vectors of all *r-cuts* in G . It can be shown (see [14]) that the system $\{M\mathbf{x} \geq \mathbf{1}, \mathbf{x} \geq \mathbf{0}\}$ is TDI. The minimal $\mathbf{x} \in \{0, 1\}^A$ that satisfy $M\mathbf{x} \geq \mathbf{1}$ are exactly the *r-arborescences*.

With this result, we can handily prove Fulkerson's optimal arborescence theorem (see [17]): Let $\mathbf{c} \in \mathbb{Z}_{\geq 0}^m$ be a non-negative integer cost vector. Because $\{\mathbf{x} \in \mathbb{R}^m : M\mathbf{x} \geq \mathbf{1}, \mathbf{x} \geq \mathbf{0}\}$ is TDI and because of LP duality, we get that

$$\min\{\mathbf{c}^\top \mathbf{x} : M\mathbf{x} \geq \mathbf{1}, \mathbf{x} \geq \mathbf{0}\} = \max\{\mathbf{1}^\top \mathbf{y} : M^\top \mathbf{y} \leq \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}$$

and both optima are assumed by integral points x^* and y^* , see Theorem 4.31. This means that the minimum cost of an r -arborescence is equal to the maximum size of a set of r -cuts such that every $e \in A$ is contained in at most c_a r cuts. \triangle

Remark 4.42. It can be checked in polynomial time whether a matrix is totally unimodular (see [40]). The same applies to recognizing balanced matrices (see [11]). On the other hand, it is NP-complete to decide whether a system is TDI, even for $Ax \leq \mathbf{1}$, $A \in \{0, 1\}^{m \times n}$ (see [13]).

5 Cutting Planes

We consider a general mixed-integer program as introduced in Definition 1.1, i.e.,

$$\begin{aligned} \max \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} \leq \mathbf{b}, \\ & \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}, \end{aligned} \tag{5.1}$$

with $A \in \mathbb{Q}^{m \times n}$, $\mathbf{c} \in \mathbb{Q}^n$, $\mathbf{b} \in \mathbb{Q}^m$ and $p \in \{0, \dots, n\}$. Note that, to ensure that the integer hull is a polyhedron, we restrict ourselves to rational data (see Chapter 2).

In Chapter 4 we investigated sufficient conditions for the LP-relaxation of (5.1) to have integral optima, and Chapter 3 we saw how to treat the general case via branch-and-bound. In this chapter, we investigate how to strengthen the LP-relaxation by adding additional constraints that shrink the feasible region of the LP without excluding feasible solutions of the MIP.

In particular, if a computed optimum vertex solution \mathbf{x}^* of the LP-relaxation does not fulfill the integer conditions (i.e., $\mathbf{x}^* \notin \mathbb{Z}^p \times \mathbb{R}^{n-p}$), there must be an inequality (since \mathcal{P}_1 is a polyhedron by Theorem 2.30), called *cutting plane*, that separates \mathbf{x}^* from $\mathcal{P}_1 = \text{conv}(\{\mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p} : A\mathbf{x} \leq \mathbf{b}\})$ – finding such an inequality amounts to solving a separation problem (see *Introduction to Optimization*). If we find such a separating inequality, we can add it to strengthen the LP relaxation and repeat the process for as long as the LP relaxation yields infeasible optima. This procedure is often referred to as the *cutting plane method*.

In this chapter, we will devise strategies for generating cutting planes that, under suitable conditions, provably produce a feasible optimum solution. These approaches generally work for arbitrary IPs. Because (5.1) is NP-hard to solve (Theorem 3.1) and the separation problem is polynomially equivalent to the optimization problem (see *Introduction to Optimization*), we cannot hope for efficient procedures that are always successful (assuming $P \neq NP$).

We will also examine the generation of cuts that can be used to accelerate the branch-and-bound method. Such cutting planes often exploit the underlying problem structure.

5.1 General cutting planes

In this section, we consider a class of inequalities that is valid for \mathcal{P}_1 and which can be applied irrespective of the specific problem structure. We will see that this class is rich enough to provide a complete description of \mathcal{P}_1 . We begin by introducing these inequalities from a geometric perspective and later assume a more algorithmic approach.

5.1.1 Geometric approach: Chvátal-Gomory inequalities

We consider a rational polyhedron $\mathcal{P} := \mathcal{P}(A, \mathbf{b})$, with $A \in \mathbb{Q}^{m \times n}$, $\mathbf{b} \in \mathbb{Q}^m$. Our goal is to describe $\mathcal{P}_1 := \text{conv}(\mathcal{P} \cap \mathbb{Z}^n)$ by linear inequalities.

Let $\mathbf{c}^\top \mathbf{x} \leq \delta$ with integer vector $\mathbf{c} \in \mathbb{Z}^n$ be a *valid inequality*, i.e., $\mathcal{P} \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} \leq \delta\}$ (and $\mathcal{P}_1 \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} \leq \delta\}$). Since $\mathbf{c}^\top \mathbf{x} \in \mathbb{Z}$ for all $\mathbf{x} \in \mathcal{P} \cap \mathbb{Z}^n$, it follows that

$$\mathcal{P}_1 \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} \leq \lfloor \delta \rfloor\}.$$

This observation suggests taking all valid inequalities with an integer normal vector and rounding down the right-hand sides in order to obtain a stronger formulation. Of course, it suffices to take supporting inequalities, i.e., those for which $\mathcal{P} \cap \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} = \delta\} \neq \emptyset$. To do this, define

$$\mathcal{P}^1 := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} \leq \lfloor \delta \rfloor \text{ for all valid supporting inequalities } \mathbf{c}^\top \mathbf{x} \leq \delta \text{ of } \mathcal{P} \text{ with } \mathbf{c} \in \mathbb{Z}^n\}. \quad (5.2)$$

The resulting inequalities $\mathbf{c}^\top \mathbf{x} \leq \lfloor \delta \rfloor$ are called *Chvátal-Gomory inequalities* and \mathcal{P}^1 is called *elementary closure* of \mathcal{P} .

At first glance, it is not clear whether \mathcal{P}^1 is a polyhedron again, because there are infinitely many supporting hyperplanes. However, we will prove that \mathcal{P}^1 is indeed a polyhedron. This allows us to iterate and define

$$\mathcal{P}^0 := \mathcal{P} \text{ and } \mathcal{P}^{t+1} := (\mathcal{P}^t)^1 \text{ for } t \in \mathbb{N}. \quad (5.3)$$

Then \mathcal{P}^t is the *Chvátal-Gomory closure of rank t* . We immediately have

$$\mathcal{P} = \mathcal{P}^0 \supseteq \mathcal{P}^1 \supseteq \dots \supseteq \mathcal{P}_1. \quad (5.4)$$

This raises the question of whether this method is finite. We will show that, indeed, a $t \in \mathbb{N}$ exists with $\mathcal{P}^t = \mathcal{P}_1$. The resulting \mathcal{P}^t provides the desired description of \mathcal{P}_1 by linear inequalities.

Remark 5.1. By rounding the right-hand side, we move the hyperplane $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} = \delta\}$ for $\mathbf{c} \in \mathbb{Z}^n$ closer to \mathcal{P}_1 . If $\gcd(c_1, \dots, c_n) = 1$, the resulting hyperplane then contains an integer point: Otherwise, according to Lemma 4.25, we would have $\lambda \in \mathbb{Q}$ with $\lambda \mathbf{c} \in \mathbb{Z}^n$ but $\lambda \lfloor \delta \rfloor \notin \mathbb{Z}$. Since $\gcd(c_1, \dots, c_n) = 1$ and $\lambda \mathbf{c} \in \mathbb{Z}^n$, it follows that $\lambda \in \mathbb{Z}$, a contradiction.

If the resulting hyperplanes define supporting hyperplanes for \mathcal{P}_1 we have reached $\mathcal{P} = \mathcal{P}_1$, since, by Lemma 2.33, we have the following.

Observation 5.2. Every minimal face of a polyhedron $\mathcal{P} \subseteq \mathbb{R}^n$ contains an integer point if and only if every supporting hyperplane of \mathcal{P} contains an integer point of \mathcal{P} .

We start with the proof that \mathcal{P}^1 is a polyhedron. According to Theorem 4.38, every rational polyhedron can be represented by a TDI system of the form $A\mathbf{x} \leq \mathbf{b}$ with integer matrix A . This allows to calculate \mathcal{P}^1 directly.

Theorem 5.3. Let $\mathcal{P} = \mathcal{P}(A, \mathbf{b}) \neq \emptyset$ with $A\mathbf{x} \leq \mathbf{b}$ TDI and $A \in \mathbb{Z}^{m \times n}$. Then, $\mathcal{P}^1 = \mathcal{P}(A, \lfloor \mathbf{b} \rfloor)$, i.e., \mathcal{P}^1 is a polyhedron.^a

^a $\lfloor \mathbf{b} \rfloor$ is to be understood component-wise.

Proof. We have $\mathcal{P}^1 \subseteq \mathcal{P}(A, \lfloor \mathbf{b} \rfloor)$, since $A_i \cdot \mathbf{x} \leq \lfloor b_i \rfloor$ is a Chvátal-Gomory inequality.

To prove the converse, let $\mathbf{c}^\top \mathbf{x} \leq \delta$ be a valid supporting inequality for \mathcal{P} and \mathbf{c} be integral. By strong duality, we have

$$\delta = \max \{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}\} = \min \{\mathbf{y}^\top \mathbf{b} : A^\top \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}.$$

Since $Ax \leq \mathbf{b}$ is TDI, there is an integral optimum solution \mathbf{y}^* of the dual problem, i.e., $\mathbf{c}^\top \mathbf{x} \leq \delta$ can be expressed as integral linear combination of the rows of $Ax \leq \mathbf{b}$. For $\hat{\mathbf{x}} \in \mathcal{P}(A, \lfloor \mathbf{b} \rfloor)$, we have

$$\mathbf{c}^\top \hat{\mathbf{x}} = (A^\top \mathbf{y}^*)^\top \hat{\mathbf{x}} = (\mathbf{y}^*)^\top (A\hat{\mathbf{x}}) \stackrel{\mathbf{y}^* \geq 0}{\leq} (\mathbf{y}^*)^\top \lfloor \mathbf{b} \rfloor = \lfloor (\mathbf{y}^*)^\top \mathbf{b} \rfloor \leq \lfloor (\mathbf{y}^*)^\top \mathbf{b} \rfloor = \lfloor \delta \rfloor.$$

It follows that $\mathcal{P}(A, \lfloor \mathbf{b} \rfloor) \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} \leq \lfloor \delta \rfloor\}$. Let $S \subseteq \mathbb{Z}^n \times \mathbb{R}$ denote the set of all pairs (\mathbf{c}, δ) such that $\mathbf{c}^\top \mathbf{x} \leq \delta$ is a valid supporting inequality. We obtain

$$\mathcal{P}(A, \lfloor \mathbf{b} \rfloor) \subseteq \bigcap_{(\mathbf{c}, \delta) \in S} \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} \leq \lfloor \delta \rfloor\} = \mathcal{P}^1. \quad \square$$

We show the useful fact that the operations of taking a face and taking the elementary closure of a polyhedron commute.

Corollary 5.4. For every face \mathcal{F} of \mathcal{P} it holds that $\mathcal{F}^1 = \mathcal{P}^1 \cap \mathcal{F}$.

Proof. By Theorem 4.38, we may let $\mathcal{P} = \mathcal{P}(A, \mathbf{b})$ with $Ax \leq \mathbf{b}$ TDI and A integral. Let \mathcal{F} be a face of \mathcal{P} and let $\mathbf{c} \in \mathbb{Z}^n$, $\delta \in \mathbb{Z}$ such that $\mathcal{F} = \{\mathbf{x} \in \mathcal{P} : \mathbf{c}^\top \mathbf{x} = \delta\}$ and $\mathbf{c}^\top \mathbf{x} \leq \delta$ is valid for \mathcal{P} . Theorem 4.35 implies that the system $Ax \leq \mathbf{b}$, $\mathbf{c}^\top \mathbf{x} \leq \delta$ is also TDI, since the corresponding polyhedron has no additional faces. Therefore, by Corollary 4.36, the system $Ax \leq \mathbf{b}$, $\mathbf{c}^\top \mathbf{x} \leq \delta$, $-\mathbf{c}^\top \mathbf{x} \leq -\delta$ is TDI as well. Since δ is an integer, Theorem 5.3 yields

$$\begin{aligned} \mathcal{P}^1 \cap \mathcal{F} &= \mathcal{P}^1 \cap \mathcal{P} \cap \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} = \delta\} \\ &\stackrel{\text{Thm 5.3}}{=} \{\mathbf{x} \in \mathbb{R}^n : Ax \leq \lfloor \mathbf{b} \rfloor, \mathbf{c}^\top \mathbf{x} = \delta\} \\ &\stackrel{\delta \in \mathbb{Z}}{=} \{\mathbf{x} \in \mathbb{R}^n : Ax \leq \lfloor \mathbf{b} \rfloor, \mathbf{c}^\top \mathbf{x} \leq \lfloor \delta \rfloor, -\mathbf{c}^\top \mathbf{x} \leq \lfloor -\delta \rfloor\} \\ &\stackrel{\text{Thm 5.3}}{=} \mathcal{F}^1. \end{aligned} \quad \square$$

Observation 5.5.

- (a) \mathcal{F}^1 is a (possibly empty) face of \mathcal{P}^1 , since $\mathcal{F}^1 = \mathcal{P}^1 \cap \mathcal{F} = \mathcal{P}^1 \cap \mathcal{P} \cap \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} = \delta\} = \mathcal{P}^1 \cap \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} = \delta\}$ (note, that $\mathcal{P}^1 \subseteq \mathcal{P} \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} \leq \delta\}$).
- (b) Therefore, $\mathcal{F}^2 = (\mathcal{F}^1)^1 = (\mathcal{P}^1)^1 \cap \mathcal{F}^1 = \mathcal{P}^2 \cap \mathcal{P}^1 \cap \mathcal{F} = \mathcal{P}^2 \cap \mathcal{F}$. Repeated application yields

$$\mathcal{F}^t = \mathcal{P}^t \cap \mathcal{P}^{t-1} \cap \dots \cap \mathcal{P}^1 \cap \mathcal{F} = \mathcal{P}^t \cap \mathcal{F} \quad \forall t \in \mathbb{N}.$$

We now have everything in place to show the finiteness of the rounding scheme.

Theorem 5.6 ([9]). For every rational polyhedron \mathcal{P} there exists $t \in \mathbb{Z}_{\geq 0}$ with $\mathcal{P}^t = \mathcal{P}_1$.

We first sketch the proof of Theorem 5.6. It uses induction on the dimension of the polyhedron. We show, that it suffices to consider the progress of supporting hyperplanes of a finite number of orientations – namely all normal directions needed for a complete description of the polyhedra \mathcal{P} and \mathcal{P}_1 . For every fixed direction, we argue that we can make progress until the corresponding supporting hyperplane contains an integer point of \mathcal{P} : As long as this is not the case, we know by induction that, after a finite number of iterations, the Chvátal-Gomory closure no longer contains points of this supporting hyperplane. Eventually, all relevant supporting hyperplanes contain an integer point of \mathcal{P} and thus we have found a description of \mathcal{P}_1 (cf. Observation 5.2).

Proof of Theorem 5.6. If $\mathcal{P} = \emptyset$, then $\mathcal{P}^0 = \mathcal{P} = \mathcal{P}_1$ follows. We may therefore assume $\mathcal{P} \neq \emptyset$ in the following. We further prove that we may assume \mathcal{P} to be full-dimensional. To that end, let $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\} := \text{aff}(\mathcal{P}) \neq \mathbb{R}^n$ be the affine hull of \mathcal{P} , i.e., $\mathcal{P} \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}$. We can assume that A is an integer matrix of full row rank (otherwise we scale and eliminate redundant rows), i.e., $A \in \mathbb{Z}^{(n-d) \times n}$ with $d := \dim(\mathcal{P})$.

If $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ contains no integer point, then, according to Lemma 4.25, there is $\mathbf{y} \in \mathbb{Q}^{n-d}$ with $\mathbf{c} := A^\top \mathbf{y} \in \mathbb{Z}^n$ and $\delta := \mathbf{b}^\top \mathbf{y} \notin \mathbb{Z}$. Every $\hat{\mathbf{x}} \in \mathcal{P}$ satisfies $A\hat{\mathbf{x}} = \mathbf{b}$ and therefore $\mathbf{c}^\top \hat{\mathbf{x}} = (A^\top \mathbf{y})^\top \hat{\mathbf{x}} = \mathbf{y}^\top A\hat{\mathbf{x}} = \mathbf{y}^\top \mathbf{b} = \delta$. Thus $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} = \delta\}$ is a supporting hyperplane of \mathcal{P} . We conclude that

$$\begin{aligned} \mathcal{P}_1 &\subseteq \mathcal{P}^1 \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} \leq \lfloor \delta \rfloor, -\mathbf{c}^\top \mathbf{x} \leq \lfloor -\delta \rfloor\} \\ &= \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} \leq \lfloor \delta \rfloor, \mathbf{c}^\top \mathbf{x} \geq \lceil \delta \rceil\} \stackrel{\delta \notin \mathbb{Z}}{=} \emptyset, \end{aligned}$$

and thus $\mathcal{P}_1 = \mathcal{P}^1 = \emptyset$.

Now let $\hat{\mathbf{x}} \in \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ be integral. Theorem 5.6 is invariant under translation by the vector $\hat{\mathbf{x}}$, thus we can assume $\text{aff}(\mathcal{P}) = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{0}\}$. By Theorem 4.21 and Remark 4.24 (b), as well as the fact that the row rank of A is $n - d$, a square unimodular matrix U with $AU = [B, 0]$ in Hermite normal form exists. By Corollary 4.6, U^{-1} is unimodular as well and $U\mathbb{Z}^n = \mathbb{Z}^n$. Via the (bijective) variable transformation $\mathbf{x} = U\mathbf{z}$, we can therefore assume that

$$\text{aff}(\mathcal{P}) = \{\mathbf{z} \in \mathbb{R}^n : [B, 0]\mathbf{z} = \mathbf{0}\} = \{0\}^{n-d} \times \mathbb{R}^d.$$

Every supporting hyperplane $\mathcal{H} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}^\top \mathbf{x} = \delta\}$ of \mathcal{P} can be transformed into the form

$$\mathcal{H}' = \{\mathbf{x} \in \mathbb{R}^n : \sum_{i=1}^{n-d} 0x_i + \sum_{i=n-d+1}^n c_i x_i = \delta\}$$

by adding suitable multiples of the rows of $[B, 0]$ to \mathbf{c} (this is possible, because $[B, 0]\mathbf{x} = \mathbf{0}$ for $\mathbf{x} \in \mathcal{P} \subseteq \text{aff}(\mathcal{P})$; δ and c_{n-d+1}, \dots, c_n are unaffected). In the construction of \mathcal{P}^1 we can restrict ourselves to support hyperplanes of the form \mathcal{H}' . We may therefore assume $n - d = 0$, i.e., \mathcal{P} is full dimensional.

Towards the proof of the theorem, we employ induction over the dimension d of \mathcal{P} . If $d = 0$, then $\mathcal{P} = \{\bar{\mathbf{x}}\}$ for some $\bar{\mathbf{x}} \in \mathbb{R}^n$. If $\bar{\mathbf{x}} \in \mathbb{Z}^n$, we obtain $\mathcal{P}^0 = \mathcal{P} = \mathcal{P}_1$. Otherwise, $\mathcal{P}^1 = \emptyset = \mathcal{P}_1$. Now, consider $d > 0$ and let the theorem hold for all polyhedra of smaller dimensions.

By Theorem 2.30, by Corollary 2.21, and by scaling, we can find $W \in \mathbb{Z}^{m \times n}$ and $\mathbf{w} \in \mathbb{Z}^m$ with $\mathcal{P}_1 = \mathcal{P}(W, \mathbf{w})$. We may further assume that $\mathbf{W}_i \neq \mathbf{0}^\top$ and $w'_i := \max\{\mathbf{W}_i \mathbf{x} : \mathbf{x} \in \mathcal{P}\}$ is bounded for all $i \in \{1, \dots, m\}$: If $\mathcal{P}_1 \neq \emptyset$, this holds since zero-rows are redundant and by Corollary 2.31. Now assume $\mathcal{P}_1 = \emptyset$. Since $\mathcal{P} \neq \emptyset$ is rational, Corollary 2.21 yields a representation $\mathcal{P} = \mathcal{V} + \mathcal{E}$. We have $\dim(\mathcal{E}) < n$, otherwise \mathcal{P} would contain integral points (exercise). Hence, there is $\mathbf{d} \in \mathbb{Z}^n$ with $\mathbf{d}^\top \mathbf{r} = 0$ for all $\mathbf{r} \in \mathcal{E}$, and both $\max\{\mathbf{d}^\top \mathbf{x} : \mathbf{x} \in \mathcal{P}\} = \max\{\mathbf{d}^\top \mathbf{x} : \mathbf{x} \in \mathcal{V}\}$ and $\max\{-\mathbf{d}^\top \mathbf{x} : \mathbf{x} \in \mathcal{P}\} = -\min\{\mathbf{d}^\top \mathbf{x} : \mathbf{x} \in \mathcal{V}\}$ are bounded. We can thus set $W := \begin{pmatrix} \mathbf{d}^\top \\ -\mathbf{d}^\top \end{pmatrix}$ and $\mathbf{w} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ with $\mathcal{P}_1 = \mathcal{P}(W, \mathbf{w}) = \emptyset$. In either case, we can find \mathbf{w}' with $\mathcal{P} \subseteq \mathcal{P}(W, \mathbf{w}')$.

Now, consider an inequality $\mathbf{a}^\top \mathbf{x} \leq \beta$ of the system $W\mathbf{x} \leq \mathbf{w}$ with $\mathbf{a} \in \mathbb{Z}^n \setminus \{\mathbf{0}\}$ and $\beta \in \mathbb{Z}$ and let $\mathbf{a}^\top \mathbf{x} \leq \beta'$ be the corresponding inequality of $W\mathbf{x} \leq \mathbf{w}'$ (feasible for \mathcal{P}). We claim that we can find $s \in \mathbb{N}$ with $\mathcal{P}^s \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} \leq \beta\}$. Suppose this was not the case, i.e., we cannot find any Chvátal-Gomory inequality that implies $\mathbf{a}^\top \mathbf{x} \leq \beta$. Then, $\mathcal{P}^1 \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} \leq \lfloor \beta' \rfloor\}$, but, by assumption, $\mathcal{P}^1 \not\subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} \leq \beta\}$ holds. It follows that $\beta < \lfloor \beta' \rfloor$. Since \mathbf{w} and therefore β is integral, there is $\beta'' \in \mathbb{Z}$ and $r \in \mathbb{Z}_{\geq 0}$ with

$$\begin{aligned} \beta &< \beta'' \leq \lfloor \beta' \rfloor, \\ \mathcal{P}^t &\subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} \leq \beta''\} \text{ for all } t \geq r, \text{ and} \\ \mathcal{P}^t &\not\subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} \leq \beta'' - 1\} \text{ for all } t \geq r. \end{aligned}$$

Due to the choice of r , we have that $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} = \beta''\}$ is a supporting hyperplane for \mathcal{P}^r . Since $\mathbf{a} \neq \mathbf{0}$, we have that $\mathcal{F} := \mathcal{P}^r \cap \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} = \beta''\}$ has lower dimension than n . We further conclude from $\mathcal{P}_1 \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} \leq \beta\}$ and $\beta'' > \beta$ that $\mathcal{F} \cap \mathbb{Z}^n = \emptyset$. Therefore, by the induction hypothesis, there is $\tilde{r} \in \mathbb{Z}_{\geq 0}$ with $\mathcal{F}^{\tilde{r}} = \mathcal{F}_1 = \emptyset$. By Remark 5.5 (b),

$$\emptyset = \mathcal{F}^{r+\tilde{r}} = \mathcal{P}^{r+\tilde{r}} \cap \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} = \beta''\}.$$

Therefore, $\mathcal{P}^{r+\tilde{r}} \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} < \beta''\}$, which means that $\mathcal{P}^{r+\tilde{r}+1} \subseteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} \leq \beta'' - 1\}$, which is a contradiction to the choice of β'' and r .

Since $\mathbf{a}^\top \mathbf{x} \leq \beta$ was chosen arbitrarily and the system $W\mathbf{x} \leq \mathbf{w}$ is finite, we can thus find $s \in \mathbb{Z}_{\geq 0}$ with $\mathcal{P}^s \subseteq \mathcal{P}(W, \mathbf{w}) = \mathcal{P}_1$. It follows that $\mathcal{P}^s = \mathcal{P}_1$ as desired (see (5.4)). \square

Let us review what we have shown so far. The procedure to obtain a linear description of the integral polyhedron $\mathcal{P}_1 = \text{conv}(\{\mathbf{x} \in \mathbb{Z}^n : A\mathbf{x} \leq \mathbf{b}\})$ proceeds as follows. We start with the linear relaxation $\mathcal{P}^0 = \mathcal{P} = \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}\}$ of \mathcal{P}_1 . Next, we consider every supporting hyperplane of \mathcal{P} whose left-hand side has integer coefficients, and round the right-hand side down to the next integer. This operation, for all such supporting hyperplanes results in the polyhedron \mathcal{P}^1 . Theorem 5.3 shows that there is no need to consider all supporting hyperplanes – all we need is a TDI system $D\mathbf{x} \leq \mathbf{d}$ that describes \mathcal{P} . For this TDI system, we need to round down the vector \mathbf{d} of right-hand sides. In the proof of Theorem 4.38, we explicitly constructed the TDI system for a rational polyhedron \mathcal{P} by explicitly constructing, for every face \mathcal{F} of \mathcal{P} a generating Hilbert basis of $\text{cone}((A_{\text{eq}(\mathcal{F})})^\top)$. Overall, the method for constructing a Hilbert basis for a polyhedral cone yields an algorithm for computing \mathcal{P}^1 . According to Theorem 5.6, this algorithm only needs to be performed a finite number of times in order to obtain a linear description of \mathcal{P}_1 .

The entire procedure is hardly practicable. First of all, the number $t \in \mathbb{N}$ in Theorem 5.6 is possibly exponential in the coding length of the input variables A, \mathbf{b} (see exercise). Secondly, in every iteration, we have to determine Hilbert bases for all cones generated by the faces of \mathcal{P}^{i-1} . In general, not only is the number of faces exponential, we cannot even compute a single Hilbert basis in polynomial time.

On the other hand, a complete description of $\mathcal{P}_1 = \text{conv}(\{\mathbf{x} \in \mathbb{Z}^n : A\mathbf{x} \leq \mathbf{b}\})$ is not necessary for solving (5.1). We “merely” need to find an optimum solution. Put another way, we are only interested in one face, namely the one that contains the optimum solutions. Even for this face, it is not necessary to find a Hilbert basis for the corresponding cone. We only need to be able to find integral vectors in this cone.

More formally, assume we have the solution of the LP relaxation $\max\{\mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$. Let \mathbf{x}^* be an optimum solution that is not integral. The challenge is to find an integral vector in $\text{cone}((A_{\text{eq}(\mathbf{x}^*)})^\top)$ that separates the current fractional optimum solution \mathbf{x}^* , i.e., to find $\mathbf{d} \in \mathbb{Z}^n \cap \text{cone}((A_{\text{eq}(\mathbf{x}^*)})^\top)$ with $\mathbf{d}^\top \mathbf{x}^* \notin \mathbb{Z}$.

5.1.2 Algorithmic approach: Gomory inequalities

We can systematically find cutting planes by evaluating information gained during the course of the simplex algorithm. We first consider the purely integral case ($p = n$) and assume that all coefficients are integers. We can transform (5.1) into standard form by splitting \mathbf{x} into $\mathbf{x}^+, \mathbf{x}^- \geq \mathbf{0}$ with $\mathbf{x} = \mathbf{x}^+ - \mathbf{x}^-$ and introducing (integral) slack variables (see *Introduction to Optimization*). We obtain a problem of the form

$$\begin{aligned} \max \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0} \\ & \mathbf{x} \in \mathbb{Z}^n. \end{aligned} \tag{5.5}$$

with $A \in \mathbb{Z}^{m \times n}$ and $\mathbf{b} \in \mathbb{Z}^m$ and the LP relaxation

$$\begin{aligned} \max \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Solving the latter via the simplex method, we obtain an optimum vertex solution \mathbf{x}^* and an optimum basis B with $B \subseteq \{1, \dots, n\}$, $|B| = m$, and A_B regular. As usual, we denote the non-basis by $N = \{1, \dots, n\} \setminus B$. This means that $\mathbf{x}_N^* = \mathbf{0}$ and $\mathbf{x}_B^* = \bar{\mathbf{b}} - \bar{A}_N \mathbf{x}_N^* = A_B^{-1} \mathbf{b}$ and $\mathbf{x}_N^* = \mathbf{0}$ with $\bar{\mathbf{b}} := A_B^{-1} \mathbf{b}$ and $\bar{A} := A_B^{-1} A$. If \mathbf{x}_B^* is integral, we have found an optimum solution of (5.5). Otherwise, there exists $i \in B$ with $x_i^* \notin \mathbb{Z}$.

For every feasible solution \mathbf{x} of (5.5) it holds that $\mathbf{x}_B = \bar{\mathbf{b}} - \bar{A}_N \mathbf{x}_N$ or

$$x_i = \bar{b}_i - \sum_{j \in N} \bar{A}_{ij} x_j. \quad (5.6)$$

With $\mathbf{x}_N \geq \mathbf{0}$ we get

$$x_i \leq \bar{b}_i - \sum_{j \in N} \lfloor \bar{A}_{ij} \rfloor x_j.$$

Because $x_i \in \mathbb{Z}$ and $\lfloor \bar{A}_{ij} \rfloor x_j \in \mathbb{Z}$ for $j \in N$, we can round the right-hand side to obtain

$$x_i \leq \lfloor \bar{b}_i \rfloor - \sum_{j \in N} \lfloor \bar{A}_{ij} \rfloor x_j. \quad (5.7)$$

This valid inequality is called (*fractional*) *Gomory cut*. It cuts off \mathbf{x}^* from the current LP relaxation because $\mathbf{x}_N^* = \mathbf{0}$ and $x_i^* = \bar{b}_i \notin \mathbb{Z}$. Furthermore, (5.7) is an inequality of the form $\mathbf{c}^\top \mathbf{x} \leq \lfloor \delta \rfloor$, $\mathbf{c} \in \mathbb{Z}^n$ for a supporting hyperplane $\mathbf{c}^\top \mathbf{x} \leq \delta$, thus it appears in the approach of the previous section.

The inequality (5.7) is converted by means of a slack variable into an equation and added to the system $A\mathbf{x} = \mathbf{b}$, $\mathbf{x} \geq \mathbf{0}$, maintaining the property that all coefficients are integers. Therefore, the slack variable introduced for the new inequality can also be regarded as an integer and thus the procedure can be iterated. It can be shown that, with a certain choice of optimum solution and choice of fractional basis element, we obtain a finite algorithm (see [20]), i.e., after a finite number of inequalities have been added, an integral solution is found. The inequality (5.7) can be rewritten as follows. Let $f_i \in (0, 1)$ be the fractional part of \bar{b}_i , i.e. $\bar{b}_i = \lfloor \bar{b}_i \rfloor + f_i$. Analogously, let $f_{ij} \in [0, 1)$ be the fractional parts of \bar{A}_{ij} , i.e., $\bar{A}_{ij} = \lfloor \bar{A}_{ij} \rfloor + f_{ij}$. From (5.6) we thus obtain

$$x_i = f_i - \sum_{j \in N} f_{ij} x_j + \left(\lfloor \bar{b}_i \rfloor - \sum_{j \in N} \lfloor \bar{A}_{ij} \rfloor x_j \right).$$

If we subtract this equation from (5.7), we obtain the equivalent inequality

$$\sum_{j \in N} f_{ij} x_j \geq f_i. \quad (5.8)$$

Example 5.7. Consider the IP

$$\begin{aligned} \max \quad & x_2 \\ \text{s.t.} \quad & 4x_1 + x_2 \leq 4 \\ & -4x_1 + x_2 \leq 0 \\ & x_1, x_2 \geq 0 \\ & x_1, x_2 \in \mathbb{Z}. \end{aligned}$$

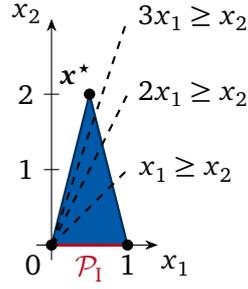


Figure 5.1: Illustration of Examples 5.7 and 5.9.

See Figure 5.1 for an illustration. Adding slack variables results in

$$\begin{aligned}
 \max \quad & x_2 \\
 \text{s.t.} \quad & 4x_1 + x_2 + x_3 = 4 \\
 & -4x_1 + x_2 + x_4 = 0 \\
 & x_1, x_2, x_3, x_4 \geq 0 \\
 & x_1, x_2, x_3, x_4 \in \mathbb{Z}.
 \end{aligned}$$

The optimal basis for the LP relaxation is $B = \{1, 2\}$ with

$$A_{.B} = \begin{pmatrix} 4 & 1 \\ -4 & 1 \end{pmatrix}, \quad A_{.N} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad A_{.B}^{-1} = \bar{A}_{.N} = \begin{pmatrix} 1/8 & -1/8 \\ 1/2 & 1/2 \end{pmatrix}.$$

Hence, $\mathbf{x}_B^* = \bar{\mathbf{b}} = (A_{.B})^{-1} \mathbf{b} = \begin{pmatrix} 1/2 \\ 2 \end{pmatrix}$. The only fractional component is x_1 and for $i = 1$ (5.6) becomes

$$x_1 = \frac{1}{2} - \frac{1}{8}x_3 + \frac{1}{8}x_4 \in \mathbb{Z},$$

which is true for all integral solutions \mathbf{x} . Inequality (5.7) then becomes

$$x_1 \leq \lfloor \frac{1}{2} \rfloor - \lfloor \frac{1}{8} \rfloor x_3 - \lfloor -\frac{1}{8} \rfloor x_4 = 0 - 0x_3 + x_4,$$

which reduces to $x_1 - x_4 \leq 0$. With respect to the original variables x_1, x_2 , by replacing the slack variable $x_4 = 4x_1 - x_2$, we obtain

$$-3x_1 + x_2 \leq 0.$$

This is a valid inequality for

$$\mathcal{P}_1 = \text{conv}(\{\mathbf{x} \in \mathbb{Z}_{\geq 0}^2 : 4x_1 + x_2 \leq 4, -4x_1 + x_2 \leq 0\}).$$

The alternative representation (5.8) is, with $x_3 = 4 - 4x_1 - x_2$,

$$\frac{1}{8}x_3 + \frac{7}{8}x_4 \geq \frac{1}{2} \iff \frac{1}{2} - \frac{1}{2}x_1 - \frac{1}{8}x_2 + \frac{7}{2}x_1 - \frac{7}{8}x_2 \geq \frac{1}{2} \iff 3x_1 - x_2 \geq 0. \quad \triangle$$

5.1.3 Gomory's mixed-integer cuts

The fractional Gomory cuts presented in the last section are dominated by the so-called mixed-integer Gomory cuts, which we now derive. We need the following observation.

Observation 5.8. Let $(\mathbf{a}^{(k)})^\top \mathbf{x} \geq \beta^{(k)}$ be a valid inequality for the polyhedron $\mathcal{P}^{(k)} \subseteq \mathbb{R}_{\geq 0}^n$ for $k \in \{1, 2\}$. Then,

$$\sum_{i=1}^n \max(a_i^{(1)}, a_i^{(2)}) x_i \geq \min(\beta^{(1)}, \beta^{(2)})$$

is valid for $\mathcal{P}^{(1)} \cup \mathcal{P}^{(2)}$ and $\text{conv}(\mathcal{P}^{(1)} \cup \mathcal{P}^{(2)})$.

As in the previous section, we start with (5.6) and split the sum on the right to obtain

$$x_i = \bar{b}_i - \sum_{j \in N: f_{ij} \leq f_i} \bar{A}_{ij} x_j - \sum_{j \in N: f_{ij} > f_i} \bar{A}_{ij} x_j, \quad (5.9)$$

where we again use the fractional values $f_i \in (0, 1)$ of \bar{b}_i (since $x_i^* \notin \mathbb{Z}$) and $f_{ij} \in [0, 1)$ of \bar{A}_{ij} . For j with $f_{ij} \leq f_i$ we use $\bar{A}_{ij} = \lfloor \bar{A}_{ij} \rfloor + f_{ij}$. For j with $f_{ij} > f_i$, and thus $f_{ij} > 0$, we use $\bar{A}_{ij} = \lfloor \bar{A}_{ij} \rfloor - 1 + f_{ij}$. Inserted in (5.9) these yield that all feasible \mathbf{x} satisfy

$$\sum_{j \in N: f_{ij} \leq f_i} f_{ij} x_j + \sum_{j \in N: f_{ij} > f_i} (f_{ij} - 1) x_j = r + f_i,$$

where

$$r = \lfloor \bar{b}_i \rfloor - \sum_{f_{ij} \leq f_i} \lfloor \bar{A}_{ij} \rfloor x_j - \sum_{f_{ij} > f_i} \lfloor \bar{A}_{ij} \rfloor x_j - x_i \in \mathbb{Z}.$$

Now, either $r \geq 0$ or $r \leq -1$ must hold. Therefore, every feasible point $\mathbf{x} \in \mathbb{Z}_{\geq 0}^n$ with $A\mathbf{x} = \mathbf{b}$ either satisfies the inequality

$$\sum_{j \in N: f_{ij} \leq f_i} \frac{f_{ij}}{f_i} x_j - \sum_{j \in N: f_{ij} > f_i} \frac{1-f_{ij}}{f_i} x_j \geq 1$$

or

$$- \sum_{j \in N: f_{ij} \leq f_i} \frac{f_{ij}}{1-f_i} x_j + \sum_{j \in N: f_{ij} > f_i} \frac{1-f_{ij}}{1-f_i} x_j \geq 1.$$

According to Observation 5.8, we can combine both inequalities to obtain the *Gomory mixed-integer cut*

$$\sum_{j \in N: f_{ij} \leq f_i} \frac{f_{ij}}{f_i} x_j + \sum_{j \in N: f_{ij} > f_i} \frac{1-f_{ij}}{1-f_i} x_j \geq 1. \quad (5.10)$$

We have just shown that this cutting plane is valid for all feasible \mathbf{x} .

Because

$$\frac{1-f_{ij}}{1-f_i} < 1 < \frac{f_{ij}}{f_i} \quad \text{for } f_{ij} > f_i,$$

the cut (5.10) dominates the fractional Gomory cuts (5.8). In particular, (5.10) cuts off the current solution of the LP relaxation.

Example 5.9. We consider the same IP as in Example 5.7 and again get

$$x_1 = \frac{1}{2} - \frac{1}{8}x_3 + \frac{1}{8}x_4.$$

We have $f_{13} = \frac{1}{8}$, $f_{14} = \frac{7}{8}$, $f_1 = \frac{1}{2}$, i.e., the two cases $f_{13} \leq f_1$ and $f_{14} > f_1$ occur. We obtain the mixed-integer Gomory cut (5.10) as

$$\frac{1}{8}x_3 + \frac{1-\frac{7}{8}}{1-\frac{1}{2}}x_4 \geq 1 \quad \Leftrightarrow \quad \frac{1}{4}x_3 + \frac{1}{4}x_4 \geq 1 \quad \Leftrightarrow \quad x_3 + x_4 \geq 4.$$

Elimination of $x_3 = 4 - 4x_1 - x_2$ and $x_4 = 4x_1 - x_2$ results in

$$x_3 + x_4 \geq 4 \iff 4 - 4x_1 - x_2 + 4x_1 - x_2 \geq 4 \iff -2x_2 \geq 0 \iff x_2 \leq 0.$$

In this case, we directly reach the integer hull \mathcal{P}_1 (see Figure 5.1).

For the two inequalities of the disjunction $r \geq 0$ or $r \leq -1$ we get $x_3 - x_4 \geq 4$ or $-x_3 + x_4 \geq 4$ respectively. With regard to the original variables, this corresponds to $x_1 \leq 0$ and $x_1 \geq 1$. The mixed-integer Gomory cut is exactly valid for this disjunction: In the left part, the intersection of the corresponding inequality with the polyhedron is $\left\{\begin{pmatrix} 0 \\ 0 \end{pmatrix}\right\}$, in the right part it is $\left\{\begin{pmatrix} 1 \\ 0 \end{pmatrix}\right\}$. \triangle

The name ‘‘mixed-integer Gomory cut’’ indicates that these inequalities were developed for the case in which, in addition to integer variables, continuous variables may occur (i.e., $p < n$). If $I = \{1, \dots, p\} \cap N$ is the set of indices of integer variables in the non-basis N , the general form of the cut is

$$\sum_{j \in I: f_{ij} \leq f_i} \frac{f_{ij}}{f_i} z_j + \sum_{j \in I: f_{ij} > f_i} \frac{1-f_{ij}}{1-f_i} z_j + \sum_{j \in N \setminus I: \bar{A}_{ij} > 0} \frac{\bar{A}_{ij}}{f_i} z_j - \sum_{j \in N \setminus I: \bar{A}_{ij} < 0} \frac{\bar{A}_{ij}}{1-f_i} z_j \geq 1. \quad (5.11)$$

Remark 5.10. In the mixed-integer case, the ideas of the previous two sections do not work. Chvátal’s approach does not work because the right-hand side in (5.3) cannot be rounded down. Gomory’s approach fails because the argument for rounding (5.7) no longer holds.

It can be shown (see [19]) that an algorithm based on iterative addition of the inequalities (5.11) solves $\min \{\mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \mathcal{X}\}$ with $\mathcal{X} = \{\mathbf{x} \in \mathbb{Z}_{\geq 0}^p \times \mathbb{R}_{\geq 0}^{n-p} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ in a finite number of steps if $\mathbf{c}^\top \mathbf{x} \in \mathbb{Z}$ holds for all $\mathbf{x} \in \mathcal{X}$. In general, it is currently not known whether a finite cutting plane method exists.

Various other inequalities have been considered independently of a specific problem structure. Among these are for example, mixed-integer rounding cuts (MIR) (see [34]) and so-called lift-and-project cuts (see [3]).

5.2 Specialized cuts

After dealing with valid inequalities for general IPs and MIP’s let us focus our attention on a single or a small subset of such constraints, where a problem might exhibit some local structure. For example, all variables in a constraint might be binary, or a MIP may contain a network flow problem as a part of its constraints. We are looking for ways to obtain stronger inequalities by exploiting such local structures.

5.2.1 Inequalities for knapsack problems

Every row of $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ defines a knapsack problem (see Example 1.6). We first attempt to derive inequalities for this problem that can then be used for the entire problem. To simplify the representation, we restrict ourselves to binary problems, i.e., problems of the form

$$\begin{aligned} \max \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & \mathbf{x} \in \{0, 1\}^n, \end{aligned}$$

with $\mathbf{A} \in \mathbb{Q}^{m \times n}$, $\mathbf{b} \in \mathbb{Q}^m$, $\mathbf{c} \in \mathbb{Q}^n$.

Now, we select a row $\mathbf{a}^\top \mathbf{x} \leq \beta$ from $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ and consider the polytope

$$\mathcal{R} := \{\mathbf{x} \in [0, 1]^n : \mathbf{a}^\top \mathbf{x} \leq \beta\}.$$

We want to find valid inequalities for \mathcal{R}_1 . Because

$$\mathcal{R}_1 = \text{conv}(\{\mathbf{x} \in \{0, 1\}^n : \mathbf{a}^\top \mathbf{x} \leq \beta\}) \supseteq \text{conv}(\{\mathbf{x} \in \{0, 1\}^n : \mathbf{A}\mathbf{x} \leq \mathbf{b}\}),$$

these inequalities are also valid for the original problem.

To simplify presentation, we bring the knapsack problem into a specific format. Recall that the knapsack problem is NP-hard (see *Algorithmic Discrete Mathematics*).

Observation 5.11. We may assume that

- (a) \mathbf{a} and β are integral, which can be achieved by multiplication by the least common multiple of the denominators,
- (b) $\mathbf{a} \geq \mathbf{0}$, otherwise we replace x_j by $1 - x_j$ for all j with $a_j < 0$,
- (c) $0 \leq \beta < \mathbf{1}^\top \mathbf{a}$, otherwise the problem is not interesting, because either $\beta \geq a_1 + \dots + a_n$ and all $\mathbf{x} \in \{0, 1\}^n$ are feasible, or $\beta < 0$ and no \mathbf{x} is feasible because we ensured $\mathbf{a} \geq \mathbf{0}$ above.
- (d) $a_j \leq \beta$ for all $j \in \{1, \dots, n\}$, otherwise $x_j = 0$ in every feasible solution and x_j can be ignored,
- (e) $\mathbf{c} \geq \mathbf{0}$ when determining $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{a}^\top \mathbf{x} \leq \beta\}$, otherwise $x_j = 0$ for every optimum solution (because of $a \geq 0$) for all $j \in \{1, \dots, n\}$ with $c_j < 0$. Furthermore, we can assume by scaling that \mathbf{c} is integral.

Definition 5.12. An index set $C \subseteq \{1, \dots, n\}$ is a *cover*, if

$$a(C) := \sum_{j \in C} a_j > \beta.$$

The cover C is *minimal* if no $C' \subsetneq C$ is a cover. The corresponding inequality

$$x(C) := \sum_{j \in C} x_j \leq |C| - 1$$

is called *cover inequality*.

Note that the cover inequality is valid for \mathcal{R}_1 , because for integral values it demands that $x_j = 0$ for at least one $j \in C$.

Two questions arise: How strong are these inequalities? How to find a (minimal) cover? We first focus on the former question.

The potential of an inequality for accelerating the solution process is difficult to assess. Presumably, however, facets of \mathcal{R}_1 are particularly important, because once we know them, we can solve the integer optimization problem easily. So we could hope that cover inequalities define facets of \mathcal{R}_1 . It will turn out that this is not always the case. We can show the following, however.

Theorem 5.13 ([2, 22, 36, 43]). A minimal cover inequality $x(C) \leq |C| - 1$ defines a facet of the restricted knapsack polytope

$$\mathcal{R}_1(C) := \text{conv}(\{\mathbf{x} \in \{0, 1\}^C : \mathbf{a}_C^\top \mathbf{x} \leq \beta\}).$$

Proof. We first need to determine $\dim \mathcal{R}_1(C)$. According to Observation 5.11, we may assume $a_j \leq \beta$ for all $j \in C$ and thus the vectors \mathbf{e}^j are feasible for all $j \in C$. Together with 0 they form $|C| + 1$ affine independent vectors. Thus $\dim \mathcal{R}_1(C) = |C|$ follows.

For the cover inequality to define a facet, there must be a number of $\dim \mathcal{R}_1(C) = |C|$ affinely independent vectors of $\mathcal{R}_1(C)$ that satisfy it with equality. Vectors with this property are $\mathbf{1} - \mathbf{e}^j \in \mathbb{Z}^C$ for $j \in C$, which are in $\mathcal{R}_1(C)$, because C is a minimal cover. \square

As mentioned above, cover inequalities do not always define facets of the general knapsack polytope \mathcal{R}_1 . However, this can be addressed by *lifting* the inequalities. For this, we, in general, calculate the largest possible coefficient d_k for a variable x_k with coefficient 0 in the inequality, i.e., in our case $k \notin C$, and add it to the inequality to obtain

$$\sum_{j \in C} x_j + d_k x_k \leq |C| - 1. \quad (5.12)$$

The new coefficient d_k must be chosen such that the inequality is valid. However, it should also be as large as possible to strengthen the inequality as much as possible. The largest possible coefficient can be found by solving

$$d_k = |C| - 1 - \max\{\mathbf{1}^\top \mathbf{x} : \mathbf{a}_C^\top \mathbf{x} \leq \beta - a_k, \mathbf{x} \in \{0, 1\}^C\}. \quad (5.13)$$

If $x_k = 0$ holds for some point $\mathbf{x} \in \mathcal{R}_1$, then it obviously satisfies (5.12). If $x_k = 1$, then $\mathbf{a}_C^\top \mathbf{x} \leq \beta - a_k$ must hold. We optimize exactly over these constraints while demanding that (5.12) remains satisfied. Note that (5.13) again defines a knapsack problem.

Example 5.14. Let the inequality $\mathbf{a}^\top \mathbf{x} \leq \beta$ be

$$6x_1 + 5x_2 + 5x_3 + 5x_4 + 8x_5 \leq 16.$$

The set $C = \{1, 2, 3, 4\}$ forms a (minimal) cover because $a(C) = 21 > 16$. The corresponding cover inequality is $x(C) \leq 3$. We now consider the only candidate $k = 5$ and calculate

$$d_5 = 3 - \max\{\{x_1 + x_2 + x_3 + x_4 : 6x_1 + 5x_2 + 5x_3 + 5x_4 \leq 8, x_1, \dots, x_4 \in \{0, 1\}\}\}.$$

The maximum is 1 and therefore $d_5 = 2$. We obtain the lifted inequality

$$x_1 + x_2 + x_3 + x_4 + 2x_5 \leq 3. \quad \triangle$$

If we iterate the procedure described above, we obtain a set $C' \subseteq \{1, \dots, n\} \setminus C$ and the valid inequality

$$\sum_{j \in C} x_j + \sum_{j \in C'} d_j x_j \leq |C| - 1. \quad (5.14)$$

The lifting problem (5.13) changes for $k \notin C \cup C'$ to

$$d_k = |C| - 1 - \max\{\sum_{j \in C} x_j + \sum_{j \in C'} d_j x_j : \sum_{j \in C \cup C'} a_j x_j \leq \beta - a_k, \mathbf{x} \in \{0, 1\}^{C \cup C'}\}. \quad (5.15)$$

By construction, $d_k \in \{0, \dots, |C| - 1\}$ holds. We obtain the following theorem.

Theorem 5.15. The lifted inequality (5.14) defines a facet of $\mathcal{R}_1(C \cup C')$. If $C \cup C' = \{1, \dots, n\}$, then (5.14) defines a facet of \mathcal{R}_1 .

Proof. We proceed by induction over $|C'|$. The result follows for $|C'| = 0$ from Theorem 5.13. Now let k be the index that is to be added to C' . Because, by induction, (5.14) defines a facet of $\mathcal{R}_1(C \cup C')$, there are affinely independent vectors $\mathbf{v}^1, \dots, \mathbf{v}^t \in \{0, 1\}^{C \cup C'}$ of \mathcal{R}_1 satisfying it with equality, where $t = |C| + |C'|$ since $\mathcal{R}_1(C \cup C')$ is full-dimensional as in Theorem 5.13. Now we define vectors $\hat{\mathbf{v}}^1, \dots, \hat{\mathbf{v}}^t \in \{0, 1\}^{C \cup C' \cup \{k\}}$ via

$$\hat{v}_i^j := \begin{cases} 0 & \text{if } i = k, \\ v_i^j, & \text{otherwise,} \end{cases}$$

for $j' \in \{1, \dots, t\}$ and $i \in C \cup C' \cup \{k\}$. Then, $\hat{\mathbf{v}}^1, \dots, \hat{\mathbf{v}}^t$ each fulfill the lifted inequality

$$\sum_{j \in C} x_j + \sum_{j \in C'} d_j x_j + d_k x_k \leq |C| - 1 \quad (5.16)$$

with equality.

Now let \mathbf{x}^* be an optimum solution for the maximum in (5.15) when calculating d_k . We define the vector $\hat{\mathbf{x}}^*$ via

$$\hat{x}_i^* := \begin{cases} 1 & \text{if } i = k, \\ x_i^* & \text{otherwise,} \end{cases}$$

for $i \in C \cup C' \cup \{k\}$. By construction, $\hat{\mathbf{x}}^*$ satisfies the inequality (5.16) with equality. Therefore $\hat{\mathbf{v}}^1, \dots, \hat{\mathbf{v}}^t, \hat{\mathbf{x}}^*$ together form $t + 1$ affinely independent vectors, which proves the theorem because of $\dim \mathcal{R}_1(C \cup C' \cup \{k\}) = |C| + |C'| + 1$. \square

To summarize, from minimal covers we can, step-by-step, determine facets of \mathcal{R}_1 . However, we have not yet discussed how to find minimal covers or how to solve the lifting problems. Unfortunately, solving knapsack problems is NP-hard (see *Algorithmic Discrete Mathematics*). However, we can use the dynamic program (DP) given below. Note that, by Observation 5.11, we may assume to $\mathbf{a}, \mathbf{c} \in \mathbb{N}^n$, $\beta \in \mathbb{N}$.

Algorithm: Dynamic program for the knapsack problem.

input: $\mathbf{a}, \mathbf{c} \in \mathbb{N}^n$; $\beta \in \mathbb{N}$, upper bound $\bar{c} \in \mathbb{N}$ on optimal value

output: optimum value of $\{\mathbf{c}^\top \mathbf{x} : \mathbf{a}^\top \mathbf{x} \leq \beta, \mathbf{x} \in \{0, 1\}^n\}$

$z(0, 0) \leftarrow 0, z(0, k) \leftarrow \infty \forall k \in \{1, \dots, \bar{c}\}$

for $j \leftarrow 1, \dots, n$:

for $k \leftarrow 0, \dots, c_j - 1$:

$z(j, k) \leftarrow z(j - 1, k)$

for $k \leftarrow c_j, \dots, \bar{c}$:

$z(j, k) \leftarrow \min\{z(j - 1, k - c_j) + a_j, z(j - 1, k)\}$

return $\max\{k \in \{0, \dots, \bar{c}\} : z(n, k) \leq \beta\}$

We prove correctness of the DP.

Theorem 5.16. The above DP solves the knapsack problem in time $\mathcal{O}(n \cdot \bar{c})$ when given an upper bound \bar{c} on optimum value, e.g., $\bar{c} = \mathbf{1}^\top \mathbf{c}$.

Proof. The algorithm runs in pseudo polynomial time $\mathcal{O}(n \cdot \bar{c})$. It uses a DP table where entry $z(j, k)$ is intended to hold the smallest weight of a partial solution only using the first j variables and reaching objective function value $k \in \{0, \dots, \bar{c}\}$, and ∞ if no such solution exists, i.e.,

$$z(j, k) = \min \{a(S) : a(S) \leq \beta, c(S) = k, S \subseteq \{1, \dots, j\}\}.$$

We have $z(j, k) \leq z(j-1, k)$, because it is always possible to discard item j . A solution of smaller weight is available when using the j -th item if and only if $k \geq c_j$, $z(j-1, k-c_j) + a_j \leq \beta$ and $z(j-1, k-c_j) + a_j < z(j-1, k)$. Hence, the DP computes the entries of the table correctly and therefore obtains the correct result overall.

An optimum solution can be found, as usual, via backtracking (see *Algorithmic Discrete Mathematics*). \square

The smallest minimal cover can be found by solving the knapsack problem

$$\begin{aligned} & \min \{ \mathbf{1}^\top \mathbf{x} : \mathbf{a}^\top \mathbf{x} \geq \beta + 1, \mathbf{x} \in \{0, 1\}^n \} \\ & = n - \max \{ \mathbf{1}^\top \mathbf{y} : \mathbf{a}^\top \mathbf{y} \leq \mathbf{1}^\top \mathbf{a} - \beta - 1, \mathbf{y} \in \{0, 1\}^n \}, \end{aligned}$$

where we substituted $\mathbf{y} := \mathbf{1} - \mathbf{x}$. Because in this case $\bar{c} = n$ is polynomial in the input size, the DP runs in polynomial time. Note that, we can also obtain minimal (but not always the smallest) cover by iteratively removing elements as long as possible.

The lifting problem (5.15) again involves a knapsack problem, which can again be solved by our dynamic program. The optimum objective function value is polynomially bounded by $\bar{c} := |C| - 1$. Overall, using Theorem 5.16 we can lift cover inequalities in polynomial time.

Corollary 5.17. The computation of a fully lifted cover inequality is possible in time $\mathcal{O}(n^2 |C|) \subseteq \mathcal{O}(n^3)$.

However, it is difficult to guarantee that the corresponding inequality truncates a given fractional point, in fact is NP-complete to decide whether a lifted cover inequality exists that excludes a given point (see [31]). Note that this does not follow from the equivalence of optimization and separation for the NP-complete knapsack problem (see *Introduction to Optimization*), since we restrict ourselves to inequalities of a specific type. The hardness of the specific separation problem means that, in practice, we have to resort to heuristics for lifting and have to terminate the lifting procedure after some number of lifting operations.

Remark 5.18. There are many other classes of inequalities for the knapsack polytope, e.g., the so-called $(1, k)$ -configuration or extended-weight inequalities (see [37, 42, 33]).

5.2.2 Inequalities for set-packing problems

Mixed-integer problems often contain constraints that only involve binary variables and coefficients. In particular, many binary programs contain logical inequalities of the form $x_i + x_j \leq 1$ or $x_i \leq x_j$. Furthermore, preprocessing routines (see Section 3.2.3) for integer problems can automatically extract implicit logical conditions. Such conditions occur, for example, when the binary variables x_i model whether a certain combination of resources is used. Conflicts between two combinations x_i and x_j are then expressed by $x_i + x_j \leq 1$.

This leads to binary programs with binary coefficient matrices. The study of such problems and, in particular, set-packing problems plays an important role in combinatorial optimization. A deep theory has been developed for these problems, which deals with notions such as perfect, ideal, or balanced matrices, perfect graphs, blocking and anti-blocking polyhedra, independence systems and semidefinite optimization.

The focus of this section is on the (partial) description of the associated polyhedra by means of inequalities. Since relaxations of many integer problems lead to set-packing problems, understanding these polyhedra can often lead to improved formulations of the initial problem.

Definition 5.19. The *set-packing problem* is given by

$$\max \{ \mathbf{c}^\top \mathbf{x} : \mathbf{A}\mathbf{x} \leq \mathbf{1}, \mathbf{x} \in \{0, 1\}^n \}$$

with $A \in \{0, 1\}^{m \times n}$ and $\mathbf{c} \in \mathbb{R}^n$, cf. Example 1.7.

Remark 5.20. Every column $A_{\cdot j}$ of A can be viewed as the incidence vector of a subset $F_j \subseteq \{1, \dots, m\}$, i.e., $F_j := \{i \in \{1, \dots, m\} : A_{ij} = 1\}$. With this interpretation, the set-packing problem consists in finding a maximum (wrt. \mathbf{c}) selection of sets from F_1, \dots, F_n that are pairwise disjoint.

Feasible solutions of the set-packing problem have a nice graph-theoretic interpretation. We introduce a vertex for every column of A and an edge $\{j, \ell\}$ between two vertices j and ℓ if their corresponding columns have a 1 in the same row.

Definition 5.21. The *(column) conflict graph* $G(A) = (V, E)$ of $A \in \{0, 1\}^{m \times n}$ is defined by $V := \{1, \dots, n\}$ and

$$E := \{\{j, \ell\} : \exists i \in \{1, \dots, m\} : A_{ij} = A_{i\ell} = 1\}.$$

Obviously, every feasible $\mathbf{x} \in \{0, 1\}^n$ that satisfies the inequality $\mathbf{A}\mathbf{x} \leq \mathbf{1}$ is the incidence vector of a stable set in the graph $G(A)$. Conversely, the incidence vector of a stable set in $G(A)$ is a feasible solution of the set-packing problem $\mathbf{A}\mathbf{x} \leq \mathbf{1}$. In other words, the study of stable sets in graphs is equivalent to the study of the set-packing problem.

Now consider $A \in \{0, 1\}^{m \times n}$ and denote the *set-packing polytope* by

$$\mathcal{P}(A) = \text{conv}(\{\mathbf{x} \in \{0, 1\}^n : \mathbf{A}\mathbf{x} \leq \mathbf{1}\}).$$

Let $G(A) = (V, E)$ be the conflict graph of A . From our previous considerations it follows that

$$\mathcal{P}(A) = \text{conv}(\{\mathbf{x} \in \{0, 1\}^n : x_j + x_\ell \leq 1 \ \forall \{j, \ell\} \in E\}) =: \mathcal{P}(G),$$

where the set in the convex hull consists of the stable sets in G , and $\mathcal{P}(G)$ is the *stable-set polytope*. Expressed differently, two matrices A and A' yield the same set-packing problem if and only if their corresponding conflict graphs coincide. We can therefore consider $\mathcal{P}(A)$ via the graph G and denote both the set-packing and the stable set polytope by $\mathcal{P}(G)$.

First, we make a few simple observations regarding $\mathcal{P}(G)$.

Proposition 5.22.

- (a) $\mathcal{P}(G)$ is full-dimensional, i.e., $\dim(\mathcal{P}(G)) = n$.
- (b) $\mathcal{P}(G)$ is downwards monotone, i.e., $\mathbf{x} \in \mathcal{P}(G)$ implies $\mathbf{y} \in \mathcal{P}(G)$ for all $\mathbf{0} \leq \mathbf{y} \leq \mathbf{x}$.
- (c) All facets of $\mathcal{P}(G)$ that are not facets of $[0, 1]^n$ have non-negative coefficients, i.e., if $\mathbf{a}^\top \mathbf{x} \leq \beta$ defines such a facet, then $\mathbf{a} \geq \mathbf{0}$ and $\beta \geq 0$.
- (d) The non-negativity conditions $x_j \geq 0$ induce facets of $\mathcal{P}(G)$.

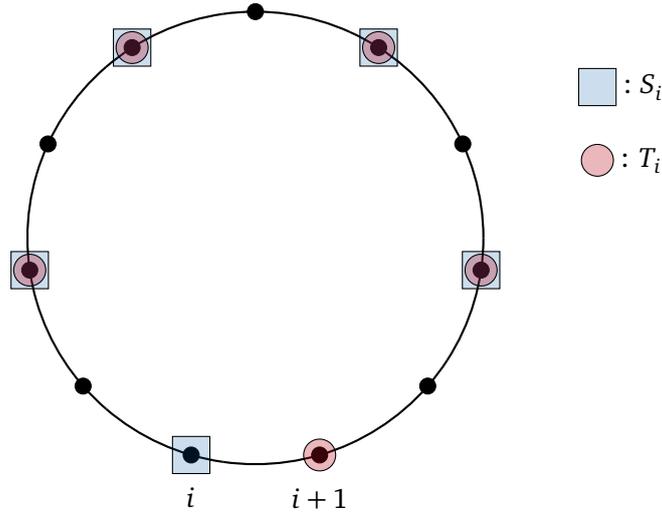


Figure 5.2: Illustration of the sets S_i and T_i in the proof of Theorem 5.23. The sets differ only by the elements i and $i + 1$ (modulo k).

Proof. exercise. □

From Theorems 4.15 and 4.8 we know that if G is bipartite, then the edge and the non-negativity conditions are sufficient to fully describe $\mathcal{P}(G)$. Conversely, non-bipartite graphs contain odd cycles generating new valid inequalities that are not induced by linear combinations of edge inequalities. By lifting, similarly to Theorem 5.15, we obtain facets of $\mathcal{P}(G)$.

Theorem 5.23 ([35]). Let $K = (V_K, E_K) \subseteq G$ be a cycle of odd length in an undirected graph G . Then, the *odd-cycle inequality*

$$\sum_{v \in V_K} x_v \leq \frac{|E_K| - 1}{2}, \quad (5.17)$$

is valid for $\mathcal{P}(G)$. It defines a facet of $\mathcal{P}(G[V_K])$ for the induced subgraph $G[V_K]$ if and only if K is a *hole*^a.

^aA *hole* is an induced subgraph that is a cycle.

Proof. The validity of the inequality is clear.

Consider an odd hole $K = (V_K, E_K)$ with, without loss of generality, $V_K = \{0, \dots, k-1\}$ and $E_K = \{\{0, 1\}, \{1, 2\}, \dots, \{k-2, k-1\}, \{k-1, 0\}\}$.

We first observe that

$$\mathcal{F} := \{\mathbf{y} \in \mathcal{P}(G[V_K]) : \mathbf{1}^\top \mathbf{y} = \frac{k-1}{2}\}$$

is a nontrivial face of $\mathcal{P}(G[V_K])$. If it is not a facet, \mathcal{F} is part of a facet $\mathcal{F}' := \{\mathbf{y} \in \mathcal{P}(G[V_K]) : \mathbf{b}^\top \mathbf{y} = \beta\} \supseteq \mathcal{F}$ of $\mathcal{P}(G[V_K])$. We consider the following sets for $i \in V_K$ (see Figure 5.2):

$$\begin{aligned} S_0 &= \{0, 3, 5, \dots, k-2\}, & S_i &= (S_0 + i) \pmod k, \\ T_0 &= \{1, 3, 5, \dots, k-2\}, & T_i &= (T_0 + i) \pmod k. \end{aligned}$$

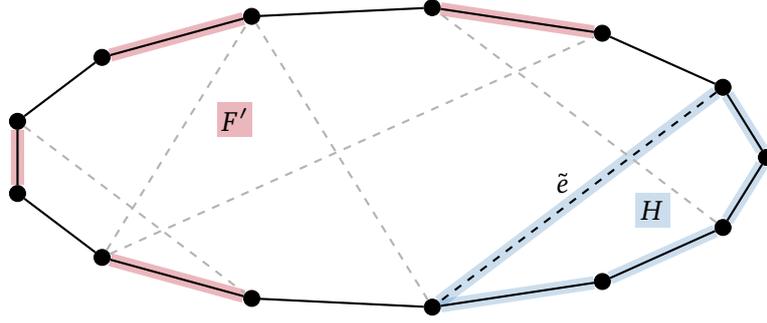


Figure 5.3: Example for the sets H and F' in the proof of Theorem 5.23. In the example, the cycle K has length 13.

(Example: For $k = 7$, $S_0 = \{0, 3, 5\}$, $S_1 = \{1, 4, 6\}$, $S_2 = \{0, 2, 5\}$, \dots , $S_6 = \{2, 4, 6\}$, $T_0 = \{1, 3, 5\}$, $T_1 = \{2, 4, 6\}$, $T_2 = \{0, 3, 5\}$, \dots , $T_6 = \{1, 2, 4\}$).

Since K has no chords, all sets are stable with respect to G and $\mathbf{1}^\top \chi^{S_i} = \mathbf{1}^\top \chi^{T_i} = \frac{k-1}{2}$, where $\chi^S := \sum_{i \in S} \mathbf{e}^i \in \{0, 1\}^{V_K}$ denotes the characteristic vector of the set S . From this it follows that $\chi^{S_i}, \chi^{T_i} \in \mathcal{F} \subseteq \mathcal{F}'$ and thus, for $i < k-1$,

$$0 = \mathbf{b}^\top \chi^{S_i} - \mathbf{b}^\top \chi^{T_i} = \mathbf{b}^\top \chi^{S_i \setminus T_i} - \mathbf{b}^\top \chi^{T_i \setminus S_i} = \mathbf{b}^\top \mathbf{e}^i - \mathbf{b}^\top \mathbf{e}^{i+1} = b_i - b_{i+1},$$

and $0 = b_i - b_0$ for $i = k-1$. Since i was chosen arbitrarily, it follows that $b_i = b_j$ for all $i, j \in V_K$. This means that $\mathbf{b}^\top \mathbf{x} \leq \beta$ is the odd-cycle inequality up to a positive factor, since both have the same left-hand side up to a positive (by Proposition 5.22 (c)) factor and $\mathcal{F} \subseteq \mathcal{F}'$. It follows that $\mathcal{F} = \mathcal{F}'$, hence \mathcal{F} is a facet of $\mathcal{P}(G')$.

Now assume that the cycle K contains a chord $\tilde{e} \in G[V_K] \setminus E_K$ in $G[V_K]$ (see Figure 5.3). Then, $E_K \cup \{\tilde{e}\}$ consists of two cycles that have \tilde{e} in common, where exactly one of the cycles has odd length. Let $H = (V_H, E_H)$ be this odd cycle and let $E_F = \{e \in K : e \cap V_H = \emptyset\}$. Then, $|E_F| = (|E_K| + 1) - |E_H| - 2 = |E_K| - |E_H| - 1$ is odd. Hence, there are pairwise disjoint edges, i.e., a matching, $M \subseteq E_F$ with $|M| = \frac{|E_F|+1}{2} = \frac{|E_K|-|E_H|}{2}$. Now, the odd-cycle inequality (5.17) is the sum of the valid inequalities $\sum_{v \in V(H)} x_v \leq \frac{|E_H|-1}{2}$ and $x_u + x_v \leq 1$ for all $\{u, v\} \in M$. The corresponding faces therefore each contain the face belonging to (5.17) (if (5.17) is satisfied with equality, each of the other inequalities must also be satisfied with equality). Thus, (5.17) cannot induce a facet of $\mathcal{P}(G[V_K])$. \square

Remark 5.24. Odd-cycle inequalities can be separated in polynomial time as follows (see [21, Lemma 9.1.11]). We create from the conflict graph $G(A) = (V, E)$ a bipartite auxiliary graph $G' = (V \cup V', E')$, where V' contains a vertex v' for every vertex $v \in V$ and $E' := \{\{u, v'\} : \{u, v\} \in E\}$, i.e, we replace every edge $\{u, v\} \in E$ by the two edges $\{u, v'\}$ and $\{v, u'\}$. Let $\mathbf{x}^* \in [0, 1]^V$ be an optimum solution of the LP relaxation of the set-packing problem. We define edge weights via $w_{\{u, v'\}} = \frac{1}{2}(1 - x_u^* - x_{v'}^*) \geq 0$ for all $\{u, v'\} \in E'$ with $u \in V$ and $v' \in V'$.

Odd cycles $K = (V_K, E_K)$ in $G(A)$ correspond to minimal paths K' from a vertex $s \in V$ to the corresponding vertex $s' \in V'$ in G' , in particular, K' thus contains at most one of the vertices v, v' for all $v \in V \setminus \{s\}$. The weight of K' is

$$w(K') = \sum_{e \in K'} w_e = \sum_{\{u, v'\} \in K'} \frac{1}{2}(1 - x_u^* - x_{v'}^*) = \frac{|E_K|}{2} - \sum_{v \in V_K} x_v^*,$$

because every x_v^* occurs twice in the sum. This means that

$$w(K') < \frac{1}{2} \iff \frac{|E_K| - 1}{2} < \sum_{v \in V_K} x_v^*.$$

We can therefore find a violated odd-cycle inequality by looking for a v - v' -path in G' of weight less than $\frac{1}{2}$, e.g., using Dijkstra's algorithm (see *Algorithmic Discrete Mathematics*).

Graphs $G = (V, E)$, for which $\mathcal{P}(G)$ is completely characterized by the non-negativity constraints, the edge inequalities $x_i + x_j \leq 1$ for $\{i, j\} \in E$ and the odd-cycle inequalities are called *t-perfect* (see [10]). The class of *t-perfect* graphs contains a number of graphs, such as, e.g., series-parallel and bipartite graphs. In practice, odd-cycle inequalities, however, are usually of little help in finding integral optima.

Another important class of valid inequalities for the stable set polytope are clique inequalities. Recall that a clique is a complete subgraph.

Theorem 5.25. Let C be a clique in an undirected graph G . Then, the *clique inequality*

$$\sum_{j \in C} x_j \leq 1$$

is valid for $\mathcal{P}(G)$ and it defines a facet of $\mathcal{P}(G)$ if and only if C is maximal with respect to vertex inclusion.

Proof. exercise. □

Graphs G , for which $\mathcal{P}(G)$ is completely described by the non-negativity conditions and the clique inequalities, are called *perfect* (see [5]).

Unfortunately, the separation problem for the class of clique inequalities is NP-hard (see [21]). Surprisingly, there is a larger class of inequalities, the *orthonormal-representation inequalities*, that generalize the clique inequalities and can be separated in polynomial time. In addition to the cycle, clique and orthonormal-representation inequalities, a number of other inequalities are known for the stable set polytope, e.g., the blossom, odd antihole, wheel, antiweb and web, wedge, and many more inequalities (see [6]).

6 Decomposition Methods

The idea of decomposition methods is to remove a part (constraints and/or variables) from the problem and consider them separately in a so-called *master problem*. The remaining subproblem can often be solved more efficiently. Decomposition methods now work alternatingly on the master problem and the subproblem and iteratively exchange information in order to solve the initial problem optimally. In this section we look at three well-known examples of this approach: the Lagrangian relaxation/decomposition, Dantzig-Wolfe decomposition and Benders' decomposition. In Lagrangian relaxation/decomposition, parts of the constraint matrix are omitted and instead treated via the objective function. Dantzig-Wolfe decomposition and Benders' decomposition also omit part of the constraint matrix, but instead of treating it via the objective function, it is reformulated and reinserted into the constraints.

6.1 Lagrangian relaxation

In the previous chapter, we solved the mixed-integer problem by relaxing the integrality constraints and trying to increase integrality of the solution by adding cutting planes. The method considered in this section employs a different relaxation that relies on the so-called *Lagrangian function*.

The system of inequalities $Ax \leq b$ is divided into two parts $A^{(1)}x \leq b^{(1)}$ and $A^{(2)}x \leq b^{(2)}$ such that

$$A = \begin{pmatrix} A^{(1)} \\ A^{(2)} \end{pmatrix} \quad \text{and} \quad b = \begin{pmatrix} b^{(1)} \\ b^{(2)} \end{pmatrix},$$

with $A \in \mathbb{Q}^{m \times n}$, $A^{(1)} \in \mathbb{Q}^{m_1 \times n}$, $A^{(2)} \in \mathbb{Q}^{m_2 \times n}$, $b^{(1)} \in \mathbb{Q}^{m_1}$, $b^{(2)} \in \mathbb{Q}^{m_2}$ and $m_1 + m_2 = m$. The general mixed-integer program (1.1) then becomes

$$\begin{aligned} \max \quad & c^\top x \\ \text{s.t.} \quad & A^{(1)}x \leq b^{(1)}, \\ & A^{(2)}x \leq b^{(2)}, \\ & x \in \mathbb{Z}^p \times \mathbb{R}^{n-p}. \end{aligned} \tag{6.1}$$

The subdivision is chosen so that the part $A^{(1)}x \leq b^{(1)}$ contains the inequalities that are difficult to treat and should therefore be relaxed. More specifically, the relaxed subsystem will be handled via a penalty term in the objective rather than as constraints. The penalty will be set such that it ensures that the relaxed constraints will be implicitly fulfilled by every optimum solution.

For a fixed $\lambda \in \mathbb{R}^{m_1}$, $\lambda \geq 0$, consider the *Lagrangian function*

$$L(\lambda) = \max \{c^\top x + \lambda^\top (b^{(1)} - A^{(1)}x) : x \in \mathcal{X}^{(2)}\}, \tag{6.2}$$

where $\mathcal{X}^{(2)} := \{x \in \mathbb{Z}^p \times \mathbb{R}^{n-p} : A^{(2)}x \leq b^{(2)}\}$. The value $L(\lambda)$ is an upper bound for (6.1), since every feasible solution \bar{x} of (6.1) satisfies

$$c^\top \bar{x} \leq c^\top \bar{x} + \lambda^\top (b^{(1)} - A^{(1)}\bar{x}) \leq \max_{x \in \mathcal{X}^{(2)}} c^\top x + \lambda^\top (b^{(1)} - A^{(1)}x) = L(\lambda),$$

where the first inequality holds because $A^{(1)}\bar{x} \leq \mathbf{b}^{(1)}$ and $\boldsymbol{\lambda} \geq \mathbf{0}$ imply $\boldsymbol{\lambda}^\top(\mathbf{b}^{(1)} - A^{(1)}\bar{x}) \geq 0$, i.e., a violation of $A^{(1)}\mathbf{x} \leq \mathbf{b}^{(1)}$ would decrease the objective function value and thus penalize it.

Since this applies for all $\boldsymbol{\lambda} \geq \mathbf{0}$, the *Lagrangian relaxation*

$$\min_{\boldsymbol{\lambda} \geq \mathbf{0}} L(\boldsymbol{\lambda}) \quad (6.3)$$

yields the smallest upper bound for (6.1) of this kind. Sometimes, (6.3) is referred to as the *Lagrangian dual* (see Observation 6.4 for a motivation of the term).

But how can we compute an optimum solution $\boldsymbol{\lambda}^*$ of (6.3)? And how close is $L(\boldsymbol{\lambda}^*)$ to the optimum value of (6.1)? We first answer the latter question. In the following, observe that $\text{conv}(\mathcal{X}^{(2)}) = \mathcal{P}(A^{(2)}, \mathbf{b}^{(2)})_{I,p}$ is the integer hull of $\mathcal{P}(A^{(2)}, \mathbf{b}^{(2)})$.

Theorem 6.1. It holds that

$$\min_{\boldsymbol{\lambda} \geq \mathbf{0}} L(\boldsymbol{\lambda}) = \max \{ \mathbf{c}^\top \mathbf{x} : A^{(1)}\mathbf{x} \leq \mathbf{b}^{(1)}, \mathbf{x} \in \text{conv}(\mathcal{X}^{(2)}) \}.$$

Proof. We obtain

$$L(\boldsymbol{\lambda}) = \max \{ \mathbf{c}^\top \mathbf{x} + \boldsymbol{\lambda}^\top(\mathbf{b}^{(1)} - A^{(1)}\mathbf{x}) : \mathbf{x} \in \mathcal{X}^{(2)} \} = \max \{ \mathbf{c}^\top \mathbf{x} + \boldsymbol{\lambda}^\top(\mathbf{b}^{(1)} - A^{(1)}\mathbf{x}) : \mathbf{x} \in \text{conv}(\mathcal{X}^{(2)}) \},$$

because the objective function is linear in \mathbf{x} and $\text{conv}(\mathcal{X}^{(2)}) = \mathcal{P}(A^{(2)}, \mathbf{b}^{(2)})_{I,p}$ is a polyhedron (Theorem 2.32), i.e., the second maximum is attained at an extreme point, which cannot be a convex combination of other points in $\mathcal{X}^{(2)}$ (see *Introduction to Optimization*). It follows that

$$\begin{aligned} \min_{\boldsymbol{\lambda} \geq \mathbf{0}} L(\boldsymbol{\lambda}) &= \min_{\boldsymbol{\lambda} \geq \mathbf{0}} \max \{ \mathbf{c}^\top \mathbf{x} + \boldsymbol{\lambda}^\top(\mathbf{b}^{(1)} - A^{(1)}\mathbf{x}) : \mathbf{x} \in \text{conv}(\mathcal{X}^{(2)}) \} \\ &= \min_{\boldsymbol{\lambda} \geq \mathbf{0}} \{ \boldsymbol{\lambda}^\top \mathbf{b}^{(1)} + \max \{ (\mathbf{c}^\top - \boldsymbol{\lambda}^\top A^{(1)})\mathbf{x} : \mathbf{x} \in \text{conv}(\mathcal{X}^{(2)}) \} \}. \end{aligned}$$

If $\mathcal{X}^{(2)} = \emptyset$, the inner maximum is defined as $-\infty$ for all $\boldsymbol{\lambda}$ and thus $\min_{\boldsymbol{\lambda} \geq \mathbf{0}} L(\boldsymbol{\lambda}) = -\infty$. Otherwise, by Theorems 2.32 and 2.19, vectors $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)} \in \mathbb{R}^n$ and $\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(\ell)} \in \mathbb{R}^n$ exist such that $\text{conv}(\mathcal{X}^{(2)}) = \text{conv}(\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}\}) + \text{cone}(\{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(\ell)}\})$. This means that

$$\max \{ (\mathbf{c}^\top - \boldsymbol{\lambda}^\top A^{(1)})\mathbf{x} : \mathbf{x} \in \text{conv}(\mathcal{X}^{(2)}) \} = \begin{cases} +\infty & \text{if } \exists i \in \{1, \dots, \ell\} : (\mathbf{c}^\top - \boldsymbol{\lambda}^\top A^{(1)})\mathbf{r}^{(i)} > 0, \\ \max_{j \in \{1, \dots, k\}} (\mathbf{c}^\top - \boldsymbol{\lambda}^\top A^{(1)})\mathbf{v}^{(j)} & \text{otherwise.} \end{cases} \quad (6.4)$$

It follows that

$$\min_{\boldsymbol{\lambda} \geq \mathbf{0}} L(\boldsymbol{\lambda}) = \min \{ \boldsymbol{\lambda}^\top \mathbf{b}^{(1)} + \max_{j \in \{1, \dots, k\}} (\mathbf{c}^\top - \boldsymbol{\lambda}^\top A^{(1)})\mathbf{v}^{(j)} : \boldsymbol{\lambda} \geq \mathbf{0}, (\mathbf{c}^\top - \boldsymbol{\lambda}^\top A^{(1)})\mathbf{r}^{(i)} \leq 0 \forall i \in \{1, \dots, \ell\} \}. \quad (6.5)$$

Note that if the problem is infeasible, i.e. no $\boldsymbol{\lambda} \geq \mathbf{0}$ with $(\mathbf{c}^\top - \boldsymbol{\lambda}^\top A^{(1)})\mathbf{r}^{(i)} \leq 0$ for all $i \in \{1, \dots, \ell\}$ exists, then the value of the minimum is ∞ by convention, in accordance with (6.4). Otherwise, the value is finite and attained in a vertex.

Using an additional variable $\eta \in \mathbb{R}$ with $\eta \geq \mathbf{c}^\top \mathbf{v}^{(j)} + \boldsymbol{\lambda}^\top(\mathbf{b}^{(1)} - A^{(1)}\mathbf{v}^{(j)})$ for all $j \in \{1, \dots, k\}$, we can reformulate (6.5) as

$$\begin{aligned} \min_{\boldsymbol{\lambda} \geq \mathbf{0}} L(\boldsymbol{\lambda}) &= \min_{\boldsymbol{\lambda} \geq \mathbf{0}, \eta \in \mathbb{R}} \eta & (6.6) \\ \text{s.t. } \eta - \boldsymbol{\lambda}^\top(\mathbf{b}^{(1)} - A^{(1)}\mathbf{v}^{(j)}) &\geq \mathbf{c}^\top \mathbf{v}^{(j)} & \forall j \in \{1, \dots, k\}, \\ \boldsymbol{\lambda}^\top A^{(1)}\mathbf{r}^{(i)} &\geq \mathbf{c}^\top \mathbf{r}^{(i)} & \forall i \in \{1, \dots, \ell\}. \end{aligned}$$

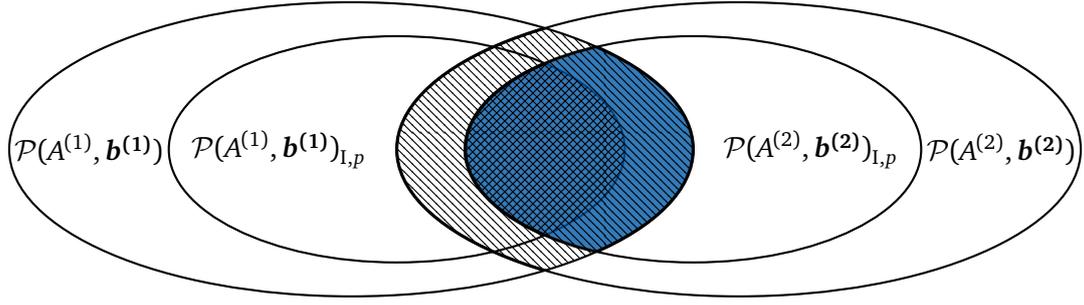


Figure 6.1: Illustration of the feasible region of the Lagrangian relaxation (blue). It contains the intersection of the integer hulls of $\mathcal{P}(A^{(1)}, \mathbf{b}^{(1)})$ and $\mathcal{P}(A^{(2)}, \mathbf{b}^{(2)})$ (crosshatch pattern) which contains the feasible region of the MIP. Conversely, the feasible region of the Lagrangian relaxation is contained in the feasible region of the LP relaxation (line pattern).

By strong duality, we finally obtain

$$\begin{aligned} \min_{\lambda \geq 0} L(\lambda) &= \max \quad \mathbf{c}^\top \left(\sum_{j=1}^k \mathbf{v}^{(j)} \alpha_j + \sum_{i=1}^{\ell} \mathbf{r}^{(i)} \beta_i \right) \\ \text{s.t.} \quad &\sum_{j=1}^k \alpha_j = 1, \\ &A^{(1)} \left(\sum_{j=1}^k \mathbf{v}^{(j)} \alpha_j + \sum_{i=1}^{\ell} \mathbf{r}^{(i)} \beta_i \right) \leq \mathbf{b}^{(1)} \sum_{j=1}^k \alpha_j, \\ &\alpha_j \geq 0 \quad \forall j \in \{1, \dots, k\}, \\ &\beta_i \geq 0 \quad \forall i \in \{1, \dots, \ell\}. \end{aligned}$$

Since $\sum_{j=1}^k \alpha_j = 1$, this is equivalent to

$$\begin{aligned} \min_{\lambda \geq 0} L(\lambda) &= \max \{ \mathbf{c}^\top \mathbf{x} : A^{(1)} \mathbf{x} \leq \mathbf{b}^{(1)}, \mathbf{x} \in \text{conv}(\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}\}) + \text{cone}(\{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(\ell)}\}) \} \\ &= \max \{ \mathbf{c}^\top \mathbf{x} : A^{(1)} \mathbf{x} \leq \mathbf{b}^{(1)}, \mathbf{x} \in \text{conv}(\mathcal{X}^{(2)}) \}. \end{aligned} \quad \square$$

Because of

$$\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}\} \supseteq \{\mathbf{x} \in \mathbb{R}^n : A^{(1)}\mathbf{x} \leq \mathbf{b}^{(1)}, \mathbf{x} \in \text{conv}(\mathcal{X}^{(2)})\} \supseteq \text{conv}(\{\mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p} : A\mathbf{x} \leq \mathbf{b}\}),$$

Theorem 6.1 implies the following, see Figure 6.1.

Corollary 6.2. It holds that $z_{\text{MIP}} \leq \min_{\lambda \geq 0} L(\lambda) \leq z_{\text{LP}}$, where z_{MIP} and z_{LP} denote the optimum solution values of the MIP (6.1) and its LP relaxation, respectively.

In particular, when $p = 0$, we have $z_{\text{MIP}} = z_{\text{LP}}$ and the Lagrangian relaxation coincides with the LP-relaxation. The latter is also true if the non-relaxed part of the problem induces an integral feasible region.

Observation 6.3. If $\mathcal{P}(A^{(2)}, \mathbf{b}^{(2)})$ is integral, i.e., $\text{conv}(\mathcal{X}^{(2)}) = \mathcal{P}(A^{(2)}, \mathbf{b}^{(2)})$, Theorem 6.1 yields that the Lagrangian relaxation coincides with the LP relaxation, i.e.,

$$\min_{\lambda \geq 0} L(\lambda) = \max \{ \mathbf{c}^\top \mathbf{x} : \mathbf{x} \in \mathcal{P}(A^{(1)}, \mathbf{b}^{(1)}) \cap \mathcal{P}(A^{(2)}, \mathbf{b}^{(2)}) \} = \max \{ \mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b} \}.$$

We justify why $\min_{\lambda \geq 0} L(\lambda)$ is called the Lagrangian *dual*.

Observation 6.4. If $p = 0$ and all inequalities are relaxed, i.e., $\mathcal{P}(A^{(1)}, \mathbf{b}^{(1)}) = \mathcal{P}(A, \mathbf{b})$, then

$$L(\lambda) = \max \{ \mathbf{c}^\top \mathbf{x} + \lambda^\top (\mathbf{b} - A\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n \} = \max \{ (\mathbf{c}^\top - \lambda^\top A)\mathbf{x} + \lambda^\top \mathbf{b} : \mathbf{x} \in \mathbb{R}^n \}.$$

This problem is only bounded if $\mathbf{c} = A^\top \lambda$, since \mathbf{x} is unconstrained. The Lagrangian relaxation thus can be rewritten as

$$\min_{\lambda \geq 0} L(\lambda) = \min \{ \lambda^\top \mathbf{b} : A^\top \lambda = \mathbf{c}, \lambda \geq \mathbf{0} \},$$

which corresponds exactly to the LP dual of $\max \{ \mathbf{c}^\top \mathbf{x} : A\mathbf{x} \leq \mathbf{b} \}$.

Remark 6.5. Of course, the utility of the Lagrangian relaxation very much depends on the set of constraints being relaxed. On the one hand, we need to be able to compute $L(\lambda)$ reasonably quickly, which makes it desirable to relax as many difficult constraints as possible. On the other hand, the more constraints are relaxed, the worse the bound $L(\lambda^*)$ becomes.

We still need to address the computation of $\min L(\lambda)$. From a theoretical point of view, the equivalence of separation and optimization can be used to compute $\min L(\lambda)$ in polynomial time, provided that $\min \{ \tilde{\mathbf{c}}^\top \mathbf{x} : \mathbf{x} \in \text{conv}(\mathcal{X}^{(2)}) \}$ can be computed in polynomial time for every $\tilde{\mathbf{c}}$ (see [39]). In practice, the *subgradient method* is often employed to compute $\min L(\lambda)$.

6.1.1 Subgradient method

The *subgradient method* solves a general minimization problem $\min \{ f(\lambda) : \lambda \in \mathcal{S} \}$ for a convex function $f : \mathcal{S} \rightarrow \mathbb{R}$ and a convex and closed set $\mathcal{S} \subseteq \mathbb{R}^r$. It is particularly suited for the case when f is not differentiable. The subgradient method is based on the following generalization of a gradient for convex functions.

Definition 6.6. A vector $\mathbf{h} \in \mathbb{R}^r$ is a *subgradient* of a convex function $f : \mathbb{R}^r \rightarrow \mathbb{R}$ at point $\hat{\lambda}$ if, for all $\lambda \in \mathbb{R}^r$,

$$f(\lambda) \geq f(\hat{\lambda}) + \mathbf{h}^\top (\lambda - \hat{\lambda}).$$

In other words, subgradients are linear underestimators. Importantly, subgradients still allow to characterize global minima of convex functions.

Lemma 6.7. The point $\lambda^* \in \mathbb{R}^r$ minimizes a convex function $f : \mathbb{R}^r \rightarrow \mathbb{R}$ if and only if $\mathbf{0}$ is a subgradient of f at λ^* .

Proof. We have that $f(\lambda^*)$ is a (global) minimum if and only if, for all $\lambda \in \mathbb{R}^r$, it holds that $f(\lambda) - f(\lambda^*) \geq 0 = \mathbf{0}^\top (\lambda - \lambda^*)$, i.e., if and only if $\mathbf{0}$ is a subgradient of f at λ^* . \square

Much like simple gradient descent, the subgradient method relies on the fact that if $\lambda^* \in \mathcal{S}$ is a global minimum of f , then, for every subgradient $\mathbf{h} \in \mathbb{R}^r$ at every point $\lambda \in \mathcal{S}$,

$$0 \geq f(\lambda^*) - f(\lambda) \geq \mathbf{h}^\top (\lambda^* - \lambda),$$

i.e., \mathbf{h} points in opposing direction of the vector $(\lambda^* - \lambda)$, which means that $-\mathbf{h}$ points from λ (vaguely) towards λ^* .

With this in mind, the subgradient method starts with a point $\lambda^{(0)} \in \mathbb{R}^r$ and then iteratively computes $\lambda^{(k+1)}$ for $k \leftarrow 1, 2, \dots$ via $\lambda^{(k+1)} \leftarrow \Pi_{\mathcal{S}}(\lambda^{(k)} - \mu_k \mathbf{h}^{(k)} / \|\mathbf{h}^{(k)}\|)$. Here, $\mathbf{h}^{(k)}$ is a subgradient of f at point $\lambda^{(k)}$, the step size $\mu_k > 0$ is chosen such that it limits the overshoot in each step and thus improves convergence of the method, and $\Pi_{\mathcal{S}}$ is a projection onto the closed set \mathcal{S} defined by

$$\Pi_{\mathcal{S}}(\hat{\lambda}) := \arg \min_{\lambda \in \mathcal{S}} \|\lambda - \hat{\lambda}\|.$$

Overall, we obtain the following general method.

Algorithm: subgradient method

input: convex function $f : \mathcal{S} \rightarrow \mathbb{R}$ with $\mathcal{S} \subseteq \mathbb{R}^r$ convex and closed

output: $\arg \min \{f(\lambda) : \lambda \in \mathcal{S}\}$

choose $\lambda^{(0)} \in \mathcal{S}$; set $k \leftarrow 0$

compute subgradient $\mathbf{h}^{(0)}$ of f at $\lambda^{(0)}$

while $\mathbf{h}^{(k)} \neq \mathbf{0}$:

choose step size μ_k
 $\lambda^{(k+1)} \leftarrow \Pi_{\mathcal{S}}(\lambda^{(k)} - \mu_k \frac{\mathbf{h}^{(k)}}{\|\mathbf{h}^{(k)}\|})$

$k \leftarrow k + 1$

compute subgradient $\mathbf{h}^{(k)}$ of f at $\lambda^{(k)}$

return $\lambda^{(k)}$

We have already seen (*Introduction to Optimization*) that convex functions are continuous in the interior of their domains. We need the following strengthening (see lecture *Nonsmooth Optimization*).¹

Lemma 6.8. Every convex function $f : \mathcal{U} \rightarrow \mathbb{R}$ on a convex and open set $\mathcal{U} \subseteq \mathbb{R}^r$ is locally Lipschitz-continuous.

Note that it is crucial to require an open domain \mathcal{U} , since a convex functions may be discontinuous at the boundary. With this, we obtain the following convergence theorem.

Theorem 6.9. Let f be convex on an open set $\mathcal{U} \supseteq \mathcal{S}$, let $\sum_{k=0}^{\infty} \mu_k = \infty$ and let $\sum_{k=0}^{\infty} \mu_k^2 < \infty$. Then, the subgradient method either returns a minimum point of f or converges to one, provided any exists.

Proof. The statement for the case that the subgradient method terminates follows from Lemma 6.7.

We use (without proof) that projections onto convex sets are non-expansive, i.e., $\|\Pi_{\mathcal{S}}(\mathbf{x}) - \Pi_{\mathcal{S}}(\mathbf{y})\| \leq \|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^r$. For every minimum point $\lambda^* \in \mathcal{S}$ of f , we have $\Pi_{\mathcal{S}}(\lambda^*) = \lambda^*$ and thus

$$\begin{aligned} \|\lambda^{(k+1)} - \lambda^*\|^2 &= \|\Pi_{\mathcal{S}}(\lambda^{(k)} - \mu_k \frac{\mathbf{h}^{(k)}}{\|\mathbf{h}^{(k)}\|}) - \Pi_{\mathcal{S}}(\lambda^*)\|^2 \\ &\leq \|\lambda^{(k)} - \mu_k \frac{\mathbf{h}^{(k)}}{\|\mathbf{h}^{(k)}\|} - \lambda^*\|^2 \\ &= \|\lambda^{(k)} - \lambda^*\|^2 - 2 \frac{\mu_k}{\|\mathbf{h}^{(k)}\|} (\lambda^{(k)} - \lambda^*)^\top \mathbf{h}^{(k)} + \mu_k^2 \\ &\leq \|\lambda^{(k)} - \lambda^*\|^2 - 2 \frac{\mu_k}{\|\mathbf{h}^{(k)}\|} (f(\lambda^{(k)}) - f(\lambda^*)) + \mu_k^2, \end{aligned} \tag{6.7}$$

¹Recall that a function is locally Lipschitz-continuous if it is Lipschitz-continuous on every compact subset of its domain.

where the last inequality follows by definition of $\mathbf{h}^{(k)}$ being a subgradient of f at $\boldsymbol{\lambda}^{(k)}$. By applying this inequality iteratively, and because $f(\boldsymbol{\lambda}^{(i)}) \geq f(\boldsymbol{\lambda}^*)$ and $\mu_i > 0$, we obtain, for all $k \in \mathbb{N}$,

$$\|\boldsymbol{\lambda}^{(k+1)} - \boldsymbol{\lambda}^*\|^2 \leq \|\boldsymbol{\lambda}^{(0)} - \boldsymbol{\lambda}^*\|^2 - 2 \sum_{i=0}^k \frac{\mu_i}{\|\mathbf{h}^{(i)}\|} (f(\boldsymbol{\lambda}^{(i)}) - f(\boldsymbol{\lambda}^*)) + \sum_{i=0}^k \mu_i^2 \leq \|\boldsymbol{\lambda}^{(0)} - \boldsymbol{\lambda}^*\|^2 + \sum_{i=0}^k \mu_i^2 < \infty. \quad (6.8)$$

Thus, the sequence $(\boldsymbol{\lambda}^{(k)})_{k \in \mathbb{N}}$ is bounded, and we can find a compact set $\mathcal{L} \subseteq \mathcal{U}$ with $\boldsymbol{\lambda}^{(k)} \in \mathcal{L}$ in its interior for all $k \in \mathbb{N}$. Lemma 6.8 now implies Lipschitz-continuity on \mathcal{L} with some Lipschitz-constant $H \geq 0$. This means that the subgradients are bounded via $\varepsilon \|\mathbf{h}^{(k)}\|^2 = ((\boldsymbol{\lambda}^{(k)} + \varepsilon \mathbf{h}^{(k)}) - \boldsymbol{\lambda}^{(k)})^\top \mathbf{h}^{(k)} \leq f(\boldsymbol{\lambda}^{(k)} + \varepsilon \mathbf{h}^{(k)}) - f(\boldsymbol{\lambda}^{(k)}) \leq H \|\varepsilon \mathbf{h}^{(k)}\|$, i.e., $\|\mathbf{h}^{(k)}\| \leq H$, for $\varepsilon > 0$ sufficiently small to stay within \mathcal{U} .

Now, let $f^* := f(\boldsymbol{\lambda}^*)$ and $\bar{f}_k := \min \{f(\boldsymbol{\lambda}^{(0)}), \dots, f(\boldsymbol{\lambda}^{(k)})\}$. Because of $\|\boldsymbol{\lambda}^{(k+1)} - \boldsymbol{\lambda}^*\|^2 \geq 0$, rearranging (6.8) yields

$$2 \left(\sum_{i=0}^k \frac{\mu_i}{\|\mathbf{h}^{(i)}\|} \right) (\bar{f}_k - f^*) \leq 2 \sum_{i=0}^k \frac{\mu_i}{\|\mathbf{h}^{(i)}\|} (f(\boldsymbol{\lambda}^{(i)}) - f(\boldsymbol{\lambda}^*)) \leq \|\boldsymbol{\lambda}^{(0)} - \boldsymbol{\lambda}^*\|^2 + \sum_{i=0}^k \mu_i^2,$$

which implies that, for every minimum point $\boldsymbol{\lambda}^* \in \mathcal{S}$ of f ,

$$\bar{f}_k - f^* \leq \frac{\|\boldsymbol{\lambda}^{(0)} - \boldsymbol{\lambda}^*\|^2 + \sum_{i=0}^k \mu_i^2}{2 \left(\sum_{i=0}^k \frac{\mu_i}{\|\mathbf{h}^{(i)}\|} \right)} \leq \frac{H}{2} \cdot \frac{\|\boldsymbol{\lambda}^{(0)} - \boldsymbol{\lambda}^*\|^2 + \sum_{i=0}^k \mu_i^2}{\sum_{i=0}^k \mu_i}.$$

Because of $\sum_{i=0}^{\infty} \mu_i^2 < \infty$ and $\sum_{i=0}^{\infty} \mu_i = \infty$ it follows that $\bar{f}_k \rightarrow f^*$ for $k \rightarrow \infty$.

It remains to show that $(\boldsymbol{\lambda}^{(k)})_{k \in \mathbb{N}}$ also converges. Since the subgradient method does not terminate and since $(\bar{f}_k)_{k \in \mathbb{N}}$ converges to f^* , there is an infinite subset $K \subseteq \mathbb{N}$ such that $(f(\boldsymbol{\lambda}^{(k)}))_{k \in K}$ converges to f^* . Because $(\boldsymbol{\lambda}^{(k)})_{k \in \mathbb{N}}$ and thus $(\boldsymbol{\lambda}^{(k)})_{k \in K}$ is bounded, there is a convergent subsequence $(\boldsymbol{\lambda}^{(k)})_{k \in K'}$. Since $(f(\boldsymbol{\lambda}^{(k)}))_{k \in K'}$ converges to f^* and f is continuous on $\mathcal{U} \supseteq \mathcal{S}$ (Lemma 6.8), $(\boldsymbol{\lambda}^{(k)})_{k \in K'}$ converges to a minimum point $\boldsymbol{\lambda}^*$. For every $\varepsilon > 0$ there exists $k_0 \in K'$ such that $\|\boldsymbol{\lambda}^{(k_0)} - \boldsymbol{\lambda}^*\| < \varepsilon/2$ and $\sum_{k=k_0}^{\infty} \mu_k^2 < \varepsilon/2$. We obtain that, for all $k \geq k_0$,

$$\|\boldsymbol{\lambda}^{(k+1)} - \boldsymbol{\lambda}^*\|^2 \stackrel{(6.7)}{\leq} \|\boldsymbol{\lambda}^{(k_0)} - \boldsymbol{\lambda}^*\|^2 + \sum_{i=k_0}^k \mu_i^2 < \varepsilon.$$

Thus, $\lim_{k \rightarrow \infty} \boldsymbol{\lambda}^{(k)} = \boldsymbol{\lambda}^*$. □

Remark 6.10. We give some additional information pertaining to the subgradient method.

- (a) The function values $f(\boldsymbol{\lambda}^{(k)})$ are typically not monotone.
- (b) A possible choice for μ_k is $\frac{1}{k}$.
- (c) Requiring $\sum_{k=0}^{\infty} \mu_k = \infty$ and $\mu_k \rightarrow 0$ is, in general, insufficient for $(\boldsymbol{\lambda}^{(k)})_{k \in \mathbb{N}}$ or $(f(\boldsymbol{\lambda}^{(k)}))_{k \in \mathbb{N}}$ to converge.
- (d) The method is often aborted early when no more significant progress is made.

The subgradient method is easy to implement and well suited when subgradients can be efficiently calculated. However, the step size must be chosen carefully. A disadvantage of the method is that it does not provide an indication as to the progress of convergence towards the optimum value. More advanced methods are discussed in the lecture *Nonsmooth Optimization*.

6.1.2 Solving the Lagrangian relaxation

We now describe how the subgradient method can be used to compute the Lagrangian relaxation $\min_{\lambda \geq 0} L(\lambda)$. To apply the method, we need to establish that L is a convex function that can be extended convexly to an open set.

Lemma 6.11. The function $L(\lambda)$ is piecewise (affine) linear and convex on the domain where it is bounded.

Proof. First, we recall that $L(\hat{\lambda}) = \max\{c^\top x + \hat{\lambda}^\top (b^{(1)} - A^{(1)}x) : x \in \mathcal{X}^{(2)}\}$ and again consider (6.4) with $\text{conv}(\mathcal{X}^{(2)}) = \text{conv}(\{v^{(1)}, \dots, v^{(k)}\}) + \text{cone}(\{r^{(1)}, \dots, r^{(\ell)}\})$. We can see that $L(\hat{\lambda})$ is finite if and only if $\hat{\lambda}$ falls in the polyhedron $\{\lambda \in \mathbb{R}_{\geq 0}^{m_1} : \lambda^\top A^{(1)}r^{(i)} \geq c^\top r^{(i)} \forall i \in \{1, \dots, \ell\}\}$, and that, in this polyhedron, we have $L(\lambda) = \max_{j \in \{1, \dots, k\}} \{c^\top v^{(j)} + \lambda^\top (b^{(1)} - A^{(1)}v^{(j)})\}$, i.e., $L(\lambda)$ is the maximum over a finite number of affine linear functions. In particular, $L(\lambda)$ is convex. \square

The key ingredient of the subgradient method is the computation of subgradients. For the Lagrangian relaxation, we can obtain subgradients by finding optimum solutions x^* of $L(\lambda)$ for fixed λ .

Lemma 6.12. Consider $\hat{\lambda} \in \mathbb{R}_{\geq 0}^{m_1}$ and an optimum solution \hat{x} of $L(\hat{\lambda})$. Then, $h = b^{(1)} - A^{(1)}\hat{x}$ is a subgradient of L at $\hat{\lambda}$.

Proof. Let $\lambda \in \mathbb{R}_{\geq 0}^{m_1}$ and $x^{(\lambda)}$ be an optimum solution of $L(\lambda)$. Then,

$$\begin{aligned} L(\lambda) - L(\hat{\lambda}) &= c^\top x^{(\lambda)} + \lambda(b^{(1)} - A^{(1)}x^{(\lambda)}) - (c^\top \hat{x} + \hat{\lambda}^\top (b^{(1)} - A^{(1)}\hat{x})) \\ &\geq c^\top \hat{x} + \lambda(b^{(1)} - A^{(1)}\hat{x}) - (c^\top \hat{x} + \hat{\lambda}^\top (b^{(1)} - A^{(1)}\hat{x})) = h^\top (\lambda - \hat{\lambda}), \end{aligned}$$

where the inequality follows because $x^{(\lambda)}$ is an optimum solution with respect to λ . \square

Finally, we have $S = \mathbb{R}_{\geq 0}^{m_1}$ in the subgradient method and the projection Π_S is given by $\Pi(\lambda)_i = \max\{\lambda_i, 0\}$ for all $i \in \{1, \dots, m_1\}$.

In the following, we describe applications where the Lagrangian relaxation has successfully been used in combination with the subgradient method. In these applications a good balance between efficiency of the subgradient method and strength of the relaxation can be achieved.

Example 6.13. If the constraint matrix has (almost) block diagonal form, see Figure 6.2, we can choose $A^{(1)}x \leq b^{(1)}$ as the coupling conditions. Relaxing them in a Lagrangian relaxation decomposes the remaining matrix into k independent blocks. In this case, $L(\lambda)$ is the sum of k independent terms that can separately be determined.

Often, each individual block $M^{(i)}$ represents a network flow problem, a knapsack problem, or the like, and can therefore be solved using specialized combinatorial algorithms. Other examples are multicommodity flow problems, which arise, for example, in vehicle routing or in decompositions of stochastic mixed-integer problems. \triangle

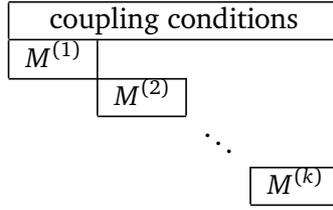


Figure 6.2: Matrix mostly in block diagonal form, except for a coupling part.

6.1.3 Application in stochastic optimization*

Two-stage stochastic optimization problems give rise to a constraint matrix with a structure as in Figure 6.2. In these problems, we consider a finite set Ω of mutually disjoint random events called *scenarios*. Each scenario $\omega \in \Omega$ occurs with some probability p_ω . Furthermore, mixed-integer first-stage decisions $\mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}$ are to be made, that are unaffected by stochasticity. Once a certain scenario ω occurs, then second-stage decisions $\mathbf{y} = \mathbf{y}(\mathbf{x}, \omega)$ have to be made based on the predetermined first-stage decisions and the event at hand. The goal is to minimize the total cost, consisting of the costs for \mathbf{x} and the *expected* cost for \mathbf{y} . If the constraints on the second-stage decisions are linear, the problem can be formulated as its so-called *deterministic equivalent*

$$\begin{aligned}
 \min \quad & \mathbf{c}^\top \mathbf{x} + \sum_{\omega \in \Omega} p_\omega (\mathbf{d}^{(\omega)})^\top \mathbf{y}^{(\omega)} \\
 \text{s.t.} \quad & W^{(\omega)} \mathbf{y}^{(\omega)} \leq \mathbf{h}^{(\omega)} - T^{(\omega)} \mathbf{x} && \forall \omega \in \Omega, \\
 & A\mathbf{x} \leq \mathbf{b}, \\
 & \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}.
 \end{aligned}$$

Here, the second term in the objective function specifies the expected costs for the second-stage decisions, where $\mathbf{d}^{(\omega)}$ denotes the costs for scenario ω . The coupling between the first and second stage for a $\omega \in \Omega$ is done by means of the *recourse matrix* $W^{(\omega)}$ and *technology matrix* $T^{(\omega)}$.

For fixed first-stage decision \mathbf{x} , the problem decomposes into $|\Omega|$ independent subproblems. In order to use a Lagrangian relaxation we introduce a copy $\mathbf{x}^{(\omega)}$ of \mathbf{x} for every scenario ω and constrain them to be all of equal value through coupling conditions $\mathbf{x}^{(\omega)} = \mathbf{x}^{(\omega')}$. We formulate these coupling conditions as $\sum_{\omega \in \Omega} H^{(\omega)} \mathbf{x}^{(\omega)} = \mathbf{0}$ and obtain the equivalent problem

$$\begin{aligned}
 \min \quad & \sum_{\omega \in \Omega} (\mathbf{c}^\top \mathbf{x} + p_\omega (\mathbf{d}^{(\omega)})^\top \mathbf{y}^{(\omega)}) \\
 \text{s.t.} \quad & \sum_{\omega \in \Omega} H^{(\omega)} \mathbf{x}^{(\omega)} = \mathbf{0}, \\
 & A\mathbf{x}^{(\omega)} \leq \mathbf{b}, && \forall \omega \in \Omega, \\
 & \mathbf{x}^{(\omega)} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}, && \forall \omega \in \Omega, \\
 & W^{(\omega)} \mathbf{y}^{(\omega)} \leq \mathbf{h}^{(\omega)} - T^{(\omega)} \mathbf{x}^{(\omega)}, && \forall \omega \in \Omega.
 \end{aligned}$$

The constraint matrix of this problem has the structure of Figure 6.2. A Lagrangian relaxation of the coupling conditions then leads to a decomposition. The equality of the different \mathbf{x}^ω is established step by step (see [7]).

6.1.4 Application in combinatorial optimization

Another application of Lagrangian relaxation is the famous traveling salesperson problem (TSP) (see *Algorithmic Discrete Mathematics*). Recall that an instance of TSP is given by an undirected complete graph $G = (V, E)$

with edge weights $c \in \mathbb{Q}_{\geq 0}^E$, and the problem consists in finding a cycle (called a *tour*) R that visits all vertices (exactly once) and minimizes $c(R) := \sum_{e \in R} c_e$. The problem is NP-hard via reduction from Hamiltonian cycle. We consider a particular relaxation for this problem that relates to graphs containing exactly one cycle.

Definition 6.14. A 1-tree with respect to vertex $v_1 \in V$ in graph $G = (V, E)$ is a subset of edges $B \subseteq E$ such that

- $B \setminus \delta(v_1)$ is a spanning tree of $G \setminus \{v_1\}$, and
- the degree of v_1 in (V, B) is 2.

In the following, we fix a vertex $v_1 \in V$. Because every tour is a 1-tree, we have

$$\min_{B \text{ 1-tree}} c(B) \leq \min_{R \text{ tour}} c(R). \quad (6.9)$$

This lower bound for TSP can be improved as follows. We consider a *vertex potential* $\pi \in \mathbb{R}^V$ and define the *reduced costs* with respect to this potential as (see *Combinatorial Optimization*)

$$c_e^\pi := c_e - \pi_v - \pi_u \quad \forall e = \{u, v\} \in E.$$

Let $\pi(V) := \sum_{v \in V} \pi_v$ and $c^\pi(B) := \sum_{e \in B} c_e^\pi$. We obtain a stronger inequality.

Theorem 6.15 ([24]). For $f(\pi) := 2\pi(V) + \min_{B \text{ 1-tree}} c^\pi(B)$ it holds that

$$\max_{\pi \in \mathbb{R}^V} f(\pi) \leq \min_{R \text{ tour}} c(R).$$

Proof. Let $\pi \in \mathbb{R}^V$ be arbitrary and R be a TSP tour. Then,

$$\begin{aligned} f(\pi) &= 2\pi(V) + \min_{B \text{ 1-tree}} c^\pi(B) \\ &\stackrel{(6.9)}{\leq} 2\pi(V) + c^\pi(R) \\ &= 2\pi(V) + c(R) - 2\pi(V) = c(R). \end{aligned}$$

Here, the first inequality uses that R is a 1-tree, and the next equality holds because every vertex occurs exactly twice in the edges of R . \square

The function f is the sum of an affine part and the minimum of affine (i.e., concave) functions, hence it is concave. This means that the subgradient method can be applied to the convex function $-f$. We can determine a subgradient as follows.

Lemma 6.16. Let $\tilde{\pi} \in \mathbb{R}^V$ and $\tilde{B} \in \arg \min_{B \text{ 1-tree}} c^{\tilde{\pi}}(B)$. Then, a subgradient \mathbf{h} for $-f$ at point $\tilde{\pi}$ is defined by

$$h_v = \deg_{(V, \tilde{B})}(v) - 2 \quad \forall v \in V.$$

Proof. We need to show that $-f(\pi) \geq -f(\tilde{\pi}) + \mathbf{h}^\top(\pi - \tilde{\pi})$ holds for all $\pi \in \mathbb{R}^V$. We have

$$\begin{aligned} f(\tilde{\pi}) &= 2\tilde{\pi}(V) + \min_{B \text{ 1-tree}} c^{\tilde{\pi}}(B) \\ &= 2\tilde{\pi}(V) + c(\tilde{B}) - \sum_{v \in V} \deg_{(V, \tilde{B})}(v) \cdot \tilde{\pi}_v \\ &= c(\tilde{B}) + \sum_{v \in V} (2 - \deg_{(V, \tilde{B})}(v)) \cdot \tilde{\pi}_v \\ &= c(\tilde{B}) - \mathbf{h}^\top \tilde{\pi}. \end{aligned}$$

Moreover,

$$f(\boldsymbol{\pi}) = 2\pi(V) + \min_{B \text{ 1-tree}} c^\pi(B) \leq 2\pi(V) + c^\pi(\tilde{B}) = c(\tilde{B}) - \mathbf{h}^\top \boldsymbol{\pi}.$$

Overall, we obtain

$$-f(\boldsymbol{\pi}) + f(\tilde{\boldsymbol{\pi}}) \geq (\mathbf{h}^\top \boldsymbol{\pi} - c(\tilde{B})) + (c(\tilde{B}) - \mathbf{h}^\top \tilde{\boldsymbol{\pi}}) = \mathbf{h}^\top (\boldsymbol{\pi} - \tilde{\boldsymbol{\pi}}). \quad \square$$

This means that we can apply the subgradient method, provided that we can efficiently compute minimum cost 1-trees. This can easily be accomplished by computing a minimum spanning tree of $G \setminus \{v_1\}$ (see *Algorithmic Discrete Mathematics*) and then adding the two edges of smallest weight incident to vertex v_1 to the result.

To relate this approach to a Lagrangian relaxation, we need the following result (without proof). Here $\boldsymbol{\chi}^B \in \{0, 1\}^E$ denotes the incidence vector of a set $B \subseteq E$, i.e., $\chi_e^B = 1 \Leftrightarrow e \in B$.

Lemma 6.17 ([23]). It holds that

$$\text{conv}(\{\boldsymbol{\chi}^B : B \subseteq E \text{ is a 1-tree wrt. } v_1\}) = \{\mathbf{x} \in [0, 1]^E : \begin{aligned} \sum_{e \in \delta(v_1)} x_e &= 2 \\ \sum_{e \in G[S]} x_e &\leq |S| - 1 \quad \forall S \subseteq V \setminus \{v_1\}, S \neq \emptyset, \\ \sum_{e \in E} x_e &= |V| \end{aligned}\}.$$

In particular, the polyhedron on the right-hand side is integral.

Together with Theorem 6.1 we can show a characterization of the strength of the lower bound of Theorem 6.15.

Theorem 6.18 ([23]). It holds that

$$\begin{aligned} \max_{\boldsymbol{\pi} \in \mathbb{R}^V} f(\boldsymbol{\pi}) &= \min \sum_{e \in E} c_e x_e \\ \text{s.t.} \quad \sum_{e \in \delta(v)} x_e &= 2 \quad \forall v \in V && \text{(degree constraints)} \\ \sum_{e \in G[S]} x_e &\leq |S| - 1 \quad \forall S \subset V, S \neq \emptyset && \text{(subtour-elimination)} \\ \mathbf{x} &\in [0, 1]^E. \end{aligned}$$

Proof. Let $A^{(2)}\mathbf{x} \leq \mathbf{b}^{(2)}$ be the system of Lemma 6.17 and let $A^{(1)}\mathbf{x} = \mathbf{b}^{(1)}$ describe the degree conditions

$$\sum_{e \in \delta(v)} x_e = 2 \quad \forall v \in V \setminus \{v_1\}.$$

We observe that, together, $A^{(1)}\mathbf{x} = \mathbf{b}^{(1)}$ and $A^{(2)}\mathbf{x} \leq \mathbf{b}^{(2)}$ yield the system in the statement of the theorem: The degree conditions are all present, the equation $\sum_{e \in E} x_e = |V|$ is redundant (sum of all degree conditions

divided by 2), and the missing subtour elimination constraints for $S \subseteq V$ with $v_1 \in S$ are implied via

$$\begin{aligned}
\sum_{e \in G[S]} x_e &= \sum_{e \in G[S]} x_e + \sum_{e \in E} x_e - |V| \\
&= 2 \sum_{e \in G[S]} x_e + \sum_{e \in \delta(S)} x_e + \sum_{e \in G[V \setminus S]} x_e - |V| \\
&= \sum_{v \in S} \sum_{e \in \delta(v)} x_e + \sum_{e \in G[V \setminus S]} x_e - |V| \\
&\stackrel{v_1 \in S}{\leq} 2|S| + (|V \setminus S| - 1) - |V| \\
&= |S| - 1,
\end{aligned}$$

where the inequality uses the degree constraints and the subtour elimination constraints for $v_1 \notin S$. In order to inspect the Lagrangian relaxation, we consider the negated integer program

$$\begin{aligned}
\max \quad & -\mathbf{c}^\top \mathbf{x} \\
\text{s.t.} \quad & A^{(1)} \mathbf{x} = \mathbf{b}^{(1)} \\
& A^{(2)} \mathbf{x} \leq \mathbf{b}^{(2)} \\
& \mathbf{x} \in \{0, 1\}^E,
\end{aligned}$$

where $A^{(1)}$ is the vertex-edge incidence matrix of G and $\mathbf{b}^{(1)} = 2 \cdot \mathbf{1}$. Note that the system is rational. Since we relax a system of equations $A^{(1)} \mathbf{x} = \mathbf{b}^{(1)}$, the Lagrangian function (6.2) may be considered without requiring non-negativity and may be written as

$$\begin{aligned}
L(\boldsymbol{\pi}) &= \max\{-\mathbf{c}^\top \mathbf{x} + \boldsymbol{\pi}^\top (A^{(1)} \mathbf{x} - \mathbf{b}^{(1)}) : \mathbf{x} \in \mathcal{P}(A^{(2)}, \mathbf{b}^{(2)}) \cap \{0, 1\}^E\} \\
&\stackrel{6.17}{=} \max\{-\mathbf{c}^\top \mathbf{x} + \boldsymbol{\pi}^\top (A^{(1)} \mathbf{x} - \mathbf{b}^{(1)}) : \mathbf{x} \in \text{conv}(\{\boldsymbol{\chi}^B : B \subseteq E \text{ is a 1-tree wrt. } v_1\})\} \\
&\stackrel{2.7}{=} -2\pi(V) - \min\{(\mathbf{c}^\top - \boldsymbol{\pi}^\top A^{(1)}) \boldsymbol{\chi}^B : B \subseteq E \text{ is a 1-tree wrt. } v_1\} \\
&= -2\pi(V) - \min_{1\text{-tree } B} \sum_{e=\{u,v\} \in E} (c_e - \pi_u - \pi_v) \chi_e^B \\
&= -f(\boldsymbol{\pi}).
\end{aligned}$$

Lemma 6.17 implies that $\mathcal{P}(A^{(2)}, \mathbf{b}^{(2)})$ is integral, and Observation 6.3 thus yields

$$\begin{aligned}
\max_{\boldsymbol{\pi} \in \mathbb{R}^V} f(\boldsymbol{\pi}) &= -\min_{\boldsymbol{\pi} \in \mathbb{R}^V} L(\boldsymbol{\pi}) \\
&\stackrel{6.3}{=} -\max\{-\mathbf{c}^\top \mathbf{x} : A^{(1)} \mathbf{x} = \mathbf{b}^{(1)}, A^{(2)} \mathbf{x} \leq \mathbf{b}^{(2)}, \mathbf{x} \in [0, 1]^E\} \\
&= \min\{\mathbf{c}^\top \mathbf{x} : A^{(1)} \mathbf{x} = \mathbf{b}^{(1)}, A^{(2)} \mathbf{x} \leq \mathbf{b}^{(2)}, \mathbf{x} \in [0, 1]^E\}. \quad \square
\end{aligned}$$

The proof of Theorem 6.18 shows that the 1-tree relaxation is a Lagrangian relaxation of the system in the theorem, where the degree constraints are relaxed for all vertices except v_1 . Note that $L(\boldsymbol{\pi}) = -f(\boldsymbol{\pi})$ also implies that Lemma 6.16 follows immediately from Lemma 6.12.

Remark 6.19. If we add to the LP in Theorem 6.18 the condition $\mathbf{x} \in \{0, 1\}^E$, we obtain an IP formulation of the TSP problem. The degree conditions guarantee that every vertex has degree 2 in the solution. The subtour elimination constraints exclude partial tours that are disconnected from the remaining graph. Theorem 6.18 shows that the Lagrangian relaxation of this formulation is as good as the corresponding LP relaxation.

6.2 Dantzig-Wolfe decomposition

The Dantzig-Wolfe decomposition (see [12]) is a different method for treating MIPs with two parts, i.e., of the form

$$\begin{aligned}
 \max \quad & \mathbf{c}^\top \mathbf{x} \\
 \text{s.t.} \quad & A^{(1)} \mathbf{x} \leq \mathbf{b}^{(1)} \\
 & A^{(2)} \mathbf{x} \leq \mathbf{b}^{(2)} \\
 & \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p},
 \end{aligned} \tag{6.10}$$

with $A^{(1)} \in \mathbb{Q}^{m_1 \times n}$, $A^{(2)} \in \mathbb{Q}^{m_2 \times n}$, $\mathbf{b}^{(1)} \in \mathbb{Q}^{m_1}$, $\mathbf{b}^{(2)} \in \mathbb{Q}^{m_2}$ and $m_1 + m_2 = m$.

We first assume that $p = 0$, i.e., we have an LP. Consider the polyhedron $\mathcal{P}^{(2)} = \mathcal{P}(A^{(2)}, \mathbf{b}^{(2)})$. By Theorem 2.19, there exist vectors $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}$ and $\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(\ell)}$ with $\mathcal{P}^{(2)} = \text{conv}(\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}\}) + \text{cone}(\{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(\ell)}\})$. In other words, $\mathbf{x} \in \mathcal{P}^{(2)}$ can be expressed in the form

$$\mathbf{x} = \sum_{i=1}^k \lambda_i \mathbf{v}^{(i)} + \sum_{j=1}^{\ell} \mu_j \mathbf{r}^{(j)} \tag{6.11}$$

with $\lambda, \mu \geq \mathbf{0}$ and $\mathbf{1}^\top \lambda = 1$. With (6.11), we can rewrite (6.10) as

$$\begin{aligned}
 \max \quad & \mathbf{c}^\top \left(\sum_{i=1}^k \lambda_i \mathbf{v}^{(i)} + \sum_{j=1}^{\ell} \mu_j \mathbf{r}^{(j)} \right) \\
 \text{s.t.} \quad & A^{(1)} \left(\sum_{i=1}^k \lambda_i \mathbf{v}^{(i)} + \sum_{j=1}^{\ell} \mu_j \mathbf{r}^{(j)} \right) \leq \mathbf{b}^{(1)}, \\
 & \mathbf{1}^\top \lambda = 1, \\
 & \lambda, \mu \geq \mathbf{0},
 \end{aligned}$$

which is equivalent to

$$\begin{aligned}
 \max \quad & \sum_{i=1}^k (\mathbf{c}^\top \mathbf{v}^{(i)}) \lambda_i + \sum_{j=1}^{\ell} (\mathbf{c}^\top \mathbf{r}^{(j)}) \mu_j \\
 \text{s.t.} \quad & \sum_{i=1}^k (A^{(1)} \mathbf{v}^{(i)}) \lambda_i + \sum_{j=1}^{\ell} (A^{(1)} \mathbf{r}^{(j)}) \mu_j \leq \mathbf{b}^{(1)}, \\
 & \mathbf{1}^\top \lambda = 1, \\
 & \lambda \geq \mathbf{0}, \mu \geq \mathbf{0}.
 \end{aligned} \tag{6.12}$$

The formulation (6.12) is called the *master problem* of (6.10). When comparing the two formulations, we realize that we reduced the number of constraints from m to $m_1 + 1$, but we now have $k + \ell$ variables instead of n . The number $k + \ell$ can be very large compared to n , even exponential, e.g., for the unit cube in \mathbb{R}^n with $2n$ constraints and 2^n vertices. At first glance, the advantage of using (6.12) is not apparent.

The key insight is that we can use the simplex method to solve (6.12) without having to create all variables explicitly. We abbreviate (6.12) by

$$\max \{ \tilde{\mathbf{c}}^\top \mathbf{z} : D\mathbf{z} = \mathbf{d}, \mathbf{z} \geq \mathbf{0} \},$$

where (including additional slack variables)

$$D = \left(\begin{array}{c|c|c} A^{(1)}\mathbf{v}^{(i)} & A^{(1)}\mathbf{r}^{(j)} & \mathbb{I} \\ \hline \mathbf{1}^\top & \mathbf{0}^\top & \mathbf{0}^\top \end{array} \right) \in \mathbb{Q}^{(m_1+1) \times (k+\ell+m_1)}, \quad \tilde{\mathbf{c}} = \begin{pmatrix} \mathbf{c}^\top \mathbf{v}^{(i)} \\ \mathbf{c}^\top \mathbf{r}^{(j)} \\ \mathbf{0} \end{pmatrix} \in \mathbb{Q}^{k+\ell+m_1}, \quad \mathbf{d} = \begin{pmatrix} \mathbf{b}^{(1)} \\ \mathbf{1} \end{pmatrix} \in \mathbb{Q}^{m_1+1}.$$

Recall that the simplex method (for LPs in standard form) works with a feasible basis $B \subseteq \{1, \dots, k+\ell+m_1\}$, $|B| = m_1 + 1$, where D_B is regular and yields a feasible solution $\mathbf{z}_B^* = D_B^{-1}\mathbf{d}$ and $\mathbf{z}_N^* = \mathbf{0}$ with $N = \{1, \dots, k+\ell+m_1\} \setminus B$. Note that $D_B \in \mathbb{R}^{(m_1+1) \times (m_1+1)}$ is smaller than a basis for the original system (6.10) would have to be, and that only a part of the variables ($m_1 + 1$ of $k + \ell + m_1$) are non-zero. In addition, the only operation in the simplex method that uses all columns is the *pricing* step, which checks whether the reduced costs satisfy $\tilde{\mathbf{c}}_N^\top - \tilde{\mathbf{y}}^\top D_N \leq \mathbf{0}$, where $\tilde{\mathbf{y}} = (D_B)^{-1}\tilde{\mathbf{c}}_B$ is the solution of $\mathbf{y}^\top D_B = \tilde{\mathbf{c}}_B^\top$. Recall that $\tilde{\mathbf{y}}$ is the dual solution, and the condition on the reduced costs asks whether $\tilde{\mathbf{y}}$ is dually feasible. We need to find a variable of positive reduced cost that can enter the basis, or, equivalently, we are looking for a violated inequality of the dual. This can be accomplished via the following *pricing problem*:

$$\begin{aligned} \max \quad & (\mathbf{c}^\top - \tilde{\mathbf{y}}^\top A^{(1)})\mathbf{x} \\ \text{s.t.} \quad & A^{(2)}\mathbf{x} \leq \mathbf{b}^{(2)}, \\ & \mathbf{x} \in \mathbb{R}^n, \end{aligned} \tag{6.13}$$

where $\tilde{\mathbf{y}}$ represents the first m_1 components of the solution $\tilde{\mathbf{y}} = \begin{pmatrix} \tilde{\mathbf{y}} \\ \tilde{y}_{m_1+1} \end{pmatrix}$. The following cases can occur.

- (a) A slack variable has positive reduced cost. In that case, we can let the variable enter the basis. This can be checked directly, because there are only $m_1 + 1$ slack variables.
- (b) (6.13) has a vertex solution $\bar{\mathbf{x}}$ with $(\mathbf{c}^\top - \tilde{\mathbf{y}}^\top A^{(1)})\bar{\mathbf{x}} > \tilde{y}_{m_1+1}$.

In this case, $\bar{\mathbf{x}} = \mathbf{v}^{(i)}$ for some $i \in \{1, \dots, k\}$, corresponding to $\boldsymbol{\lambda} = \mathbf{e}^i$ and $\boldsymbol{\mu} = \mathbf{0}$ in (6.12). The associated column D_i has reduced costs

$$\tilde{c}_i - \tilde{\mathbf{y}}^\top D_i = \mathbf{c}^\top \mathbf{v}^{(i)} - \tilde{\mathbf{y}}^\top \begin{pmatrix} A^{(1)}\mathbf{v}^{(i)} \\ \mathbf{1} \end{pmatrix} = \mathbf{c}^\top \mathbf{v}^{(i)} - \tilde{\mathbf{y}}^\top A^{(1)}\mathbf{v}^{(i)} - \tilde{y}_{m_1+1} > 0.$$

In other words, $D_i = \begin{pmatrix} A^{(1)}\mathbf{v}^{(i)} \\ \mathbf{1} \end{pmatrix}$ can be used as the entering column in the simplex algorithm (in particular, $i \in N$).

- (c) (6.13) is unbounded.

Here we obtain a feasible extreme ray $\bar{\mathbf{r}} = \mathbf{r}^{(j)}$ with $(\mathbf{c}^\top - \tilde{\mathbf{y}}^\top A^{(1)})\bar{\mathbf{r}} > 0$ for some $j \in \{1, \dots, \ell\}$, corresponding to $\boldsymbol{\lambda} = \mathbf{0}$ and $\boldsymbol{\mu} = \mathbf{e}^j$ in (6.12). The corresponding column $D_{(k+j)}$ has reduced costs

$$\tilde{c}_{k+j} - \tilde{\mathbf{y}}^\top D_{(k+j)} = \mathbf{c}^\top \mathbf{r}^{(j)} - \tilde{\mathbf{y}}^\top \begin{pmatrix} A^{(1)}\mathbf{r}^{(j)} \\ \mathbf{0} \end{pmatrix} = \mathbf{c}^\top \mathbf{r}^{(j)} - \tilde{\mathbf{y}}^\top (A^{(1)}\mathbf{r}^{(j)}) > 0.$$

Thus, $D_{(k+j)} = \begin{pmatrix} A^{(1)}\mathbf{r}^{(j)} \\ \mathbf{0} \end{pmatrix}$ can be chosen as the entering column.

- (d) (6.13) has an optimum solution $\bar{\mathbf{x}}$ with $(\mathbf{c}^\top - \tilde{\mathbf{y}}^\top A^{(1)})^\top \bar{\mathbf{x}} \leq \tilde{y}_{m_1+1}$.

In this case, using the same arguments as in (b) and (c), we obtain that $\tilde{c}_i - \tilde{\mathbf{y}}^\top D_i \leq 0$ for all $i \in \{1, \dots, k+\ell\}$. Since no slack variable has positive reduced cost either, $\tilde{\mathbf{c}}_N^\top - \tilde{\mathbf{y}}^\top D_N \leq \mathbf{0}$ holds, and thus \mathbf{z}^* is an optimum solution of the master problem (6.12).

With these arguments, a variable of positive reduced costs can be computed, if any exists, without having to evaluate the reduced costs of all (non-basis) variables.

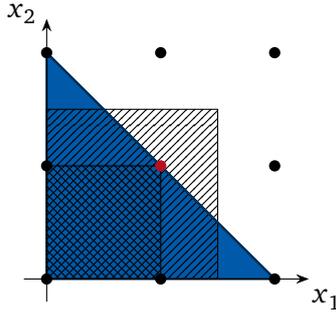


Figure 6.3: Example for the extension of the Dantzig-Wolfe decomposition to ILPs in Remark 6.20. The integral optimum solution $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ cannot be written as an integer linear combination of the vertices of $\mathcal{P}^{(2)}$ (blue).

Note that the whole problem (6.1) is split into two subproblems, i.e., into (6.12) and (6.13) and the approach alternates between working at the higher level (6.12) and at the lower level (6.13). The procedure starts with a feasible solution for (6.12) and generates new promising columns on demand by solving (6.13). Such methods are called *(delayed) column-generation algorithms*.

The approach can also be applied to ILPs with some caution. In this case, the problem (6.13) changes from an LP to an ILP. In addition we have to make sure that in (6.11) all feasible integral solutions x of (6.1) can be generated by linear combinations of the vectors $v^{(1)}, \dots, v^{(k)}$ and $r^{(1)}, \dots, r^{(\ell)}$ with

$$\text{conv}(\{x \in \mathbb{Z}^n : A^{(2)}x \leq b^{(2)}\}) = \text{conv}(\{v^{(1)}, \dots, v^{(k)}\}) + \text{cone}(\{r^{(1)}, \dots, r^{(\ell)}\}).$$

On the other hand, we need to restrict to linear combinations that yield integer vectors x .

Remark 6.20. It is not sufficient to require that λ and μ be integral. As a counterexample, consider

$$A^{(1)} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad b^{(1)} = \begin{pmatrix} 3 \\ 3 \end{pmatrix} \quad A^{(2)} = \begin{pmatrix} 1 & 1 \\ -1 & 0 \\ 0 & -1 \end{pmatrix} \quad b^{(2)} = \begin{pmatrix} 2 \\ 0 \\ 0 \end{pmatrix},$$

and the problem

$$\max\{x_1 + x_2 : A^{(1)}x \leq b^{(1)}, A^{(2)}x \leq b^{(2)}, x \in \mathbb{Z}^2\}.$$

Then, $\mathcal{P}^{(2)} = \text{conv}(\{(0,0), (2,0), (0,2)\})$, see Figure 6.3. But the optimum solution $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ of the integer problem is not an integer linear combination of the vertices of $\mathcal{P}^{(2)}$. Thus, we lose optimum solutions if we do not handle the integrality of the variables in the master problem in some other way.

However, if all variables are binary, this difficulty does not arise, since every binary solution of the LP relaxation of a binary MIP is a vertex of the corresponding polyhedron (as it is an extreme point even for the hypercube and can thus not be written as nontrivial convex combination). In this case, it suffices to require $\lambda \in \{0, 1\}^k$ (no variables μ occur, since binary programs have a bounded feasible region). In fact, column-generation algorithms are not only used for solving large linear problems, but especially for large binary problems.

Of course, the Dantzig-Wolfe decomposition for linear or binary integer problems is only one possibility for a column-generation algorithm. Other algorithms do not solve the lower-order problem using general techniques for LPs or ILPs, but by combinatorial algorithms or exhaustive enumeration. Furthermore, the problems are often not modeled via (6.1), but rather directly as in (6.12). This is the case, for example, when the feasible region allows for a concise interior representation.

6.2.1 Application in combinatorial optimization*

As an example for the application of the Dantzig-Wolfe decomposition, we consider the the *bin-packing problem*: How many bins of fixed size $W \geq 0$ do we need to fit n items of sizes $w_i > 0$, $i \in \{1, \dots, n\}$? For clarity, let $m := n$ denote an upper bound on the number of bins required.

We can model this problem as a binary problem using variables $\mathbf{z} \in \{0, 1\}^m$ with $z_j = 1$ if and only if the j -th potential bin is used, and $\mathbf{x} \in \{0, 1\}^{\{1, \dots, n\} \times \{1, \dots, m\}}$ with $x_{ij} = 1$ if and only if item i is packed in bin j . We obtain the following assignment problem (see Example 1.5):

$$\begin{aligned} \min \quad & \mathbf{1}^\top \mathbf{z} \\ \text{s.t.} \quad & \sum_{j=1}^m x_{ij} = 1 & \forall i \in \{1, \dots, n\}, \end{aligned} \quad (6.14)$$

$$\sum_{i=1}^n w_i x_{ij} \leq W z_j \quad \forall j \in \{1, \dots, m\}, \quad (6.15)$$

$$\mathbf{x} \in \{0, 1\}^{n \times m}, \mathbf{z} \in \{0, 1\}^m.$$

However, this model can only be solved very inefficiently. The reason for this is the weakness of the LP relaxation: Optimum solutions $(\mathbf{x}^*, \mathbf{z}^*)^\top$ of the LP relaxation satisfy (6.15) with equality, therefore

$$\mathbf{1}^\top \mathbf{z}^* = \sum_{j=1}^m \sum_{i=1}^n \frac{w_i x_{ij}^*}{W} = \sum_{i=1}^n \frac{w_i}{W} \sum_{j=1}^m x_{ij}^* = \sum_{i=1}^n \frac{w_i}{W} \rightarrow 0 \text{ for } W \rightarrow \infty. \quad (6.16)$$

So, if W becomes very large, the optimum value of the LP relaxation is small (value 0) and says almost nothing about the integral solution (value 1 for $W \geq \sum_{i=1}^n w_i$).

As a remedy, we can use Dantzig-Wolfe decomposition to eliminate the knapsack conditions (6.15). Every opened bin yields the same knapsack polytope

$$\mathcal{P} := \text{conv}(\{\mathbf{y} \in \{0, 1\}^n : \mathbf{w}^\top \mathbf{y} \leq W\}) = \text{conv}(\{\mathbf{0}, \mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}\}) \subset \mathbb{R}^n,$$

where $\mathbf{0}, \mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)} \in \{0, 1\}^n$ are the vertices of the (bounded) polytope \mathcal{P} , while every unopened bin only allows the vector $\mathbf{0}$. Each such vertex describes one *configuration* of items that can be packed into a bin. We introduce variables $\lambda_{j1}, \dots, \lambda_{jk}$ for every bin $j \in \{1, \dots, m\}$, that indicate, which vertex/configuration is selected for each bin, i.e., $\mathbf{x}_j = \sum_{\ell=1}^k \lambda_{j\ell} \mathbf{v}^{(\ell)}$. We obtain the master problem

$$\begin{aligned} \min \quad & \sum_{j=1}^m \sum_{\ell=1}^k \lambda_{j\ell} \\ \text{s.t.} \quad & \sum_{j=1}^m \sum_{\ell=1}^k \lambda_{j\ell} \mathbf{v}_i^{(\ell)} = 1 & \forall i \in \{1, \dots, n\} \end{aligned} \quad (6.17)$$

$$\begin{aligned} & \sum_{\ell=1}^k \lambda_{j\ell} \leq 1 & \forall j \in \{1, \dots, m\} \\ & \boldsymbol{\lambda} \in \{0, 1\}^{n \times k}. \end{aligned} \quad (6.18)$$

Note that in (6.18) we did not require equality to allow the vertex $\mathbf{0}$. Furthermore, we replaced the variable z_j by $\sum_{\ell=1}^k \lambda_{j\ell}$, as it makes no sense to use a bin that does not contain any items, i.e., we encode unopened bins j

via $\sum_{\ell=1}^k \lambda_{j\ell} = 0$. Similarly, we can observe that, because of $\mathbf{v}^{(\ell)} \neq \mathbf{0}$, it follows from (6.17) that $\sum_j \sum_{\ell} \lambda_{j\ell} \leq n$, since each $\lambda_{j\ell}$ belongs to a single bin and occurs in exactly one row of (6.17). Therefore, in each feasible solution, a maximum of n vertices (except $\mathbf{0}$) are used. Because of $m \geq n$, we can therefore omit (6.18), since we can convert every solution in such a way that no bin is used by more than one vertex by introducing separate bins for every vertex.

In addition, we define the set $S_i := \{\ell : v_i^{(\ell)} = 1\}$ of the indices of all vertices that contain item i , and introduce the variables

$$\delta_{\ell} := \sum_{j=1}^m \lambda_{j\ell},$$

which indicate whether or not the vertex $\mathbf{v}^{(\ell)}$ is used. Because of (6.17) we have $\boldsymbol{\delta} \in \{0, 1\}^k$. Using

$$\sum_{j=1}^m \sum_{\ell=1}^k \lambda_{j\ell} v_i^{(\ell)} = \sum_{\ell=1}^k v_i^{(\ell)} \sum_{j=1}^m \lambda_{j\ell} = \sum_{\ell \in S_i} \delta_{\ell},$$

we obtain

$$\begin{aligned} \min \quad & \mathbf{1}^{\top} \boldsymbol{\delta} \\ \text{s.t.} \quad & \sum_{\ell \in S_i} \delta_{\ell} = 1 \quad \forall i \in \{1, \dots, n\} \\ & \boldsymbol{\delta} \in \{0, 1\}^k. \end{aligned} \tag{6.19}$$

Here $\boldsymbol{\delta}$ encodes a partition of the items into configurations. The exact assignment to the bins is implicit.

Example 6.21. Let $m = n = 4$, $w_1 = \frac{3}{4}$, $w_2 = \frac{1}{4}$, $w_3 = \frac{1}{2}$, $w_4 = \frac{1}{2}$ and $W = 1$. A possible solution with two bins is

$$(x_{ij}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}.$$

The vertices of the knapsack polytope are

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \mathbf{v}^{(1)} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \mathbf{v}^{(2)} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \end{pmatrix}, \mathbf{v}^{(8)} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \end{pmatrix}.$$

Thus, $S_1 = \{1, 5\}$, $S_2 = \{2, 5, 6, 7\}$, $S_3 = \{3, 6, 8\}$ and $S_4 = \{4, 7, 8\}$. The solution is obtained by combining the vertices $\mathbf{v}^{(5)}$ and $\mathbf{v}^{(8)}$, i.e.,

$$\boldsymbol{\delta} = (0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1)^{\top},$$

which corresponds to opening 2 bins, one containing the items $\{1, 2\}$ (since $5 \in S_1 \cap S_2$), the other containing the items $\{3, 4\}$ (since $8 \in S_3 \cap S_4$). \triangle

The pricing problem aims to find a configuration to include in the current basis, or, equivalently, a violated constraint of the dual of (6.19). The dual is given by

$$\max\{\mathbf{1}^{\top} \mathbf{y} : \sum_{i: \ell \in S_i} y_i \leq 1 \ \forall \ell \in \{1, \dots, k\}\},$$

where $\{i : \ell \in S_i\}$ are the items used by the configuration corresponding to vertex $\mathbf{v}^{(\ell)}$. We can find a vertex violating its constraint for the current dual solution $\tilde{\mathbf{y}}$ via

$$\begin{aligned} \max \quad & \tilde{\mathbf{y}}^\top \mathbf{x} \\ \text{s.t.} \quad & \mathbf{w}^\top \mathbf{x} \leq W \\ & \mathbf{x} \in \{0, 1\}^m. \end{aligned}$$

This problem is bounded. If the optimum vertex solution $\bar{\mathbf{x}}$ satisfies $\tilde{\mathbf{y}}^\top \bar{\mathbf{x}} > 1$, it yields a column to add to the basis for the master problem, otherwise $\tilde{\mathbf{y}}$ is dually feasible and the current solution of the master problem is optimal.

In practice, this Dantzig-Wolfe formulation results in a very strong bound: Often the rounded value of the LP relaxation of the master problem (at the root of the branch-and-bound tree) already gives the optimum number of bins for currently allowed configurations (cf. (6.16) for the LP relaxation of the original formulation).

6.3 Benders' decomposition

Finally, we consider the *Benders' decomposition* (see [4]). Benders' decomposition also eliminates part of the constraint matrix, but, unlike the Dantzig-Wolfe decomposition, where we delete a part of the constraints and reinsert them by means of column generation, we now delete part of the variables and reinsert them using cutting planes. From this perspective, Benders' decomposition is equivalent to Dantzig-Wolfe decomposition applied to the dual problem (see Section 6.4). Consider a MIP of the form

$$\begin{aligned} \max \quad & (\mathbf{c}^{(1)})^\top \mathbf{x}^{(1)} + (\mathbf{c}^{(2)})^\top \mathbf{x}^{(2)} \\ \text{s.t.} \quad & A^{(1)}\mathbf{x}^{(1)} + A^{(2)}\mathbf{x}^{(2)} \leq \mathbf{b}, \\ & \mathbf{x}^{(1)} \in \mathbb{Z}^{n_1}, \mathbf{x}^{(2)} \in \mathbb{R}^{n_2}, \end{aligned} \tag{6.20}$$

where $A = [A^{(1)}, A^{(2)}] \in \mathbb{Q}^{m \times n}$, $A^{(1)} \in \mathbb{Q}^{m \times n_1}$, $A^{(2)} \in \mathbb{Q}^{m \times n_2}$, $\mathbf{c}^{(1)} \in \mathbb{Q}^{n_1}$, $\mathbf{c}^{(2)} \in \mathbb{Q}^{n_2}$ with $n_1 + n_2 = n$. Note that, for the sake of simplicity, we assume that all integer variables are in $\mathbf{x}^{(1)}$. We want to eliminate the variables $\mathbf{x}^{(2)}$ via projection.

To apply projection we start by reformulating (6.20) to

$$\begin{aligned} \max \quad & z + \mathbf{0}^\top \mathbf{x}^{(2)} \\ \text{s.t.} \quad & z \leq (\mathbf{c}^{(1)})^\top \mathbf{x}^{(1)} + (\mathbf{c}^{(2)})^\top \mathbf{x}^{(2)} \\ & A^{(1)}\mathbf{x}^{(1)} + A^{(2)}\mathbf{x}^{(2)} \leq \mathbf{b} \\ & z \in \mathbb{R}, \mathbf{x}^{(1)} \in \mathbb{Z}^{n_1}, \mathbf{x}^{(2)} \in \mathbb{R}^{n_2}. \end{aligned} \tag{6.21}$$

By strong duality, this is equivalent to

$$\begin{aligned} & \max_{\mathbf{x}^{(1)} \in \mathbb{Z}^{n_1}, z \in \mathbb{R}} \left\{ z + \max_{\mathbf{x}^{(2)} \in \mathbb{R}^{n_2}} \{ \mathbf{0}^\top \mathbf{x}^{(2)} : -(\mathbf{c}^{(2)})^\top \mathbf{x}^{(2)} \leq (\mathbf{c}^{(1)})^\top \mathbf{x}^{(1)} - z, A^{(2)}\mathbf{x}^{(2)} \leq \mathbf{b} - A^{(1)}\mathbf{x}^{(1)} \} \right\} \\ = & \max_{\mathbf{x}^{(1)} \in \mathbb{Z}^{n_1}, z \in \mathbb{R}} \left\{ z + \min_{\gamma \in \mathbb{R}, \mathbf{v} \in \mathbb{R}^m} \{ \gamma((\mathbf{c}^{(1)})^\top \mathbf{x}^{(1)} - z) + \mathbf{v}^\top (\mathbf{b} - A^{(1)}\mathbf{x}^{(1)}) : \begin{pmatrix} \gamma \\ \mathbf{v} \end{pmatrix} \in \mathcal{C} \} \right\}, \end{aligned} \tag{6.22}$$

with

$$\mathcal{C} = \left\{ \begin{pmatrix} \gamma \\ \mathbf{v} \end{pmatrix} \in \mathbb{R}^{m+1} : \mathbf{v}^\top A^{(2)} - \gamma(\mathbf{c}^{(2)})^\top = \mathbf{0}^\top, \gamma \geq 0, \mathbf{v} \geq \mathbf{0} \right\}. \tag{6.23}$$

Since the primal objective function is $\mathbf{0}^\top \mathbf{x}^{(2)} = 0$, we have that (6.21) is feasible for some z if and only if the value of the inner minimization problem in (6.22) is 0. Otherwise, since $\mathbf{0} \in \mathcal{C} \neq \emptyset$ the inner minimum is unbounded. Thus, (6.22) is equivalent to

$$\max_{\mathbf{x}^{(1)} \in \mathbb{Z}^{n_1}, z \in \mathbb{R}} \{z : \gamma((\mathbf{c}^{(1)})^\top \mathbf{x}^{(1)} - z) + \mathbf{v}^\top (\mathbf{b} - A^{(1)} \mathbf{x}^{(1)}) \geq 0 \quad \forall \begin{pmatrix} \gamma \\ \mathbf{v} \end{pmatrix} \in \mathcal{C}\}. \quad (6.24)$$

Since \mathcal{C} is a polyhedral cone, by Theorem 2.17, there are vectors $\begin{pmatrix} \gamma_1 \\ \mathbf{v}^{(1)} \end{pmatrix}, \dots, \begin{pmatrix} \gamma_s \\ \mathbf{v}^{(s)} \end{pmatrix}$ with

$$\mathcal{C} = \text{cone}(\{\begin{pmatrix} \gamma_1 \\ \mathbf{v}^{(1)} \end{pmatrix}, \dots, \begin{pmatrix} \gamma_s \\ \mathbf{v}^{(s)} \end{pmatrix}\}).$$

Because $\mathcal{C} \subseteq \mathbb{R}_{\geq 0}^m$, we can assume by rescaling, that $\gamma \in \{0, 1\}^s$, i.e.

$$\mathcal{C} = \text{cone}(\{\begin{pmatrix} 0 \\ \mathbf{v}^{(k)} \end{pmatrix} : k \in K\}) + \text{cone}(\{\begin{pmatrix} 1 \\ \mathbf{v}^{(j)} \end{pmatrix} : j \in J\}),$$

with $K \cup J = \{1, \dots, s\}$ and $K \cap J = \emptyset$. With this description of \mathcal{C} , (6.24) can be rewritten as

$$\begin{aligned} \max \quad & z \\ \text{s.t.} \quad & z \leq (\mathbf{c}^{(1)})^\top \mathbf{x}^{(1)} + (\mathbf{v}^{(j)})^\top (\mathbf{b} - A^{(1)} \mathbf{x}^{(1)}) \quad \forall j \in J, \\ & 0 \leq (\mathbf{v}^{(k)})^\top (\mathbf{b} - A^{(1)} \mathbf{x}^{(1)}) \quad \forall k \in K, \\ & z \in \mathbb{R}, \mathbf{x}^{(1)} \in \mathbb{Z}^{n_1}. \end{aligned} \quad (6.25)$$

The problem (6.25) is the *Benders' master problem*. It arises from (6.21) by orthogonal projection in the sense that $z, \mathbf{x}^{(1)}$ are feasible for (6.25) if and only if $\mathbf{x}^{(2)}$ exists, so that $z, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}$ are feasible for (6.21). The master problem only has $n_1 + 1$ variables instead of the $n_1 + n_2$ variables in (6.20). Furthermore, the MIP (6.20) has been converted into an integer program (6.25) with an additional continuous variable z . Nevertheless, (6.25) contains a large number of constraints – generally exponentially many in n .

To circumvent this problem, we solve Benders' master problem by means of an iterative procedure as follows: We start with a small subset of extreme rays of \mathcal{C} (possibly with the empty set) and optimize (6.25) over this subset (e.g., via Branch-and-Bound, see Chapter 3). We obtain an optimum solution $\hat{\mathbf{x}}^{(1)}, \hat{z}$ of the relaxed problem and have to check whether this solution satisfies all other inequalities in (6.25). This can be done using the *Benders' subproblem* (see (6.23) and (6.24))

$$\begin{aligned} \min \quad & \mathbf{v}^\top (\mathbf{b} - A^{(1)} \hat{\mathbf{x}}^{(1)}) - \gamma(\hat{z} - (\mathbf{c}^{(1)})^\top \hat{\mathbf{x}}^{(1)}) \\ \text{s.t.} \quad & \mathbf{v}^\top A^{(2)} - \gamma(\mathbf{c}^{(2)})^\top = \mathbf{0}^\top \\ & \gamma \geq 0, \mathbf{v} \geq \mathbf{0}. \end{aligned} \quad (6.26)$$

This problem is feasible since $\mathbf{0} \in \mathcal{C}$, and has an optimum solution of value 0 or is unbounded (see discussion of (6.22)). In the first case, $(\hat{\mathbf{x}}^{(1)}, \hat{z})$ satisfies all inequalities in (6.25) and we have solved (6.25) optimally and thus (6.20). In the other case, we obtain (e.g., via the Simplex method) an extreme ray $\begin{pmatrix} \hat{\gamma} \\ \hat{\mathbf{v}} \end{pmatrix}$ from (6.26) with

$$\hat{\mathbf{v}}^\top (\mathbf{b} - A^{(1)} \hat{\mathbf{x}}^{(1)}) - \hat{\gamma}(\hat{z} - (\mathbf{c}^{(1)})^\top \hat{\mathbf{x}}^{(1)}) < 0,$$

which, after rescaling to $\hat{\gamma} \in \{0, 1\}$, yields an inequality for (6.25), which is violated by $\hat{\mathbf{x}}^{(1)}, \hat{z}$. We add this cut to the Benders master problem (6.25) and iterate.

6.3.1 Application in stochastic optimization*

We return to the example of stochastic optimization from Section 6.1.3 and apply Benders' decomposition directly. Here, $A^{(1)}\mathbf{x}^{(1)} + A^{(2)}\mathbf{x}^{(2)} \leq \mathbf{b}$ consists of all blocks $T^{(\omega)}\mathbf{x} + W^{(\omega)}\mathbf{y}^{(\omega)} \leq \mathbf{h}^{(\omega)}$ for all $\omega \in \Omega$. This leaves $\mathbf{x}^{(1)} = \mathbf{x}$ in the master problem and $\mathbf{x}^{(2)} = (\mathbf{y}^{(\omega)})_{\omega \in \Omega}$ goes to the subproblem. The master problem can then be written as

$$\begin{aligned} \min \quad & \mathbf{c}^\top \mathbf{x} + \sum_{\omega \in \Omega} p_\omega z_\omega \\ \text{s.t.} \quad & z_\omega \geq (\mathbf{v}^{(\omega,j)})^\top (\mathbf{h}^{(\omega)} - T^{(\omega)}\mathbf{x}), & \forall \omega \in \Omega, j \in J_\omega, \\ & 0 \leq (\mathbf{v}^{(\omega,k)})^\top (\mathbf{h}^{(\omega)} - T^{(\omega)}\mathbf{x}), & \forall \omega \in \Omega, k \in K_\omega, \\ & A\mathbf{x} \leq \mathbf{b}, \quad \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p}. \end{aligned}$$

Note that we minimize, which leads to slight adjustments. The vectors $(\mathbf{v}^{(\omega,j)})_{j \in J_\omega}$ and $(\mathbf{v}^{(\omega,k)})_{k \in K_\omega}$ are the extreme rays and vertices of

$$C_\omega := \left\{ \begin{pmatrix} \gamma \\ \mathbf{v} \end{pmatrix} \in \mathbb{R}^{m+1} : \mathbf{v}^\top W^{(\omega)} - \gamma(\mathbf{d}^{(\omega)})^\top = \mathbf{0}^\top, \gamma \geq 0, \mathbf{v} \geq \mathbf{0} \right\}.$$

The resulting *L-method* (see [41]) can be accelerated using various techniques and is very effective in many applications. However, it is applicable in the presented form only if the second-stage variables are all continuous.

6.3.2 Application in combinatorial optimization*

We consider the *maximum feasible subsystem problem*: Given an infeasible system $\tilde{A}\mathbf{x} \leq \tilde{\mathbf{b}}$ ($\tilde{A} \in \mathbb{R}^{m \times n}$, $\tilde{\mathbf{b}} \in \mathbb{R}^m$), identify a feasible subsystem of largest possible size. This can be expressed as MIP via

$$\begin{aligned} \max \quad & \mathbf{1}^\top \mathbf{x}^{(1)} + \mathbf{0}^\top \mathbf{x}^{(2)} \\ \text{s.t.} \quad & \tilde{A}\mathbf{x}^{(2)} \leq \tilde{\mathbf{b}} + M(\mathbf{1} - \mathbf{x}^{(1)}) \\ & \mathbf{x}^{(1)} \in \{0, 1\}^m, \mathbf{x}^{(2)} \in \mathbb{R}^n. \end{aligned} \tag{6.27}$$

Here, M is a very large number that guarantees the following: If $x_i^{(1)} = 0$, then $\tilde{A}_i \mathbf{x} - \tilde{\mathbf{b}}_i \leq M$ for all interesting \mathbf{x} . This means that the binary variables $\mathbf{x}^{(1)}$ activate/deactivate the respective inequalities. Note that the auxiliary variable z is not used here, because $\mathbf{c}^{(2)} = \mathbf{0}$.

We can now reformulate (6.27) to

$$\begin{aligned} & \max_{\mathbf{x}^{(1)} \in \{0,1\}^m} \left\{ \mathbf{1}^\top \mathbf{x}^{(1)} + \max_{\mathbf{x}^{(2)} \in \mathbb{R}^n} \{ \mathbf{0}^\top \mathbf{x}^{(2)} : \tilde{A}\mathbf{x}^{(2)} \leq \tilde{\mathbf{b}} + M(\mathbf{1} - \mathbf{x}^{(1)}) \} \right\} \\ & = \max_{\mathbf{x}^{(1)} \in \{0,1\}^m} \left\{ \mathbf{1}^\top \mathbf{x}^{(1)} + \min_{\mathbf{v} \in \mathbb{R}^m} \{ \mathbf{v}^\top (\tilde{\mathbf{b}} + M(\mathbf{1} - \mathbf{x}^{(1)})) : \mathbf{v}^\top \tilde{A} = \mathbf{0}^\top, \mathbf{v} \geq \mathbf{0} \} \right\}. \end{aligned}$$

We then have the cone $\mathcal{C} = \{ \mathbf{v} : \mathbf{v}^\top \tilde{A} = \mathbf{0}^\top, \mathbf{v} \geq \mathbf{0} \}$ with rays $(\mathbf{v}^{(k)})_{k \in K}$, and therefore $\mathcal{C} = \text{cone}(\{ \mathbf{v}^{(k)} : k \in K \})$. The Benders' master problem (6.25) is then

$$\begin{aligned} \max \quad & \mathbf{1}^\top \mathbf{x}^{(1)} \\ \text{s.t.} \quad & 0 \leq (\mathbf{v}^{(k)})^\top (\tilde{\mathbf{b}} + M(\mathbf{1} - \mathbf{x}^{(1)})) \quad \forall k \in K, \\ & \mathbf{x}^{(1)} \in \{0, 1\}^m, \end{aligned}$$

and the Benders' subproblem (6.26) becomes

$$\begin{aligned} \min \quad & \mathbf{v}^\top (\tilde{\mathbf{b}} + M(\mathbf{1} - \hat{\mathbf{x}}^{(1)})) \\ \text{s.t.} \quad & \mathbf{v}^\top \tilde{\mathbf{A}} = \mathbf{0}^\top, \mathbf{v} \geq \mathbf{0}. \end{aligned}$$

This can then be used to apply the solution scheme described above. We discuss this implementation further. Let $\hat{\mathbf{x}}^{(1)} \in \{0, 1\}^m$ be the solution of the (incomplete) Benders' master problem and $\hat{\mathbf{v}}$ be the solution of the corresponding Benders' subproblem. If the optimum value of the subproblem is 0, we are done. Otherwise it is smaller 0 and the following applies.

Lemma 6.22. If $\hat{\mathbf{v}}^\top (\tilde{\mathbf{b}} + M(\mathbf{1} - \hat{\mathbf{x}}^{(1)})) < 0$, then the rows of $\tilde{\mathbf{A}}\mathbf{x} \leq \tilde{\mathbf{b}}$ form an invalid subsystem with respect to $S := \{i : \hat{v}_i > 0\}$.

Proof. According to the Farkas lemma, it holds that $\tilde{\mathbf{A}}_S \mathbf{x} \leq \tilde{\mathbf{b}}_S$ is infeasible (exactly) if $\mathbf{y} \geq \mathbf{0}$ exists with $\mathbf{y}^\top \tilde{\mathbf{b}}_S < 0$ and $\mathbf{y}^\top \tilde{\mathbf{A}}_S = \mathbf{0}^\top$. With $\mathbf{y} := \hat{\mathbf{v}}_S$ we have $\mathbf{y} \geq \mathbf{0}$ and $\mathbf{y}^\top \tilde{\mathbf{A}}_S = \mathbf{0}^\top$, because of feasibility of $\hat{\mathbf{v}}$ for the Benders' subproblem. By assumption, $\hat{\mathbf{v}}^\top \tilde{\mathbf{b}} + M \hat{\mathbf{v}}^\top (\mathbf{1} - \hat{\mathbf{x}}^{(1)}) < 0$. Because M is sufficiently large, $\hat{\mathbf{x}}_S^{(1)} = \mathbf{1}$ must hold. It follows that $0 > \hat{\mathbf{v}}^\top \tilde{\mathbf{b}} + M \hat{\mathbf{v}}^\top (\mathbf{1} - \hat{\mathbf{x}}^{(1)}) = \hat{\mathbf{v}}^\top \tilde{\mathbf{b}} = \hat{\mathbf{v}}^\top \tilde{\mathbf{b}}_S = \mathbf{y}^\top \tilde{\mathbf{b}}_S$, which proves infeasibility. \square

A corresponding ray $\hat{\mathbf{v}}$ of the Benders' subproblem thus yields an infeasible subsystem in the part of $\tilde{\mathbf{A}}\mathbf{x} \leq \tilde{\mathbf{b}}$ in which $\hat{v}_i > 0$ holds. The corresponding inequalities in the master problem can therefore be reduced to

$$\sum_{i \in S} x_i^{(1)} \leq |S| - 1,$$

because at least one $x_i^{(1)}$ must be set to 0 (i.e., must be removed from the system).

6.4 Connections between these approaches

At first glance, Lagrangian relaxation, Dantzig-Wolfe decomposition and Benders' decomposition seem to be different approaches relying on different relaxations. However, they are strongly related to each other. Consider again (6.13), which, for a fixed $\bar{\mathbf{y}} \leq \mathbf{0}$, can be written as

$$\begin{aligned} \max (\mathbf{c}^\top - \bar{\mathbf{y}}^\top A^{(1)})\mathbf{x} &= \max \mathbf{c}^\top \mathbf{x} + \bar{\mathbf{y}}^\top (\mathbf{b}^{(1)} - A^{(1)}\mathbf{x}) - \bar{\mathbf{y}}^\top \mathbf{b}^{(1)} &= L(\bar{\mathbf{y}}) - \bar{\mathbf{y}}^\top \mathbf{b}^{(1)}, \\ \text{s.t. } \mathbf{x} \in \mathcal{P}^{(2)} & \text{s.t. } \mathbf{x} \in \mathcal{P}^{(2)} \end{aligned}$$

i.e., (6.3) and (6.13) define the same problems except for a constant shift $-\bar{\mathbf{y}}^\top \mathbf{b}$ in the objective. Furthermore, by replacing $\mathcal{P}^{(2)}$ with $\text{conv}(\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}\}) + \text{cone}(\{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(\ell)}\})$, it can be shown that (6.12) coincides with the right-hand side in Theorem 6.1 and thus with $\min_{\bar{\mathbf{y}}} L(\bar{\mathbf{y}})$. In other words, the Dantzig-Wolfe decomposition and the Lagrangian relaxations compute the same bound. The only differences are that, for the update of the dual variables, i.e., $\boldsymbol{\lambda}$ in the Lagrangian relaxation and $\bar{\mathbf{y}}$ in the Dantzig-Wolfe decomposition, subgradient methods are used in the former case, whereas LP techniques are applied in the latter.

Analogously, Benders' decomposition is nothing other than the Dantzig-Wolfe decomposition applied to the dual of (6.20). To see this, consider its dual LP

$$\begin{aligned} \min \quad & \mathbf{y}^\top \mathbf{b} \\ \text{s.t.} \quad & \mathbf{y}^\top A^{(1)} = (\mathbf{c}^{(1)})^\top, \\ & \mathbf{y}^\top A^{(2)} = (\mathbf{c}^{(2)})^\top, \\ & \mathbf{y} \geq \mathbf{0}. \end{aligned} \tag{6.28}$$

Now we write

$$\bar{\mathcal{P}}^{(2)} = \{\mathbf{y} \in \mathbb{R}^m : \mathbf{y}^\top A^{(2)} = (\mathbf{c}^{(2)})^\top, \mathbf{y} \geq \mathbf{0}\}$$

as

$$\bar{\mathcal{P}}^{(2)} = \text{conv}(\{\mathbf{v}^{(j)} : j \in J\}) + \text{cone}(\{\mathbf{v}^{(k)} : k \in K\}),$$

where K, J and $\mathbf{v}^{(\ell)}$ for $\ell \in K \cup J$ are exactly the quantities of (6.25) (set $\mathbf{v} = \gamma \mathbf{y}$, and note that $\text{hog}(\bar{\mathcal{P}}^{(2)}) := \{(\frac{1}{\gamma}) : \mathbf{y} \in \bar{\mathcal{P}}^{(2)}\}^{\circ\circ} = \mathcal{C}$ (exercise)), and write (6.28) as

$$\begin{aligned} \min \quad & \sum_{j \in J} ((\mathbf{v}^{(j)})^\top \mathbf{b}) \lambda_j + \sum_{k \in K} ((\mathbf{v}^{(k)})^\top \mathbf{b}) \mu_k \\ \text{s.t.} \quad & \sum_{j \in J} ((\mathbf{v}^{(j)})^\top A^{(1)}) \lambda_j + \sum_{k \in K} ((\mathbf{v}^{(k)})^\top A^{(1)}) \mu_k = (\mathbf{c}^{(1)})^\top, \\ & \sum_{j \in J} \lambda_j = 1, \\ & \boldsymbol{\lambda} \in \mathbb{R}_{\geq 0}^J, \boldsymbol{\mu} \in \mathbb{R}_{\geq 0}^K. \end{aligned} \tag{6.29}$$

We now conclude from the results of Section 6.2 that (6.29) is the master problem of (6.28). Dualization of (6.29) results in

$$\begin{aligned} \max \quad & (\mathbf{c}^{(1)})^\top \mathbf{x}^{(1)} + z \\ \text{s.t.} \quad & (\mathbf{v}^{(j)})^\top (A^{(1)} \mathbf{x}^{(1)} - \mathbf{b}) \leq -z \quad \forall j \in J, \\ & (\mathbf{v}^{(k)})^\top (A^{(1)} \mathbf{x}^{(1)} - \mathbf{b}) \leq 0 \quad \forall k \in K, \end{aligned}$$

which is equivalent to (6.25), i.e., to Benders' master problem (6.20). In other words, for LPs, Benders' and Dantzig-Wolfe decomposition result the same bound, which, according to the considerations above, is equal to the value of the Lagrangian relaxation (6.3).

7 Heuristics

In this chapter, we tackle the question how to quickly find good feasible solutions for a mixed integer program, e.g., for the acceleration of the branch-and-bound method (see Chapter 3). Also, the size of practical instances often does not allow for exact but exponential-time solution methods. In both cases, we are willing to sacrifice optimality for the sake of efficiency. Fast algorithms that come without guarantees regarding solution quality are called *heuristics*. In contrast, in Chapter 8, we will see *approximation algorithms*, which, similarly, sacrifice optimality for the sake of efficiency, but come with a provable solution quality.

Heuristics are usually tailored to exploit the specific structure of the problem at hand and to the distribution of the instances that we can expect. That being said, in the following we present some general methods, which are used as a basis for many heuristics in practice.

7.1 The greedy algorithm

The greedy algorithm constructs a solution incrementally by iteratively adding an element to the solution that increases the objective function the most / the least. In order for this approach to be applicable, we need that every feasible solution can be constructed incrementally. This can elegantly be expressed by considering set systems of feasible solutions that are closed under taking subsets.

Definition 7.1. An *independence system* is a pair (E, \mathcal{I}) consisting of a finite ground set E and a family of subsets $\mathcal{I} \subseteq 2^E$ such that $\mathcal{I} \neq \emptyset$ and $I \in \mathcal{I}$ for all $I \subseteq J \in \mathcal{I}$. A set $I \subseteq E$ is *independent* if $I \in \mathcal{I}$, and a *basis* of $J \subseteq E$ if $I' \notin \mathcal{I}$ for all $I \subsetneq I' \subseteq J$.

This abstract framework captures many classical optimization problems.

Example 7.2. The following are independence systems with the set of bases of E denoted by $\mathcal{B} \subseteq \mathcal{I}$.

- linear independence: $E = \{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}\} \subseteq \mathbb{R}^n$ and \mathcal{I} consists of all sets of linearly independent vectors. Then, $\max_{I \in \mathcal{I}} |I|$ is the problem of finding a basis of $\text{lin}(\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)})$.
- knapsack: $E = \{1, \dots, n\}$ and $\mathcal{I} = \{I \subseteq E : \sum_{i \in I} a_i \leq \beta\}$. Then, $\max_{I \in \mathcal{I}} \sum_{i \in I} c_i$ is the knapsack problem (see Example 1.6).
- stable sets: $E = V$ are the vertices of a graph G and $\mathcal{I} = \{I \subseteq E : G[I] = (I, \emptyset)\}$. Then, $\max_{I \in \mathcal{I}} |I|$ is the stable set problem (see *Algorithmic Discrete Mathematics*).
- cycles: E are the edges of a graph and $\mathcal{I} = \{I \subseteq E : I \text{ has no cycle}\}$. Then, $\min_{B \in \mathcal{B}} \sum_{e \in B} c(e)$ is the minimum spanning tree problem (see *Algorithmic Discrete Mathematics*).
- tours: E are the edges of a complete graph and $\mathcal{I} = \{I \subseteq E : \exists \text{ tour } T \supseteq I \text{ in the graph}\}$. Then, $\min_{B \in \mathcal{B}} \sum_{e \in B} c(e)$ is the traveling salesperson problem (see Section 6.1.4). \triangle

The greedy algorithm can abstractly be formulated as follows.

Algorithm: greedy algorithm for maximization or minimization

input: independence system (E, \mathcal{I}) , objective $c: E \rightarrow \mathbb{R}_{\geq 0}$

output: basis of (E, \mathcal{I})

$I \leftarrow \emptyset, E_0 \leftarrow \emptyset, j \leftarrow 0$

while $E_j \neq E$:

$e \leftarrow \arg \max[\text{or min}]\{c(e') : e' \in E \setminus E_j\}$

$E_{j+1} \leftarrow E_j \cup \{e\}, j \leftarrow j + 1$

if $I \cup \{e\} \in \mathcal{I}$:

$I \leftarrow I \cup \{e\}$

return I

By definition, the greedy algorithm maintains the invariant that, no previously considered element can be included in the solution.

Observation 7.3. The solution computed by the greedy algorithm contains a basis of E_j for all $j \in \{1, \dots, |E|\}$.

In some cases, the greedy algorithm computes an optimum solution, as we will see in more detail in the next section.

Example 7.4. In the cycle-system of Example 7.2, the greedy algorithm is equivalent to Kruskal's algorithm for computing a minimum spanning tree (see *Algorithmic Discrete Mathematics*). \triangle

On the other hand, the greedy algorithm sometimes performs very poorly.

Example 7.5. In the knapsack-system of Example 7.2 with costs $c = (n-1)\mathbf{1} + \mathbf{e}^1$, weights $\mathbf{a} = \mathbf{1} + (n-2)\mathbf{e}^1$ and capacity $\beta = n-1$, the greedy algorithm finds the solution $I = \{1\}$ of value n , while the optimum solution $I^* = \{2, \dots, n\}$ has value $(n-1)^2$. \triangle

7.1.1 Matroids

We have seen that the quality of the solutions produced by the greedy algorithm can vary wildly, depending on the application (see Examples 7.4 and 7.5). Remarkably, the cases when the greedy algorithm finds an optimum solution have an elegant characterization in terms of the following structure.

Definition 7.6. An independence system (E, \mathcal{I}) is a *matroid* if for all $F \subseteq E$ with bases $\mathcal{B}(F)$ it holds that

$$r(F) := \max\{|B| : B \in \mathcal{B}(F)\} = \min\{|B| : B \in \mathcal{B}(F)\} =: \rho(F),$$

where $r(F)$ is the *rank* of F and $\rho(F)$ is its *lower rank*.

In particular, the independence system underlying the minimum spanning tree problem is a matroid.

Example 7.7. The following independence systems of Example 7.2 are matroids.

- *vector matroid:* $E = \{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}\} \subseteq \mathbb{R}^n$ and \mathcal{I} consists of all sets of linearly independent vectors, since all bases of a subset $F \subseteq E$ have cardinality $\dim(\text{lin}(F))$.

- *cycle matroid*: E are the edges of a graph and $\mathcal{I} = \{I \subseteq E : I \text{ has no cycle}\}$, since all bases of a subset $F \subseteq E$ have cardinality $|F| - k$, where k is the number of connected components of F . \triangle

Recall that the greedy algorithm is optimal for the cycle matroid (see 7.4). On the other hand, the knapsack system, where the greedy algorithm does not yield a good solution (see 7.5), is not a matroid, and neither is the system for TSP.

Example 7.8. The following independence systems of Example 7.2 are not matroids.

- knapsack system: Not all maximal packings involve the same number of items, i.e., not all bases of E have the same cardinality.
- tours system: While all bases of E have the same cardinality, the same is not true for the bases of subsets of E . To see this, consider the complete graph on 4 vertices. The subgraph obtained by omitting an arbitrary edge has a cycle of length 4 and a path of length 3 as its bases. \triangle

We now more generally characterize the solution quality of the greedy algorithm via the *rank quotient* $q := \max_{F \subseteq E: r(F) > 0} r(F)/\rho(F)$.

Theorem 7.9 ([27]). Let (E, \mathcal{I}) be an independence system, $c: E \rightarrow \mathbb{R}_{\geq 0}$ be an objective to be maximized, $I^* \in \arg \max_{I \in \mathcal{I}} c(I)$ be an optimum solution with $c(I^*) > 0$, and I be the solution found by the greedy algorithm. Then,

$$\frac{c(I^*)}{c(I)} \leq q := \max_{F \subseteq E: r(F) > 0} \frac{r(F)}{\rho(F)},$$

and there exists $c: E \rightarrow \{0, 1\}$ such that equality holds.

Proof. For all $j \in \{1, \dots, |E|\}$, Observation 7.3 implies that $|I \cap E_j| \geq \rho(E_j)$ and $I^* \cap E_j \in \mathcal{I}$ implies $|I^* \cap E_j| \leq r(E_j)$. Hence, we obtain (even if $r(E_j) = 0$)

$$|I \cap E_j| \geq (1/q) \cdot |I^* \cap E_j|. \quad (7.1)$$

Let $\{e_j\} := E_j \setminus E_{j-1}$ for $j \in \{1, \dots, |E|\}$ and $c(e_{|E|+1}) \equiv 0$. By definition of the greedy algorithm and since $c(e_{|E|}) \geq 0$, we have $c(e_j) \geq c(e_{j+1})$ for all $j \in \{1, \dots, |E|\}$. This allows to apply (7.1) to obtain

$$\begin{aligned} c(I) &= \sum_{j=1}^{|E|} (|I \cap E_j| - |I \cap E_{j-1}|) c(e_j) \\ &= \sum_{j=1}^{|E|} |I \cap E_j| (c(e_j) - c(e_{j+1})) \\ &\stackrel{(7.1)}{\geq} (1/q) \cdot \sum_{j=1}^{|E|} |I^* \cap E_j| (c(e_j) - c(e_{j+1})) \\ &= (1/q) \cdot \sum_{j=1}^{|E|} (|I^* \cap E_j| - |I^* \cap E_{j-1}|) c(e_j) \\ &= (1/q) \cdot c(I^*), \end{aligned}$$

as claimed.

For the second part of the statement, observe that $c(I^*) > 0$ implies that $F \subseteq E$ exists with $q = \frac{r(F)}{\rho(F)}$. Set $c(e) = 1$ for $e \in F$ and $c(e) = 0$ for $e \notin F$. Since the greedy algorithm does not specify how to break ties in $\arg \max$, we may assume that $E_{\rho(F)}$ is a basis of F . But because $c(e) = 0$ for all $e \notin F$, we have $c(I) = \rho(F) = (1/q) \cdot r(F) = (1/q) \cdot c(I^*)$. \square

With this we obtain that matroids capture exactly the independence systems for which the greedy algorithm is optimal.

Corollary 7.10. The following are equivalent for an independence system (E, \mathcal{I}) :

- (a) (E, \mathcal{I}) is a matroid.
- (b) For all $c : E \rightarrow \mathbb{R}_{\geq 0}$, the greedy algorithm yields an optimum solution to $\max\{c(I) : I \in \mathcal{I}\}$.
- (c) For all $c : E \rightarrow \{0, 1\}$, the greedy algorithm yields an optimum solution to $\max\{c(I) : I \in \mathcal{I}\}$.
- (d) For all $c : E \rightarrow \mathbb{R}_{\geq 0}$, the greedy algorithm (for minimization) yields an optimum solution to $\min\{c(B) : B \in \mathcal{B}\}$.

Proof. The equivalence of the first three points is immediate from Theorem 7.9.

For the last point, observe that $\max\{\sum_{e \in I} c'(e) : I \in \mathcal{I}\}$ with $c'(e) := -c(e) + \sum_{e \in E} c(e) \geq 0$ has optimum solutions that are a bases of E (this may not be true for all optima, since $c(e) = 0$ is permitted). If (E, \mathcal{I}) is a matroid, all bases have the same cardinality, hence $c(B)$ and $c'(B)$ differ by the same constant for all bases B of E . This means that a basis is optimal for the maximization problem if and only if it is optimal for

$$\max\{-c(B) : I \in \mathcal{I}\} = -\min\{c(B) : B \in \mathcal{B}\}.$$

Finally, observe that the minimization version of the greedy algorithm for c produces the same solution as the maximization version for c' and vice-versa. Hence, the greedy algorithm is optimal for minimization for all $c : E \rightarrow \mathbb{R}_{\geq 0}$ if and only if it is optimal for maximization for all $c' : E \rightarrow \mathbb{R}_{\geq 0}$. \square

It is remarkable that we obtain a full characterization for the optimality of the greedy algorithm. In particular, Corollary 7.10 yields an alternate proof for the correctness of Kruskal's algorithm for the minimum spanning tree problem (see *Algorithmic Discrete Mathematics*). Note that the bound of Theorem 7.9 only carries over to minimization problems on matroids. The TSP problem has bounded rank quotient, but the greedy algorithm can produce arbitrarily bad solutions (*exercise*).

We further give a sufficient condition for an independence system to have bounded rank quotient, namely, when it arises as an intersection of matroids.

Theorem 7.11. Let $(E, \mathcal{I}_1), \dots, (E, \mathcal{I}_k)$ be matroids. Then, $(E, \mathcal{I}) := (E, \bigcap_{i=1}^k \mathcal{I}_i)$ is an independence system with rank quotient $q = \max_{F \subseteq E: r(F) > 0} r(F)/\rho(F) \leq k$.

Proof. The fact that $(E, \mathcal{I}_1), \dots, (E, \mathcal{I}_k)$ are independence systems immediately implies that (E, \mathcal{I}) is an independence system. Let $F \subseteq E$ and let $B, B' \subseteq F$ be any two bases of F with respect to (E, \mathcal{I}) . It is sufficient to show that $|B'| \leq k \cdot |B|$.

Since $B \in \mathcal{I} = \bigcap_{i=1}^k \mathcal{I}_i$, we can choose bases \bar{B}_i with $B \subseteq \bar{B}_i \subseteq B \cup B'$ with respect to (E, \mathcal{I}_i) for $i \in \{1, \dots, k\}$. We claim that $\bigcap_{i=1}^k \bar{B}_i = B$. To see this, suppose there exists $e \in \bigcap_{i=1}^k \bar{B}_i \setminus B$. Then, $B \cup \{e\} \subseteq \bigcap_{i=1}^k \bar{B}_i \in \mathcal{I}_j$ for all $j \in \{1, \dots, k\}$ and thus $B \cup \{e\} \in \mathcal{I}$. Because of $B \cup \{e\} \in B \cup B' \subseteq F$, this contradicts B being a basis of F .

Now, our claim implies that every $e \in B' \setminus B$ is contained in at most $k - 1$ of the sets $\bar{B}_i \setminus B$, which means that

$$\left(\sum_{i=1}^k |\bar{B}_i| \right) - k|B| = \sum_{i=1}^k |\bar{B}_i \setminus B| \leq (k-1)|B' \setminus B| \leq (k-1)|B'|. \quad (7.2)$$

As before, we can choose additional bases \bar{B}'_i with $B' \subseteq \bar{B}'_i \subseteq B \cup B'$ with respect to (E, \mathcal{I}_i) for $i \in \{1, \dots, k\}$. Since (E, \mathcal{I}_i) is a matroid, we have

$$|\bar{B}_i| = |\bar{B}'_i| \quad \forall i \in \{1, \dots, k\}. \quad (7.3)$$

Overall, we obtain

$$|B'| \leq |B'| + \sum_{i=1}^k |\bar{B}'_i \setminus B'| = \left(\sum_{i=1}^k |\bar{B}'_i| \right) - (k-1)|B'| \stackrel{(7.3)}{=} \left(\sum_{i=1}^k |\bar{B}_i| \right) - (k-1)|B'| \stackrel{(7.2)}{\leq} k|B|. \quad \square$$

For most practical problems, the greedy algorithm can produce arbitrarily bad solutions. Nevertheless, due to its simplicity and efficiency, it is a popular choice in practice.

7.2 Local search

Once we found some feasible solution, we can try to improve it by modifications. In other words, we search a neighborhood for better solutions, update our solution, and iterate.

Algorithm: local search

input: feasible solutions \mathcal{F} , neighborhoods $N: \mathcal{F} \rightarrow \mathcal{F}$
objective $c: \mathcal{F} \rightarrow \mathbb{R}$, initial solution $S \in \mathcal{F}$

output: solution $S' \in \mathcal{F}$

while $\exists S' \in N(S)$ with $c(S') \geq c(S)$:

$S \leftarrow S'$

return S

The quality of this algorithm critically depends on the definition of the neighborhood N .

Example 7.12.

- 2-opt heuristic for TSP: Here \mathcal{F} is the set of all tours and the neighborhood of a tour T consists of the tours T' obtained from T via a pairwise exchange of two edges $\{u_1, v_1\}, \{u_2, v_2\}$ by the edges $\{u_1, v_2\}, \{u_2, v_1\}$.
- 1-exchange/2-exchange for the equipartition problem: We are given an undirected, edge-weighted graph $G = (V, E)$ with $|V|$ even and are looking for $S \subseteq V$ with $|S| = |V \setminus S| = |V|/2$, while maximizing the weights of the cut $\delta(S)$. The neighborhood of a solution $S \subseteq V$ could then, for example, be all $S' = S \setminus \{v\} \cup \{v'\}$ for $v \in S, v' \notin S$ (2-exchange). Alternatively, we can add a penalty term of the form $\alpha(|S| - |V \setminus S|)^2$ to the objective and allow 1-exchanges $S' = S \setminus \{v\}, v \in S$ or $S' = S \cup \{v'\}, v' \notin S$. \triangle

A general problem of local search is that, inherently, it gets stuck in local extrema. We now discuss modifications of local search that alleviate this flaw to some extent.

7.2.1 Tabu search

The idea of *tabu search* is to escape local minima by forcing, in every step, to switch to a solution in the neighborhood that we didn't already visit – even if that means getting worse in the objective. However, it is generally computationally too expensive to track all previously encountered solutions, and we instead limit ourselves to a *tabu list* of limited length that only contains recently encountered solutions.

Algorithm: tabu search

input: feasible solutions \mathcal{F} , neighborhoods $N: \mathcal{F} \rightarrow \mathcal{F}$
objective $c: \mathcal{F} \rightarrow \mathbb{R}$, initial solution $S^{(0)} \in \mathcal{F}$

output: solution $S' \in \mathcal{F}$

$T \leftarrow \emptyset$ (tabu list)

for $i \leftarrow 1, 2, \dots$:

$S^{(i)} \leftarrow \arg \max \text{ or } \min \{c(S) : S \in N(S^{(i-1)}) \setminus T\}$

$T \leftarrow (T_2, T_3, \dots) \oplus (S^{(i)})$ (discard oldest)

return $\arg \max \text{ or } \min \{S^{(0)}, \dots, S^{(i)}\}$

Example 7.13. Consider the problem of, for given $k \in \mathbb{N}$, finding a coloring $\phi: V \rightarrow \{1, 2, \dots, k\}$ of the vertices of an undirected graph $G = (V, E)$, such that $\phi(u) \neq \phi(v)$ for all edges $\{u, v\} \in E$. We minimize

$$c(\phi) := |\{\{u, v\} \in E : \phi(u) = \phi(v)\}|,$$

and define the neighborhood of ϕ to consist of the colorings that differ from ϕ in the color of a single vertex, i.e.,

$$N(\phi) := \{\varphi: V \rightarrow \{1, \dots, k\} : |\{v \in V : \phi(v) \neq \varphi(v)\}| = 1\}. \quad \triangle$$

The parameters of tabu search are the length of the tabu list and the stopping condition. A longer tabu list or execution potentially improves the result, but comes at the cost of deteriorating the running time. Natural stopping conditions are fixing the number of iterations without improvement or in total. The best choices of the associated parameters are highly problem-dependent and usually determined by experimental tuning.

7.2.2 Simulated annealing

Simulated annealing also allows to deteriorate the objective, but does so probabilistically instead of using a deterministic tabu list. More specifically, we allow worse solutions with a certain probability that depends on how much worse they are and decreases over the course of the execution of the algorithm. The method is inspired by the physics of crystal formation, which requires a gradual decrease in temperature to allow time for the individual particles to find their place in the crystalline structure. Accordingly, it maintains a temperature and uses a probability according to the Boltzmann distribution.

Algorithm: simulated annealing (for maximization)

input: feasible solutions \mathcal{F} , neighborhoods $N: \mathcal{F} \rightarrow \mathcal{F}$
objective $c: \mathcal{F} \rightarrow \mathbb{R}$, initial solution $S^{(0)} \in \mathcal{F}$
initial temperature $T \geq 1$, cooling rate θ

output: solution $S' \in \mathcal{F}$

for $i \leftarrow 1, 2, \dots$:

 randomly pick $S^{(i)} \in N(S^{(i-1)})$ and $\alpha \in [0, 1]$

if $c(S^{(i)}) - c(S^{(i-1)}) < 0$ and $e^{(c(S^{(i)}) - c(S^{(i-1)}))/T} < \alpha$:

$S^{(i)} \leftarrow S^{(i-1)}$ (discard solution)

$T \leftarrow T \cdot \theta$

return $\arg \max \{c(S) : S \in \{S^{(0)}, \dots, S^{(i)}\}\}$

Remark 7.14. It can be shown that there always exists parameters for which, with high probability, simulated annealing eventually finds an optimum solution (see [30]).

As for the previous heuristics, a lot of experimentation is required for tuning the parameters of simulated annealing, such as the starting temperature, the cooling rate, and the stopping criterion.

7.2.3 Genetic algorithms

While simulated annealing is physics inspired, *genetic algorithms* take their inspiration from biology. The basic idea is to simultaneously consider a set (called *population*) of candidate solutions and to improve that set by combining individual solutions. In each iteration (called *generation*) the method selects a series pairs of solutions (called *parents*) which are *crossed* to generate new solutions (called *offspring*). In addition, the resulting solutions are randomly perturbed (*mutated*) with a small probability to avoid getting stuck in local optima. A subset of the generated solutions is then selected, favoring solutions with better objective values (their *fitness*).

Algorithm: genetic algorithm

input: feasible solutions \mathcal{F} , objective $c: \mathcal{F} \rightarrow \mathbb{R}$
initial population $\mathcal{S}^{(0)} \subset \mathcal{F}$

output: solution $S \in \mathcal{F}$

for generations $i \leftarrow 1, 2, \dots$:

for parents $\{S^\circ, S^\bullet\} \subseteq \mathcal{S}^{(i-1)}$:

$S^\circ \leftarrow S^\circ \times S^\bullet$ (reproduction)

$\tilde{S}^\bullet \leftarrow \begin{cases} \text{random perturbation of } S^\bullet, & \text{with low probability,} \\ S^\bullet & \text{otherwise.} \end{cases}$ (mutation)

$\mathcal{S}^{(i)} \leftarrow \mathcal{S}^{(i)} \cup \{S^\circ\}$

while $|\mathcal{S}^{(i)}|$ too large :

choose $S \in \mathcal{S}^{(i)}$ with probability depending on $c(S)$ (selection)

$\mathcal{S}^{(i)} \leftarrow \mathcal{S}^{(i)} \setminus \{S\}$

return $\arg \max [\text{or } \min] \{c(S) : S \in \mathcal{S}^{(0)} \cup \dots \cup \mathcal{S}^{(i)}\}$

As with the previous methods, the parameters such as choice of parents, crossing method, mutation rate and process, selection criterion, stopping criterion, etc., are heavily problem dependent and require careful tuning.

8 Approximation Algorithms

The solution methods for MIPs covered in Chapters 3, 5, and 6 generally require exponential running times, and we cannot hope for more efficient methods due to the NP-hardness of solving integer programs (Theorem 3.1). In practice, exponential runtimes are computationally infeasible except for very small problem instances, and we can only hope that input instances are “friendly” in the sense that our methods terminate much faster than their theoretical worst-case running times suggest. Otherwise, we have to abort the solution process prematurely. In general, we cannot make any theoretical guarantees regarding the quality of a solution obtained this way. The same applies to the heuristic solutions covered in Chapter 7.

In this chapter, we focus on *approximation algorithms* that provide (theoretically) efficient running times as well as a provable solution quality.

Definition 8.1. For every instance I of an optimization problem, let $\text{OPT}(I)$ denote the optimum solution value and $\text{ALG}(I)$ the objective function value of the solution computed by an algorithm ALG . Then, ALG is a γ -approximation algorithm if it has polynomial running time and its solution is always within a factor of $\gamma \geq 1$ of the optimum solution, i.e., for all instances I ,

$$\text{ALG}(I) \leq \gamma \text{OPT}(I) \text{ for minimization problems, and } \text{ALG}(I) \geq (1/\gamma) \text{OPT}(I) \text{ for maximization problems.}$$

In this sense, Theorem 7.9 gives us a solution guarantee for the greedy algorithm.

Example 8.2. The greedy algorithm is a q -approximation algorithm for independence systems (E, \mathcal{I}) of rank quotient $q = \max_{F \subseteq E} r(F)/\rho(F)$. \triangle

In particular, we obtain the following approximation algorithm.

Theorem 8.3. The greedy algorithm is a 2-approximation algorithm for finding a bipartite matching of maximum weight.

Proof. Let $G = (V_1 \cup V_2, E)$ be a bipartite graph. By Theorems 7.9 and 7.11, it suffices to show the existence of two matroids (E, \mathcal{I}_1) and (E, \mathcal{I}_2) such that $\mathcal{I} = \mathcal{I}_1 \cap \mathcal{I}_2$ is the set of matchings in G . Two sets with this property are defined by, for $\ell \in \{1, 2\}$,

$$\mathcal{I}_\ell := \{M \subseteq E : |M \cap \delta_G(v)| \leq 1 \forall v \in V_\ell\}.$$

Clearly, $\emptyset \in \mathcal{I}_\ell$ holds, \mathcal{I}_ℓ is closed under taking subsets, and $\mathcal{I} = \mathcal{I}_1 \cap \mathcal{I}_2$. Also, for all $F \subseteq E$, the bases of F in (E, \mathcal{I}_ℓ) have the same cardinality $|\{v \in V_\ell : F \cap \delta_G(v) \neq \emptyset\}|$. Hence, (E, \mathcal{I}_ℓ) is a matroid. \square

8.1 Approximation for TSP

To prove lower bounds on the best possible approximation factor of an NP-hard problem, we obviously have to assume $P \neq NP$ since otherwise even optimal algorithms with polynomial running time would be possible.¹ Under this assumption we can show that TSP is not approximable in general.

Theorem 8.4. There is no γ -approximation algorithm for TSP for any $\gamma = \gamma(n) \geq 1$, provided that $P \neq NP$.

Proof. We consider the Hamiltonian cycle problem in which an undirected graph $G = (V, E)$ is given and we have to decide whether G contains a Hamiltonian cycle, i.e., a cycle that visits every vertex exactly once. This problem is one of the classical NP-complete problems (see *Algorithmic Discrete Mathematics*).

We construct a TSP instance on a complete graph $G' = (V, E')$ with edge weights $c: E' \rightarrow \mathbb{R}_{\geq 0}$, such that any γ -approximate solution of the TSP instance is sufficient to decide whether a Hamiltonian cycle exists in G . Since approximation algorithms have polynomial running time by definition, and since G' and c can be generated in polynomial time, this construction immediately implies the statement of the theorem.

To do this, we set the edge weights for $e \in E'$ to

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

Obviously, a Hamiltonian cycle in G corresponds to a tour of weight 0 in G' and vice-versa. Moreover, every tour in G' that has edges from $E' \setminus E$ has weight at least 1. By definition, every γ -approximation algorithm with $\gamma \geq 1$ must produce a solution of weight 0 if one exists, and thus decides whether a Hamiltonian cycle exists. \square

Corollary 8.5. There is no γ -approximation algorithm for TSP for any constant $\gamma \geq 1$, provided that $P \neq NP$, even if edge weights must be strictly positive.

Proof. We proceed as in the proof of Theorem 8.4, but set

$$c(e) := \begin{cases} \frac{1}{2^{\gamma n}}, & \text{if } e \in E, \\ 1, & \text{otherwise.} \end{cases}$$

Then, a Hamiltonian cycle in G corresponds to a tour of weight $\frac{1}{2^{\gamma}} < 1$ in G' and, conversely, every tour of weight less than 1 must be Hamiltonian, since it cannot afford even a single edge of $E' \setminus E$. Deciding whether a Hamiltonian cycle exists in G thus amounts to deciding whether the lightest tour in G' has weight $\frac{1}{2^{\gamma}}$ or at least 1. By definition, every γ -approximation algorithm for TSP can be used to distinguish these cases. \square

We have thus shown that we have to accept very bad solutions if we only want to invest polynomial computing time to solve TSP. This result strongly relies on the fact that the edge weights can be arbitrary. Fortunately, TSP instances often have additional structure: Most of the time the edge weights correspond to distances (e.g., in a road network) and are therefore metric. In this case, we speak of *metric TSP*. In particular, the edge weights (or lengths) fulfil the triangle inequality $c(\{u, v\}) \leq c(\{u, w\}) + c(\{w, v\})$ for all $u, v, w \in V$. Note that metric TSP is still NP-hard (see for example [32, Theorem 21.2]).

¹This statement assumes that the corresponding decision problem is in NP.

8.1.1 Greedy algorithm

We first consider the greedy algorithm that proceeds as follows. We start with a cycle on two vertices v_1, v_2 , where $\{v_1, v_2\}$ is a cheapest edge of the graph (the cycle uses this edge twice). We extend the cycle step by step by choosing a cheapest edge e that connects a vertex v of our cycle with a vertex v' outside. To obtain a cycle again, we replace an edge $\{v, v''\}$ from our cycle with $\{v', v''\}$.

Algorithm: greedy algorithm for metric TSP

input: complete graph $G = (V, E)$, metric edge weights $c: E \rightarrow \mathbb{R}_{\geq 0}$

output: tour T

$\{v_1, v_2\} \leftarrow \arg \min_{e \in E} c(e)$

$T = (V_T, E_T) \leftarrow \{(v_1, v_2), (v_2, v_1)\}$

while $V_T \neq V$:

$(v, v') \leftarrow \arg \min\{c(\{u, u'\}) : (u, u') \in V_T \times (V \setminus V_T)\}$ (edge selection)

take $v'' \in V_T$ with $(v, v'') \in T$

$T \leftarrow T \setminus \{(v, v'')\} \cup \{(v, v'), (v', v'')\}$

return T

Theorem 8.6. The greedy algorithm is a 2-approximation algorithm for metric TSP.

Proof. A closer look reveals that the edge selection $\{v, v'\}$ in each iteration corresponds exactly to the edge selection of Prim's algorithm for finding a minimum spanning tree (see *Algorithmic Discrete Mathematics*). If F is the set of all edges selected during the course of the algorithm, then F is therefore a minimum spanning tree.

In each step the length of our cycle increases by

$$c(\{v, v'\}) + c(\{v', v''\}) - c(\{v, v''\}) \leq 2c(\{v, v'\}),$$

where we used the triangle inequality. The greedy algorithm thus ensures that $c(T) \leq 2c(F)$.

Now consider an optimum TSP tour T^* . If we remove any edge from T^* , we obtain a spanning tree F' . It follows that

$$c(T) \leq 2c(F) \leq 2c(F') \leq 2c(T^*). \quad \square$$

8.1.2 Algorithm of Christofides

Starting from a minimum spanning tree F^* , we can also construct a 2-approximation more directly: We double every edge in F^* and obtain a *Eulerian (multi)graph*, i.e., a graph that has a cycle using every edge exactly once. It is easy to see (cf. *Algorithmic Discrete Mathematics*) that a graph is Eulerian if and only if it is connected and every vertex has even degree (start somewhere, move along unused edges until a cycle is closed; repeat and join obtained cycles). A Eulerian cycle in the doubled spanning tree is of course not yet a tour, since it uses vertices multiple times. However, the triangle inequality allows to eliminate redundant visits to vertices without increasing the length of the tour, i.e., we skip all but the first occurrence of each vertex.

To summarise, we started from a minimum spanning tree and extended it to a Eulerian graph while doubling its weight. We can improve this algorithm by finding the cheapest extension to a Eulerian graph. To this

end, we first observe that every graph has an even number of vertices of odd degrees, since the sum of all vertex degrees must be even (namely $2|E|$). We want to partition these vertices into pairs and join each pair via the corresponding direct edge to obtain an Eulerian graph. The problem of finding such a partition resulting in the minimum total weight of the corresponding edges is precisely the minimum weight perfect matching problem, which can be solved using the algorithm of Edmonds [15] in polynomial time (cf. lecture *Combinatorial Optimization*). The resulting algorithm is known as *Christofides' algorithm* [8].

Algorithm: Christofides' algorithm for metric TSP

input: complete graph $G = (V, E)$, metric edge weights $c: E \rightarrow \mathbb{R}_{\geq 0}$

output: tour T

$F^* \leftarrow$ minimum spanning tree in G

$V' \leftarrow \{v \in V : v \text{ has odd degree in } F^*\}$

$M^* \leftarrow$ perfect matching of minimum weight in $G[V']$

$R \leftarrow$ Eulerian cycle in $(V, F^* \cup M^*)$

$T \leftarrow$ eliminate redundant vertices in R (replace $\{u, v\}, \{v, w\}$ by $\{u, w\}$)

return T

Theorem 8.7 ([8]). Christofides' algorithm is a $3/2$ -approximation algorithm for metric TSP.

Proof. As before, we have $c(F^*) \leq c(T^*)$, where F^* denotes a minimum spanning tree and T^* an optimum TSP tour. Let V' be the set of all vertices of odd degrees in F^* , and let $M^* \subset E$ be a perfect matching of minimum weight in $G[V']$. Since, by the triangle inequality, the removal of redundant vertices does not increase the length of the tour, it suffices to show that $c(M^*) \leq c(T^*)/2$.

First, let T' be an optimum TSP tour in $G[V']$. We claim that $c(T') \leq c(T^*)$. To see this, we start with T^* and gradually remove vertices that are not in V' by replacing consecutive edges with the direct connection. By the triangle inequality the cycle does not become longer, which establishes the claim.

Now, since $|T'|$ is even, T' can be divided into two perfect matchings in $G[V']$ which each have weight at least $c(M^*)$, because M^* is optimal. It follows that $c(M^*) \leq c(T')/2 \leq c(T^*)/2$. \square

The state of the art in approximation guarantees for metric TSP is summarized below.

Theorem 8.8 ([28]). There exists a randomized algorithm for metric TSP that outputs a tour whose expected weight is at most $(\frac{3}{2} - \varepsilon) \cdot c(T^*)$, for some $\varepsilon > 10^{-36}$.

Theorem 8.9 ([29]). The metric TSP problem does not admit a γ -approximation algorithm with $\gamma < 123/122$, provided that $P \neq NP$.

Closing the gap between the best upper and lower bounds is an important open problem in the field of approximation algorithms.

8.2 Polynomial-time approximation schemes

An approximation algorithm guarantees solutions of fixed approximation quality. We can view the best-possible approximation factor as a measure of the difficulty of a (NP-hard) optimization problem. In that sense, the easiest problems allow *every* fixed approximation factor $\gamma > 1$.

Definition 8.10. A *polynomial-time approximation scheme (PTAS)* is a family $\{A_\varepsilon\}$ of algorithms, one for every $\varepsilon > 0$, such that A_ε is a $(1 + \varepsilon)$ -approximation algorithm.

A PTAS therefore yields, for every fixed $\varepsilon > 0$, an algorithm with polynomial running time. However, our definition allows an arbitrary dependence of the running time on ε , in particular the running time could increase exponentially with increasing $1/\varepsilon$. The following definition is more restrictive.

Definition 8.11. A *fully polynomial-time approximation scheme (FPTAS)* is a polynomial-time approximation scheme $\{A_\varepsilon\}$ for which the running time of A_ε can be bounded by a polynomial in $1/\varepsilon$ (as well as in the instance size).

We first show that we cannot hope for a better dependence on ε for NP-hard problems.

Theorem 8.12. If an optimization problem with integral, non-negative objective function has a PTAS $\{A_\varepsilon\}$, where A_ε has polynomial running time in $\log(1/\varepsilon)$, the problem also has a polynomial-time exact algorithm.

Proof. We show the statement for minimization problems (the proof for maximization problems is analogous). We construct a polynomial-time algorithm as follows. We first apply A_ε with $\varepsilon > 0$ arbitrarily (e.g., $\varepsilon = 1$) to the instance and obtain a solution of value $c_\varepsilon \geq c_{\text{OPT}}$. Since A_ε has polynomial running time in the input size, c_ε is at most simply exponential in the input size.

Set $\varepsilon' := \frac{1}{c_\varepsilon + 1}$. Since c_ε is at most simply exponential in the input size, $\log(1/\varepsilon') \in \mathcal{O}(\log c_\varepsilon)$ is polynomial in the the input size. We apply $A_{\varepsilon'}$ and return the computed solution of value $c_{\varepsilon'}$. The outlined algorithm has polynomial running time since the running time of $A_{\varepsilon'}$ is polynomial in $\log(1/\varepsilon')$.

It remains to show that $c_{\varepsilon'} = c_{\text{OPT}}$. Because c_{OPT} and $c_{\varepsilon'}$ are integers by assumption, this follows from

$$c_{\text{OPT}} \leq c_{\varepsilon'} \leq (1 + \varepsilon') c_{\text{OPT}} = \left(1 + \frac{1}{c_\varepsilon + 1}\right) \cdot c_{\text{OPT}} < \left(1 + \frac{1}{c_{\text{OPT}}}\right) \cdot c_{\text{OPT}} \leq c_{\text{OPT}} + 1. \quad \square$$

8.2.1 Example: knapsack problem

We now show that the knapsack problem can be approximated arbitrarily well, and that even an FPTAS exists. In addition to the pseudopolynomial algorithm from Section 5.2.1, this is an indication that the knapsack problem is an “easy” NP-hard problem.

We consider a knapsack problem of the form $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{a}^\top \mathbf{x} \leq \beta, \mathbf{x} \in \{0, 1\}^n\}$ with $\mathbf{a}, \mathbf{c} \in \mathbb{R}_{\geq 0}^n$, $\beta \in \mathbb{R}_{\geq 0}$, and $0 < a_i \leq \beta$ for all $i \in \{1, \dots, n\}$. The idea of the FPTAS is simple: We round the entries of \mathbf{c} to integer multiples of some $\mu \in \mathbb{R}$, and apply the dynamic program of Section 5.2.1 to the rounded instance. We choose μ in such a way that (a) our algorithm has a polynomial running time and (b) our solution differs by at most a factor of $(1 + \varepsilon)$ from c_{OPT} .

Algorithm: FPTAS for the knapsack problem.

input: $\mathbf{a}, \mathbf{c} \in \mathbb{R}_{\geq 0}^n; \beta \in \mathbb{R}_{\geq 0}, \varepsilon > 0$

output: approximate solution to $\max\{\mathbf{c}^\top \mathbf{x} : \mathbf{a}^\top \mathbf{x} \leq \beta, \mathbf{x} \in \{0, 1\}^n\}$

$M \leftarrow \max_{i \in \{1, \dots, n\}} c_i, \mu \leftarrow \varepsilon M / n$

$\mathbf{c}' \leftarrow \lfloor \mathbf{c} / \mu \rfloor$ (component-wise)

$\mathbf{x} \leftarrow \arg \max\{(\mathbf{c}')^\top \mathbf{x} : \mathbf{a}^\top \mathbf{x} \leq \beta, \mathbf{x} \in \{0, 1\}^n\}$ (via DP of Section 5.2.1)

return \mathbf{x}

Theorem 8.13. There exists an FPTAS for the knapsack problem with running time $\mathcal{O}(n^3/\varepsilon)$.

Proof. We first consider the rounding error. Rounding the entries of \mathbf{c} to the nearest integer multiples of μ influences the solution value by at most μn , so we need to ensure that $\mu n \leq \varepsilon c_{\text{OPT}}$. Because of $M \leq c_{\text{OPT}}$, it is sufficient to set $\mu n = \varepsilon M$, i.e. $\mu := \varepsilon M/n$.

The running time of the algorithm is determined by the running time of the dynamic program on the rounded instance. According to Theorem 5.16, the latter is $\mathcal{O}(nZ)$, where $Z := \sum_{i=1}^n \lfloor c_i/\mu \rfloor \leq n \lfloor M/\mu \rfloor \in \mathcal{O}(n^2/\varepsilon)$. Hence, the running time is polynomial in n and $1/\varepsilon$, as desired. \square

8.3 LP Rounding

Many approximation algorithms are based on rounding a solution to some LP relaxation. The applicability of this approach depends on whether we can ensure the rounded solution to be feasible and whether we are able to estimate the quality of the rounded solution. Randomization often helps, especially for the latter aspect (cf. the proof of Theorem 4.15). We illustrate the general approach on an example from machine scheduling. Scheduling problems are generally concerned with assigning jobs $j \in \{1, 2, \dots, n\}$ to machines $i \in \{1, \dots, m\}$, where assigning job j on machine i incurs a processing time of $p_{ij} \geq 0$. Many practical problems can be expressed in this way, and, accordingly, there is a large number of variants of this problem, which differ in the objective function, the processing times and other restrictions for valid solutions.

We consider a variant in which every job $j \in \{1, \dots, n\}$ may be rejected (i.e., not processed at all) for a cost s_j . Machines can only process one job at a time, but jobs can be interrupted (*preemption*), transferred to other machines while being processed (*migration*), and even run simultaneously on several machines (*parallelization*). The objective is to reduce the sum of the completion time of the last job (the *makespan*) and the rejection costs of the unprocessed jobs.

We can formulate this problem as a MIP with the following variables

x_{ij} = proportion of job j processed by machine i

$$y_j = \begin{cases} 1, & \text{if job } j \text{ is processed,} \\ 0, & \text{otherwise} \end{cases}$$

T = makespan.

The resulting MIP is

$$\begin{aligned} \min \quad & T + \sum_{j=1}^n s_j (1 - y_j), \\ \text{s.t.} \quad & \sum_{j=1}^n p_{ij} x_{ij} \leq T \quad \forall i \in \{1, \dots, m\}, \\ & \sum_{i=1}^m x_{ij} = y_j \quad \forall j \in \{1, \dots, n\}, \\ & x_{ij} \geq 0 \quad \forall i \in \{1, \dots, m\}, j \in \{1, \dots, n\}, \\ & \mathbf{y} \in \{0, 1\}^n. \end{aligned}$$

It can be shown that the above scheduling problem is NP-hard, see [26]. We describe an algorithm that determines a feasible solution by rounding the LP relaxation.

Algorithm: DETERMINISTICROUNDING(α)

input: threshold $\alpha \in (0, 1)$

determine solution $(\mathbf{x}^*, \mathbf{y}^*, T^*)$ of the LP relaxation

for $j \in \{1, \dots, n\}$:

if $y_j^* \leq \alpha$:

$\hat{y}_j \leftarrow 0$

$\hat{x}_{ij} \leftarrow 0 \quad \forall i \in \{1, \dots, m\}$

else

$\hat{y}_j \leftarrow 1$

$\hat{x}_{ij} \leftarrow x_{ij}^*/y_j^* \quad \forall i \in \{1, \dots, m\}$

$(y_j^* > \alpha)$

$\hat{T} \leftarrow \max_{i \in \{1, \dots, m\}} \sum_{j=1}^n p_{ij} \hat{x}_{ij}$

return $(\hat{\mathbf{x}}, \hat{\mathbf{y}}, \hat{T})$

We establish an approximation guarantee for this rounding method.

Theorem 8.14. DETERMINISTICROUNDING($1/2$) is a 2-approximation algorithm.

Proof. The computed solution $(\hat{\mathbf{x}}, \hat{\mathbf{y}}, \hat{T})$ is feasible since either $\hat{y}_j = 0 = \sum_{i=1}^m \hat{x}_{ij}$ or, by feasibility of $(\mathbf{x}^*, \mathbf{y}^*, T^*)$, $\hat{y}_j = 1 = \sum_{i=1}^m x_{ij}^*/y_j^* = \sum_{i=1}^m \hat{x}_{ij}$.

We claim that, for all $i \in \{1, \dots, m\}$, $j \in \{1, \dots, n\}$, it holds that

$$(1 - \hat{y}_j) \leq \frac{1}{1 - \alpha} (1 - y_j^*), \text{ and} \quad (8.1)$$

$$\hat{x}_{ij} \leq \frac{1}{\alpha} x_{ij}^*. \quad (8.2)$$

For $y_j^* \leq \alpha$, we have $\hat{y}_j = \hat{x}_{ij} = 0$ and thus

$$\frac{1}{1 - \alpha} (1 - y_j^*) \geq \frac{1}{1 - \alpha} (1 - \alpha) = 1 = (1 - \hat{y}_j),$$

and

$$\frac{1}{\alpha} x_{ij}^* \geq 0 = \hat{x}_{ij}.$$

For $y_j^* > \alpha$, we have $\hat{y}_j = 1$ and thus

$$\frac{1}{1 - \alpha} (1 - y_j^*) \geq 0 = (1 - \hat{y}_j),$$

and

$$\hat{x}_{ij} = \frac{x_{ij}^*}{y_j^*} < \frac{1}{\alpha} x_{ij}^*.$$

Hence the claim holds.

Now, let $i \in \{1, \dots, m\}$ be a machine with completion time \hat{T} in the solution $(\hat{\mathbf{x}}, \hat{\mathbf{y}}, \hat{T})$, i.e., such that $\hat{T} = \sum_{j=1}^n p_{ij} \hat{x}_{ij}$. Then,

$$\frac{1}{\alpha} T^* \geq \frac{1}{\alpha} \sum_{j=1}^n p_{ij} x_{ij}^* \stackrel{(8.2)}{\geq} \sum_{j=1}^n p_{ij} \hat{x}_{ij} = \hat{T}. \quad (8.3)$$

It follows that

$$\begin{aligned} \hat{T} + \sum_{j=1}^n s_j (1 - \hat{y}_j) &\stackrel{(8.1),(8.3)}{\leq} \frac{1}{\alpha} T^* + \frac{1}{1-\alpha} \sum_{j=1}^n s_j (1 - y_j^*) \\ &\leq \max\left\{\frac{1}{\alpha}, \frac{1}{1-\alpha}\right\} \left(T^* + \sum_{j=1}^n s_j (1 - y_j^*)\right) \leq \max\left\{\frac{1}{\alpha}, \frac{1}{1-\alpha}\right\} \cdot c_{LP} \leq \max\left\{\frac{1}{\alpha}, \frac{1}{1-\alpha}\right\} \cdot c_{MIP}. \end{aligned}$$

Setting $\alpha = \frac{1}{2}$ completes the proof. \square

By randomizing the selection of α we can further improve the (expected) approximation factor. We can derandomize the resulting algorithm by testing a linear number of values for α .

Algorithm: DERANDOMIZEDROUNDING

$S \leftarrow \{\text{DETERMINISTICROUNDING}(\alpha) : \alpha \in ((\frac{1}{e}, 1) \cap \{y_1^*, \dots, y_n^*\}) \cup \{\frac{1}{e}\}\}$
return $\arg \min\{\hat{T} + \sum_{i=1}^m s_j (1 - \hat{y}_j) : (\hat{x}, \hat{y}, \hat{T}) \in S\}$

We obtain the following solution guarantee.

Theorem 8.15 ([26]). The derandomized rounding algorithm is a $\frac{e}{e-1}$ -approximation algorithm, where $\frac{e}{e-1} \approx 1.58$.

Proof. First, consider the randomized procedure where $\alpha \sim \mathcal{U}(\frac{1}{e}, 1)$ is selected uniformly at random from the open interval $(\frac{1}{e}, 1)$. Let $(\hat{x}(\alpha), \hat{y}(\alpha), \hat{T}(\alpha))$ denote the resulting solutions. Then, using the arguments from the proof of Theorem 8.14, we have

$$\mathbb{E}_\alpha[\hat{T}(\alpha)] = \frac{1}{1-\frac{1}{e}} \int_{\frac{1}{e}}^1 \hat{T}(\alpha) d\alpha \stackrel{(8.3)}{\leq} \frac{e}{e-1} \int_{\frac{1}{e}}^1 \frac{1}{\alpha} T^* d\alpha = \frac{e}{e-1} T^* [\ln \alpha]_{\frac{1}{e}}^1 = \frac{e}{e-1} T^*. \quad (8.4)$$

Furthermore, we have

$$\begin{aligned} \mathbb{E}_\alpha\left[\sum_{j=1}^n s_j (1 - \hat{y}_j(\alpha))\right] &= \sum_{j=1}^n s_j \Pr_\alpha[y_j^* \leq \alpha] = \sum_{j=1}^n s_j \cdot \frac{1}{1-\frac{1}{e}} \int_{\max\{\frac{1}{e}, y_j^*\}}^1 d\alpha \\ &\leq \sum_{j=1}^n s_j \cdot \frac{e}{e-1} \int_{y_j^*}^1 d\alpha = \frac{e}{e-1} \sum_{j=1}^n s_j (1 - y_j^*). \end{aligned} \quad (8.5)$$

Hence,

$$\mathbb{E}_\alpha[\hat{T}(\alpha) + \sum_{j=1}^n s_j (1 - \hat{y}_j(\alpha))] \stackrel{(8.4),(8.5)}{\leq} \frac{e}{e-1} \left(T^* + \sum_{j=1}^n s_j (1 - y_j^*)\right) = \frac{e}{e-1} c_{LP} \leq \frac{e}{e-1} c_{MIP},$$

and thus we have proven an approximation factor of $\frac{e}{e-1}$ in expectation.

For derandomization, note that all possible solutions occur for $\alpha \in \{y_1^*, \dots, y_n^*\}$ and thus it is sufficient to compute rounded solutions for $\alpha \in (\frac{1}{e}, 1) \cap \{y_1^*, \dots, y_n^*\}$ (and $\frac{1}{e}$ if this set is empty). Furthermore, the best of these solutions must be at least as good as the expected value. \square

On the other hand, it is known that the scheduling problem cannot be approximated arbitrarily well. Such problems are referred to as *APX-hard*.

Theorem 8.16 ([26]). There is no PTAS for the above scheduling problem, provided that $P \neq NP$.

Bibliography

- [1] Tobias Achterberg, Thorsten Koch, and Alexander Martin. Branching rules revisited. *Oper. Res. Lett.*, 33(1):42–54, 2005.
- [2] E. Balas. Facets of the knapsack polytope. *Mathematical Programming*, 8:146–164, 1975.
- [3] E. Balas, S. Ceria, and G. Cornuéjols. A lift-and-project cutting plane algorithm for mixed 0-1 programs. *Mathematical Programming*, 58:295–324, 1993.
- [4] J.F. Benders. Partitioning procedures for solving mixed variables programming. *Numerische Mathematik*, 4:238–252, 1962.
- [5] C. Berge. Färbung von Graphen, deren sämtliche bzw. deren ungerade Kreise starr sind (Zusammenfassung). In *Wissenschaftliche Zeitschrift, Mathematisch-Naturwissenschaftliche Reihe*. Martin Luther Universität Halle-Wittenberg, 1961.
- [6] R. Borndörfer. *Aspects of Set Packing, Partitioning, and Covering*. PhD thesis, Technische Universität Berlin, 1998.
- [7] Claus C. Carøe and Rüdiger Schultz. Dual decomposition in stochastic integer programming. *Oper. Res. Lett.*, 24(1-2):37–45, 1999.
- [8] N. Christofides. Worst-case analysis of a new heuristic for the travelling salesman problem. Technical Report 388, Graduate School of Industrial Administration, CMU, 1976.
- [9] V. Chvátal. Edmonds polytopes and a hierarchy of combinatorial problems. *Discrete Mathematics*, 4:305–337, 1973.
- [10] V. Chvátal. On certain polytopes associated with graphs. *Journal on Combinatorial Theory B*, 18:305–337, 1975.
- [11] Michele Conforti, Gérard Cornuéjols, and M. R. Rao. Decomposition of balanced matrices. *J. Comb. Theory, Ser. B*, 77(2):292–406, 1999.
- [12] G.B. Dantzig and P. Wolfe. Decomposition principle for linear programs. *Operations Research*, 8:101–111, 1960.
- [13] Guoli Ding, Li Feng, and Wenan Zang. The complexity of recognizing linear systems with certain integrality properties. *Mathematical Programming*, 114(2):321–334, 2008.
- [14] J. Edmonds and R. Giles. A min-max relation for submodular functions on graphs. In P.L. Hammer, editor, *Studies in Integer Programming*, volume 1, pages 185–204. North-Holland, Amsterdam, 1977.
- [15] Jack Edmonds. Paths, trees, and flowers. *Canadian Journal of Mathematics*, 17:449–467, 1965.
- [16] D. Fulkerson, A.J. Hoffman, and R. Oppenheim. On balanced matrices. *Mathematical Programming Study*, 1:120–132, 1974.
- [17] D.R. Fulkerson. Packing rooted directed cuts in a weighted directed graph. *Mathematical Programming*, 6:1–13, 1974.
- [18] A. Ghouila-Houri. Caractérisation des matrices totalement unimodulaires. *Comptes Rendus Hebdomadaires des Séances de l'Académie des Sciences (Paris)*, 254:1192–1194, 1962.

-
- [19] R.E. Gomory. An algorithm for the mixed integer problem. Technical Report RM-2597, The RAND Cooperation, 1960.
- [20] R.E. Gomory. An algorithm for integer solutions to linear programming. In R.L. Graves and P. Wolfe, editors, *Recent Advances in Mathematical Programming*, pages 269–302, New York, 1969. McGraw-Hill.
- [21] M. Grötschel, L. Lovász, and A. Schrijver. *Geometric Algorithms and Combinatorial Optimization*. Springer, 1988.
- [22] P.L. Hammer, E.L. Johnson, and U.N. Peled. Facets of regular 0-1 polytopes. *Mathematical Programming*, 8:179–206, 1975.
- [23] Michael Held and Richard M. Karp. The traveling-salesman problem and minimum spanning trees. *Operations Research*, 18:1138–1162, 1970.
- [24] Michael Held and Richard M. Karp. The traveling-salesman problem and minimum spanning trees: Part II. *Mathematical Programming*, 1(1):6–25, 1971.
- [25] A.J. Hoffman and J.B. Kruskal. Integral boundary points of convex polyhedra. In H. W. Kuhn and A. W. Tucker, editors, *Linear inequalities and related systems*, pages 223–246. Princeton University Press, Princeton, NJ, 1956.
- [26] Han Hoogeveen, Martin Skutella, and Gerhard J. Woeginger. Preemptive scheduling with rejection. *Mathematical Programming*, 94:361–374, 2003.
- [27] Tom A. Jenkyns. The efficacy of the “greedy” algorithm. *Proc. 7th southeast. Conf. Comb., Graph Theory, Comput.*, pages 341–350, 1976.
- [28] Anna R. Karlin, Nathan Klein, and Shayan Oveis Gharan. A (slightly) improved approximation algorithm for metric tsp. *Operations Research*, 2023. To appear.
- [29] Marek Karpinski, Michael Lampis, and Richard Schmied. New inapproximability bounds for TSP. *Journal of Computer and System Sciences*, 81(8):1665–1677, 2015.
- [30] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, 1983.
- [31] D Klabjan, G.L Nemhauser, and C Tovey. The complexity of cover inequality separation. *Operations Research Letters*, 23(1-2):35–40, 1998.
- [32] B. Korte and J. Vygen. *Combinatorial Optimization. Theory and Algorithms*, volume 21 of *Algorithms and Combinatorics*. Springer, Heidelberg, 5th edition, 2012.
- [33] H. Marchand, A. Martin, R. Weismantel, and L.A. Wolsey. Cutting planes in integer and mixed integer programming. *Discrete Applied Mathematics*, 123/124:391–440, 2002.
- [34] H. Marchand and L.A. Wolsey. The 0-1 knapsack problem with a single continuous variable. *Mathematical Programming*, 85:15–33, 1999.
- [35] M.W. Padberg. On the facial structure of set packing polyhedra. *Mathematical Programming*, 5:199–215, 1973.
- [36] M.W. Padberg. A note on zero-one programming. *Operations Research*, 23(4):833–837, 1975.
- [37] M.W. Padberg. $(1, k)$ -configurations and facets for packing problems. *Mathematical Programming*, 18:94–99, 1980.
- [38] M.W. Padberg and G. Rinaldi. A branch and cut algorithm for the resolution of large-scale symmetric traveling salesman problems. *SIAM Review*, 33:60–100, 1991.
- [39] A. Schrijver. *Theory of Linear and Integer Programming*. Wiley, Chichester, 1986.
- [40] Paul D. Seymour. Decomposition of regular matroids. *Journal of Combinatorial Theory*, 28:305–359, 1980.

-
- [41] R. Van Slyke and R. Wets. L-shaped linear programs with applications to optimal control and stochastic programming. *SIAM Journal on Applied Mathematics*, 17(4):638–663, 1969.
- [42] R. Weismantel. On the 0/1 knapsack polytope. *Mathematical Programming*, 77(1):49–68, 1997.
- [43] L.A. Wolsey. Faces of linear inequalities in 0-1 variables. *Mathematical Programming*, 8:165–178, 1975.