

INAUGURAL-DISSERTATION

zur Erlangung der Doktorwürde
der Naturwissenschaftlich-Mathematischen Gesamtfakultät
der Ruprecht-Karls-Universität Heidelberg

vorgelegt von
Master of Science Francesco Silvestri
aus Eberbach am Neckar

Tag der mündlichen Prüfung:

Representations of Partition Problems and the Method of Moments

Betreuer:

Prof. Dr. Gerhard Reinelt
Prof. Dr. Christoph Schnörr

to my parents

Zusammenfassung

Die vorliegende Dissertation verfolgt zwei Ziele. Das erste besteht darin, verschiedene Darstellungen von Partitionen zu formulieren, zu erklären und miteinander in Verbindung zu bringen, die zur Modellierung von Partitionsproblemen verwendet werden können. Das zweite besteht darin, durch Einsatz der Momentmethode, einem Ansatz aus der polynomiellen Optimierung, konvexe Relaxierungen dieser Darstellungen zu konstruieren, um die globalen Optimallösungen der zugehörigen Partitionsprobleme abzuschätzen. Für Probleme wie dem Euklidisches k -Clustering weicht diese Methodik stark von den herkömmlichen Ansätzen ab, die sich vorwiegend mit Heuristiken und lokalen Optima beschäftigen. Da die Momentmethode konvexe Probleme konstruiert, liegt der Fokus der Arbeit darauf, Darstellungen zu finden und auszunutzen, deren Lösungsräume nur eine triviale Symmetriegruppe zulassen, damit die Lösungen der Relaxierung zu zulässigen Lösungen gerundet werden können.

Die in dieser Arbeit behandelten Darstellungen sind Assignment-, Partitions- und Projektionsmatrizen, sowie simpliziale Überdeckungen für eine verallgemeinerte Version des Euklidischen k -Clustering. Es werden Verbindungen sowie Übergänge zwischen den Matrizenklassen hergestellt und mit der Literatur verglichen, und es wird explizit nachgewiesen, wie Partitionsmatrizen auf natürliche Weise durch die Momentmethode aus Assignmentmatrizen hervorgehen.

Das Ausnutzen von Projektionsmatrizen ermöglicht es uns eine neue Formulierung für die Färbungszahl einzuführen, und die aus der Momentmethode resultierenden Relaxierungen werden mit der Lovász Theta Zahl verglichen. Es wird charakterisiert, unter welchen Bedingungen beide Relaxierungen übereinstimmen, und als erstes Hauptresultat liefern Relaxierungen von Binärmatrizen in diesem Fall bessere Ergebnisse als Relaxierungen von binären Eigenwerten.

Der letzte Teil der Arbeit beschäftigt sich mit dem sogenannten affinen Euklidischen k -Clustering, welches das Euklidische k -Clustering verallgemeinert. Als zweites Hauptresultat der Arbeit wird eine neue Methodik für dieses anspruchsvolle Problem eingeführt, die simpliziale Überdeckungen des Lösungsraums ausnutzt um eindeutige Darstellungen der Optimallösungen des zugrunde liegenden Problems zu ermöglichen. Im Gegensatz zum direkten Einsatz der Momentmethode auf die Standardformulierung ermöglicht der Einsatz auf dieser Formulierung einen langsameren Anstieg der Problemgröße, bessere Parallelisierbarkeit sowie die Möglichkeit Informationen als Grundlage für eine Rundungsheuristik zu erhalten, was aus Symmetriegründen bei der Standardformulierung nicht möglich ist.

Abstract

The thesis follows two main goals. The first is to formulate, explain and link representations of partitions that can be used to model partition problems. The second is to use the method of moments, an approach from polynomial optimization, to bound the global optima of the corresponding partition problems by constructing convex relaxations of these representations. For problems like Euclidean k -clustering, this is a stark contrast to their usual treatment, which mostly involves heuristics that are content with local optima. Since the method of moments results in a convex approach, the focus lies on finding and exploiting representations that lack a non-trivial symmetry-invariant solution space in order to be able to round the relaxations to feasible solutions.

The representations considered in the thesis are assignment matrices, partition matrices, projection matrices and simplicial covers for a generalized version of Euclidean k -clustering. Connections and transformations between the matrix classes are established and compared to the literature, and it is explicitly shown how partition matrices arise naturally from assignment matrices through the method of moments.

Using projection matrices, we are able to give a new formulation of the colouring number, and the resulting relaxations from the method of moments are compared to the Lovász theta number. We characterize under which circumstances the relaxations agree and explain when they do not, indicating our first main result that in this case, relaxing binary matrix entries yields better results than relaxing binary eigenvalues.

The final part of the thesis is devoted to what we call the affine Euclidean k -clustering problem, which is a more general version of the Euclidean k -clustering problem. As our second main result of the thesis, we introduce a new method for this challenging problem, utilizing simplicial covers of the feasible region to formulate unique representations of the optimal solutions of the underlying problem. In contrast to applying the method of moments directly, applying it to our formulation yields a slower growth in size, better parallelizability and enables us to recover information that can be used for rounding, which is not possible for the standard formulation due to symmetry.

Acknowledgement

This thesis could not have been written without the support of various people:

First, I want to thank my supervisor Prof. Gerhard Reinelt for introducing me to the fascinating topic of combinatorial optimization as a student, and for guiding me along my academic career from bachelor student to Ph.D. student, while always encouraging me to follow my interests and generously supporting my endeavors.

I also want to thank my second supervisor Prof. Christoph Schnörr for his generous support and input on various topics of this thesis, and in particular for sharing his perspective on the affine Euclidean clustering problem encountered in image analysis.

Next, I would like to thank the secretaries Catherine Proux-Wieland, Evelyn Wilhelm and Barbara Werner, who did a fantastic job and helped with all administrative tasks that occurred during my time at the Faculty of Mathematics and Computer Science of Heidelberg University.

I am very grateful for the time I was able to spend with all the wonderful people from the *Combinatorial Optimization* and the *Image & Pattern Analysis* groups, resulting in great memories and fruitful discussions both in and outside of work. In particular, I would like to thank Achim Hildenbrandt, Artjom Zern, Jan Kuske, Mattia Desana, Robert Breckner and Tobias Dencker for our time working together and the nice atmosphere they helped to shape, which I will remember for a long time.

Thanks for the support and funding of the German Research Foundation (DFG). I was member of the Research Training Group (RTG) 1653 "Spatio/Temporal Probabilistic Graphical Models and Applications in Image Analysis", which I acknowledge gratefully.

I also very much enjoyed my interactions with the people from AAU in Klagenfurt during my trips to visit Prof. Rendl, and I would like to thank him for introducing me to the topics of semidefinite optimization, as well as for his support. In the same vein, I'd like to thank all people that I met during the HeKKSaGOn project that allowed me to go to Kyoto, and Prof. Fujishige for introducing me to the LP-Newton method.

I'd like to thank my friends for the time we spent together, especially during our board game nights, which greatly helped me to relax. For proofreading parts of this thesis I thank Artjom Zern, Nadine Bär, Peter Gräf and Thomas Hölters, whose comments helped a lot to improve this thesis, and to ease my mind.

My parents Rita Silvestri and Pasquale Silvestri as well as my brother Alessandro Silvestri deserve my heartfelt thanks for their unconditional support throughout my life. At the same time, I thank Nadine Bär for her loving support and encouragement, as well as for the great time we have been spending together. She contributed a lot to making my time as a Ph.D. student a precious experience.

Contents

1	Introduction	1
2	Preliminaries	5
2.1	Linear Algebra	5
2.2	Combinatorial Structures	6
2.3	Computational Complexity	8
2.4	Convex Analysis	9
2.5	Algebra	12
2.6	Method of Moments	16
3	Computational Aspects	21
3.1	Method of Moments on a Variety	21
3.2	Conic Linear Programming	24
3.3	The LP-Newton Method for CLPs	25
3.3.1	The CLP-Newton Method	26
3.3.2	The Minimum-Norm-Point Algorithm	31
3.3.3	Linear Optimization on \mathcal{K} -Zonotopes	32
3.3.4	Experiments	39
4	Partitions and Assignment Matrices	42
4.1	Overview	42
4.2	Assignment Matrices	45
4.2.1	Symmetry induced Problems	48
4.2.2	Orbitopes	50
5	Partition Matrices	55
5.1	Overview	55
5.2	Connection to Combinatorial Moment Matrices	57
5.3	Convexification	59
5.3.1	Applying the Method of Moments to Partition Matrices	59
5.3.2	Applying the Method of Moments to Orbitopes	63

6	Projection Matrices	66
6.1	Overview	66
6.2	Convexification	70
6.3	Applications	74
6.3.1	Graph Colouring	74
6.3.2	Euclidean k -Clustering	82
7	Affine Euclidean Clustering	86
7.1	Overview	86
7.1.1	Problem Formulation	86
7.2	Simplicial Covers	90
7.2.1	Separating Simplicial Covers	93
7.3	Convexification	97
7.4	Related Approaches	101
7.4.1	Moment Sequences	102
7.4.2	Mixed Linear Regression	102
7.5	Rounding	103
7.6	Modifications	109
7.7	Applications	116
8	Conclusion	121
	Bibliography	124
	Index	128

List of Figures

1.1	Euclidean clustering	1
1.2	Graph colouring	2
1.3	Symmetry in Euclidean clustering	2
3.1	Update step for the Conic LP-Newton Method	30
3.2	Runtime of the Conic LP-Newton Method	40
3.3	Newton-steps of the Conic LP-Newton Method	40
3.4	MNP computations of the Conic LP-Newton Method	41
4.1	Visualization of Theorem 4.14.	50
4.2	Plot of distance between orbitope and barycenter	54
6.1	Relaxation of projection matrix visualized as third order tensor	73
7.1	Simplicial cover vs simplicial complex	90
7.2	Illustration of Theorem 7.28.	107
7.3	Decomposition of a square into triangles	112
7.4	Simplicial decompositions of a square	113
7.5	σ -skeletons of a square	114
7.6	Comparison of simplicial covers for Euclidean Clustering I	117
7.7	Comparison of simplicial covers for Euclidean Clustering II	117
7.8	Euclidean Clustering on a discrete solution space	118
7.9	Hyperplane Clustering I	119
7.10	Hyperplane Clustering II	119
7.11	Affine Hyperplane Clustering	120

List of Tables

- 6.1 Comparison of the growth of $|\text{ProMo}_{n,t}|$ and $\dim(\mathbb{R}_t[\mathbf{R}])$ 72
- 6.2 Relaxations for $G(n_1, n_2, n_3)$ 82
- 6.3 Relaxations for $G(n_1, \mathbf{e}_m)$ 82

List of Algorithms

3.1	Conic LP-Newton Method	29
3.2	Minimum-Norm-Point Algorithm	31
3.3	Linear Optimization over $[\mathbf{l}, \mathbf{u}]_{\mathcal{L}_n}$	37
7.1	Farthest Point Clustering	108
7.2	Deterministic Rounding	109

List of Symbols

Basics

\mathbb{N}	non-negative integers
\mathbb{N}_t^d	d -dimensional non-negative integer vectors whose sum is at most t
\mathbb{R}	real numbers
\mathbb{R}_+	non-negative real numbers
2^S	power set of the finite set S
$[n]$	the set $\{1, 2, \dots, n\}$ for $n \in \mathbb{N}$, $n \geq 1$
$[U]$	orbit of U under a group action <i>or</i> equivalence class
\mathfrak{S}_n	set of permutations on $[n]$

Partitions

$\mathcal{P}^n(L)$	set of all partitions of $[n]$ with parts restricted to $L \subseteq 2^{[n]}$
$\mathcal{P}_k^n(L)$	set of all partitions of $[n]$ with k parts restricted to $L \subseteq 2^{[n]}$
$\mathcal{U}_{n,k}(L)$	set of $n \times k$ assignment matrices arising from $\mathcal{P}_k^n(L)$
$\mathcal{U}_{n,k}^{lex}(L)$	set of matrices in $\mathcal{U}_{n,k}(L)$ with lexicographic sorted columns
$\text{PM}_k^n(L)$	set of $n \times n$ partition matrices arising from $\mathcal{P}_k^n(L)$
$\text{CProM}_k^n(L)$	set of projection matrices arising from $\mathcal{P}_k^n(L)$

Algebra

\mathbf{x}^α	the monomial $\prod_{i \in [d]} x_i^{\alpha_i}$ for $\alpha \in \mathbb{N}^d$, $\mathbf{x} \in \mathbb{R}^d$
$\mathbb{R}[\mathbf{x}]$	ring of real multivariate polynomials in $\mathbf{x} = (x_1, \dots, x_d)$
$\mathbb{R}_t[\mathbf{x}]$	vector space of polynomials in $\mathbb{R}[\mathbf{x}]$ of degree at most t
$z_d(t)$	dimension of $\mathbb{R}_t[\mathbf{x}]$ where $\mathbf{x} = (x_1, \dots, x_d)$, equal to $\binom{d+t}{d}$
$v_t(\mathbf{x})$	vector of moments of $\mathbf{x} \in \mathbb{R}^d$ up to degree t
$\text{td}(f)$	truncated degree of a multivariate polynomial f given by $\left\lceil \frac{\deg(p)}{2} \right\rceil$

Combinatorial Objects

\mathcal{G}_n	set of all simple graphs with vertex set $[n]$
L	an independence system
L_*	an independence system without the empty set
\mathcal{T}	the partition $\{T_1, \dots, T_k\}$
$\left\{ \begin{matrix} n \\ k \end{matrix} \right\}$	Stirling number of the second kind

Convex Cones

\mathcal{K}	a proper cone
\mathbb{R}_+^n	cone of vectors in \mathbb{R}^n with non-negative entries
\mathcal{S}_+^n	cone of positive semidefinite $n \times n$ matrices
\mathcal{L}_n	cone of vectors $(x_0, \mathbf{x}) \in (\mathbb{R}_+ \times \mathbb{R}^n)$ satisfying $\ \mathbf{x}\ _2 \leq x_0$

Polytopes

C_k^d	k -truncated d -dimensional unit cube
Δ_Ω^d	Ω -constrained standard simplex in \mathbb{R}^d

Vectors / Matrices

\mathbf{e}_I	characteristic vector of $I \subseteq [n]$
\mathbf{e}_i	shorthand for $\mathbf{e}_{\{i\}}$
\mathbf{e}	shorthand for $\mathbf{e}_{[n]}$ if n is clear from context
\mathbf{I}_n	$n \times n$ identity matrix
$\mathbf{J}_{n,m}$	$n \times m$ matrix of all ones defined as $\mathbf{J}_{n,m} = \mathbf{e}_{[n]} \mathbf{e}_{[m]}^\top$
Ω_L	$n \times n$ binary matrix defined by independence system $L \subseteq 2^{[n]}$
\mathbf{A}_G	$n \times n$ adjacency matrix of the graph $G \in \mathcal{G}_n$
$\mathbf{M}_t(f, \mathbf{y})$	localizing moment matrix of order t of \mathbf{y} with respect to f

Relations

\subseteq	inclusion as subset (with possible equality)
\leftrightarrow	existence of a bijection between two sets
\leq	conic order induced by \mathbb{R}_+^n or subgraph relation
\leq	conic order induced by \mathcal{S}_+^n
$\leq_{\mathcal{K}}$	conic order induced by \mathcal{K}
$<_{lex}$	lexicographic monomial order on \mathbb{N}^d

Maps / Operators

\wedge	join in the context of orders, logical AND otherwise
\vee	meet in the context of orders, logical OR otherwise
\otimes	tensor product for matrices <i>or</i> product of independence systems
δ_Ψ	Kronecker delta assuming the truth values 0/1 of Ψ
*	group action
$\text{rank}(\mathbf{X})$	the rank of the matrix \mathbf{X}
$\text{col}(\mathbf{X})$	the set of column vectors of the matrix \mathbf{X}
$\text{tr}(\mathbf{X})$	trace of the quadratic matrix \mathbf{X}
$\text{spct}(\mathbf{X})$	set of eigenvalues of the quadratic matrix \mathbf{X}
$\text{diag}(\mathbf{X})$	diagonal of the square matrix \mathbf{X} as vector
$\text{Diag}(\mathbf{x})$	diagonal matrix \mathbf{X} with diagonal equal to \mathbf{x}
$\text{supp}(\mathbf{x})$	set of indices $i \in [n]$ for which $\mathbf{x} \in \mathbb{R}^n$ is non-zero
$\ \mathbf{x}\ _0$	shorthand for $ \text{supp}(\mathbf{x}) $
$\text{dist}(\mathbf{x}, C)$	Euclidean distance between $\mathbf{x} \in \mathbb{R}^n$ and the convex set $C \subseteq \mathbb{R}^n$
$\pi_C(\mathbf{x})$	orthogonal projection from $\mathbf{x} \in \mathbb{R}^n$ onto the convex set $C \subseteq \mathbb{R}^n$
$\text{bd}(C)$	boundary of the set $C \subseteq \mathbb{R}^n$
$\text{int}(C)$	interior of the set $C \subseteq \mathbb{R}^n$
$\text{lin}(C)$	linear span of the set $C \subseteq \mathbb{R}^n$
$\text{aff}(C)$	affine hull of the set $C \subseteq \mathbb{R}^n$
$\text{conv}(C)$	convex hull of the set $C \subseteq \mathbb{R}^n$
$\text{bar}(C)$	barycenter of C
Id_C	identity map on C
$\text{vert}(P)$	set of vertices of a polytope P
$\text{CC}(G)$	set of connected components of graph G

Abbreviations

MM	method of moments
psd.	positive semidefinite

Chapter 1

Introduction

Overview

In this thesis, we are interested in the general problem of optimally partitioning a set of mathematical objects according to a given criterion. A prominent geometric example for this is *Euclidean k -clustering*, where the goal is to partition a set of points $\{\mathbf{b}_i \mid i \in [n]\} \subseteq \mathbb{R}^d$ into k parts such that each part has a center \mathbf{x}^j with low average Euclidean distance towards its members, as can be seen in Figure 1.1.

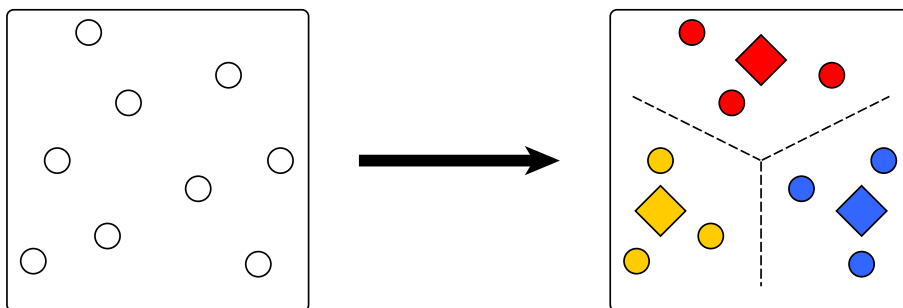


Figure 1.1: Example of Euclidean 3-clustering. Circles correspond to the input \mathbf{b}_i , diamonds to centers \mathbf{x}^j and colours to parts.

Partition problems constitute a large class of hard and well-studied problems with underlying combinatorial structure and are ubiquitous in the fields of data science. A purely combinatorial example for a notoriously hard problem belonging to this class is the one of *graph-colouring*. Given a simple graph G , the goal is to find the least amount of colours necessary to colour each node, such that the same colour is never connected with an edge. Figure 1.2 shows some minimal colourings for small graphs.

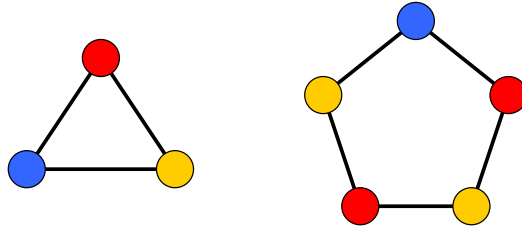


Figure 1.2: Examples of minimal colourings.

Unfortunately, our examples, as well as a lot of other interesting partition problems, are NP-hard, and a lot of work has been done to tackle these problems. Investigating these problems, it quickly becomes apparent that a key difficulty of formulating partition problems is their inherent symmetrical structure. While a partition is mathematically a *set* of parts, a computer necessarily needs to store the parts in some order, thus representing a partition as a *list* of parts. This transition introduces a fixed, but arbitrary order on the parts, which means that any solution has one representation for each permutation of its parts, leading to overblown solution spaces and naturally ill-posed problems.

These symmetries are especially problematic for approaches relying purely on convex optimization, since their optima tend to lie inside the interior of the convex hull of the optimal feasible solutions. In the worst case, such an optimum may not even provide any means to recover an actual feasible solution, as shown in Figure 1.3.

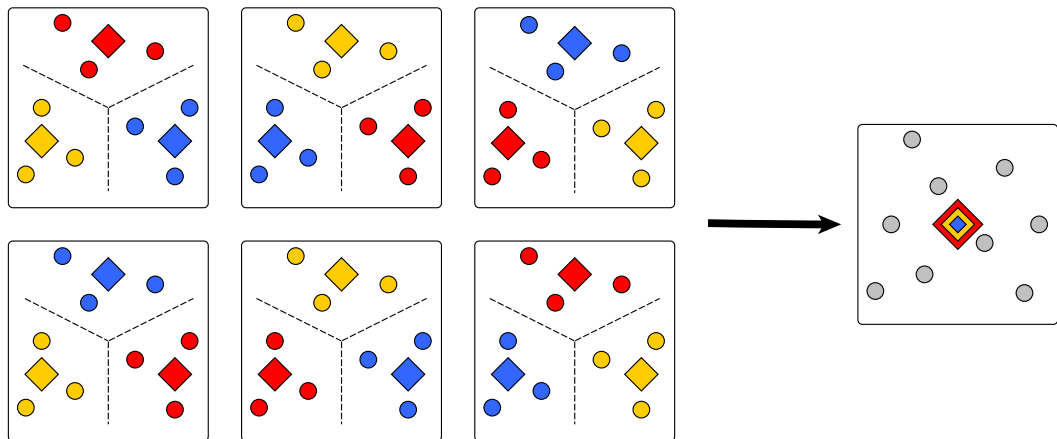


Figure 1.3: Left: All $3! = 6$ optimal solutions for the Euclidean 3-clustering instance from Figure 1.1. Right: Convex combination of the optimal solutions both in colours and center locations, with no means to go back to the left.

For this reason, it is important to find proper representations of partitions that circumvent the problems induced by symmetries and are still tractable from a computational point of view as well. There are still many open questions in this regard, and the aim of this thesis is to give some insight in how to approach them.

Organization

This thesis pursues two main goals. The first is to formulate, explain and link representations of partitions. The second is to use the *method of moments*, an approach from polynomial optimization, to construct convex relaxations of these representations, in order to bound the global optima of the corresponding partition problems. For problems like Euclidean k -clustering, this is a stark contrast to their usual treatment, which mostly involves heuristics that are content with local optima.

In Chapter 2, we recall the theoretical results that are used as mathematical foundation throughout the thesis. While mostly self-contained, we keep the explanation of results to a bare minimum, and proofs are omitted. We highly recommend to inspect the listed books there for further information.

Chapter 3 complements Chapter 2 by going into the details of the underlying computational aspects of this thesis. In particular, we explore preprocessing to reduce the problem size, explain the computational setup and describe our adaption of the LP-Newton method for conic linear programs, which may be skipped on a first reading.

Our treatment of partition problems starts in Chapter 4, where they are formally defined. We show how to represent partitions in terms of *assignment matrices* and properly illustrate the issues of symmetrical solutions in a convex setting, as well as a potential fix for this in the form of orbitopes.

In Chapter 5, we review *partition matrices* used in combinatorial optimization, our first matrix class that uniquely represents each partition. In particular, we explicitly show that the method of moments connects them to the assignment matrices from the preceding chapter, a fact that has implicitly been exploited for a long time in literature.

Chapter 6 shows how to transition from partition matrices to *projection matrices*, our second matrix class that uniquely represents each partition. We use this class to introduce a new relaxation for the graph-colouring problem and show how it relates to established relaxations like the Lovász *theta number* $\vartheta(G)$.

Finally, Chapter 7 is devoted to the *affine Euclidean k -clustering problem*, which is a more general version of the Euclidean k -clustering previewed here. As one main result of this thesis, we introduce a new method for this challenging problem, utilizing *simplicial covers* of the feasible region to formulate unique representations of the optimal solutions of the underlying problem.

Parts of this thesis have been published in the following papers:

- [SRS15] Francesco Silvestri, Gerhard Reinelt, and Christoph Schnörr. A convex relaxation approach to the affine subspace clustering problem. In *Pattern Recognition - 37th German Conference, GCPR 2015, Aachen, Germany, October 7-10, 2015, Proceedings*, pages 67–78, 2015
- [SRS16] Francesco Silvestri, Gerhard Reinelt, and Christoph Schnörr. Symmetry-free SDP relaxations for affine subspace clustering. *ArXiv e-prints*, July 2016
- [SR16] Francesco Silvestri and Gerhard Reinelt. The LP-Newton method and conic optimization. *ArXiv e-prints*, November 2016

In particular, Section 3.3 contains material from [SR16], and Chapter 7 contains material from [SRS15, SRS16].

Chapter 2

Preliminaries

This chapter recalls the underlying definitions and most important results that will be used throughout the thesis. It is assumed that the reader has some basic understanding of linear algebra, and proofs are omitted in general, but can be found in the given references.

Notation

The majority of the commonly used notation can be found in the List of Symbols, while the rest will be introduced when appropriate. In general, we will use small letters like a, x, λ for scalars or elements of a set and capital letters like S, T, U for sets. Small bold letters are used for vectors like $\mathbf{a}, \mathbf{x}, \boldsymbol{\lambda}$ and capital bold letters for matrices like $\mathbf{A}, \mathbf{X}, \boldsymbol{\Lambda}$. Instead of writing index sets of small size like $\{i\}$ or $\{i, j\}$, we will sometimes simply write i or ij respectively. For a formal statement Ψ , the *Kronecker delta* δ_Ψ assumes the value 1 if the statement Ψ is true and 0 otherwise; for a pair of objects a, b in the same space, $\delta_{a,b}$ denotes the shorthand for $\delta_{a=b}$.

For the sets $S \subseteq \mathbb{R}$ and $C, D \subseteq \mathbb{R}^n$, we use the shorthand set operations

$$\begin{aligned} S \cdot C &:= \{s \cdot c \mid s \in S, c \in C\}, \\ C + D &:= \{c + d \mid c \in C, d \in D\}, \end{aligned}$$

where the latter is known as *Minkowski sum*. In particular, $-C$ is short for $-1 \cdot C$ and corresponds to a point reflection of C at $\mathbf{0}$.

2.1 Linear Algebra

For a quadratic matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, its *trace* is defined as

$$\text{tr}(\mathbf{A}) := \sum_{i \in [n]} a_{ii},$$

which is equivalent to the sum of its eigenvalues. The trace induces an inner product on the space of $n \times m$ matrices called the *trace product*, which is defined by

$$\langle \mathbf{A}, \mathbf{B} \rangle := \text{tr}(\mathbf{A}\mathbf{B}^\top) = \sum_{i \in [n]} \sum_{j \in [m]} a_{ij} b_{ij}.$$

For $m = 1$, this reduces to the usual notion of the scalar product for vectors in \mathbb{R}^n .

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be invertible and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$. Then the rank 1 update $\mathbf{A} + \mathbf{u}\mathbf{v}^\top$ is invertible if and only if $1 + \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u} \neq 0$. Furthermore, if the inverse exists, it is given by the *Sherman-Morrison formula*

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^\top)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{u}\mathbf{v}^\top \mathbf{A}^{-1}}{1 + \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u}}. \quad (2.1)$$

2.2 Combinatorial Structures

Graphs

A good introduction for graphs can be found in the textbook [Die05]. A (*simple*) graph G is a tuple (V, E) consisting of a finite set V called *nodes* (or *vertices*) and a set

$$E \subseteq (V \times V) \setminus \{(v, v) \mid v \in V\}$$

called *edges*. A (*graph*) *homomorphism* from graph $H = (V_H, E_H)$ to graph $G = (V_G, E_G)$ is a map $\varphi : V_H \rightarrow V_G$ such that $(i, j) \in E_H$ if and only if $(\varphi(i), \varphi(j)) \in E_G$. A graph H is called a *subgraph* of a graph G , which we denote by $H \leq G$, if there is an injective homomorphism from H to G . For a graph $G = (V_G, E_G)$ and a subset $U \subseteq V_G$, the subgraph $G[U]$ induced by U is defined as the graph $G[U] := (U, E_G \cap (U \times U))$.

A bijective graph homomorphism is called a (*graph*) *isomorphism*, and a graph isomorphism from a graph to itself is called a (*graph*) *automorphism*. A graph $G = (V, E)$ is called *vertex-transitive*, if, for every pair of nodes $i \neq j \in V$, there exists an automorphism φ_{ij} of G such that $\varphi_{ij}(i) = j$.

Given a graph $G = (V, E)$ whose node set has finite size $|V| = n \in \mathbb{N}$, we can always relabel the nodes as natural numbers by constructing a graph homomorphism from V to $[n]$ with appropriate edge set. For this reason, we can define the set of graphs of size n as

$$\mathcal{G}_n := \{G \mid G = ([n], E) \text{ is a graph}\},$$

and define with $E(G)$ the set of edges of G .

Given a graph $G \in \mathcal{G}_n$, its *adjacency matrix* is defined as the unique symmetrical $n \times n$ binary matrix \mathbf{A}_G with the property that $(\mathbf{A}_G)_{ij} = 1$ if and only if $\{i, j\} \in E$. In particular, this assignment yields a bijection for

$$\mathcal{G}_n \leftrightarrow \{\mathbf{A} \in \{0, 1\}^{n \times n} \mid \mathbf{A} = \mathbf{A}^\top, \text{diag}(\mathbf{A}) = \mathbf{0}\}.$$

The *complement (graph)* of a graph $G \in \mathcal{G}_n$ is the graph $\overline{G} \in \mathcal{G}_n$ satisfying

$$\mathbf{A}_G + \mathbf{A}_{\overline{G}} + \mathbf{I}_n = \mathbf{J}_n.$$

For a pair of nodes $i, j \in [n]$ in a graph $G \in \mathcal{G}_n$, an (i, j) -path is a non-repeating sequence of nodes v_0, v_1, \dots, v_m such that $(v_{k-1}, v_k) \in E(G)$ for all $k \in [m]$ and $v_0 = i, v_m = j$. The pair (i, j) is called *connected* if there exists an (i, j) -path, and being connected is an equivalence relation. The corresponding equivalence classes are called *connected components*, and denoting the set of all connected components of G by $\text{CC}(G)$ induces a surjection

$$\text{CC} : G \mapsto \text{CC}(G) \in \mathcal{P}^n$$

from \mathcal{G}_n to \mathcal{P}^n , the set of partitions of $[n]$. A graph $G \in \mathcal{G}_n$ is called *bipartite* if its nodes $[n]$ can be partitioned into two sets U, V such that $E(G[U]) = E(G[V]) = \emptyset$.

Independence Systems

A subset $L \subseteq 2^{[n]}$ is called an *independence system* or *abstract simplicial complex* if it is closed under taking subsets; more formally, the defining quality is the implication

$$I \in L \Rightarrow 2^I \subseteq L. \quad (2.2)$$

Without loss of generality, we will assume that whenever $L \subseteq 2^{[n]}$ is an independence system, then $\{i\} \in L$ for all $i \in [n]$; otherwise, we could consider L as independence system of a strictly contained subset of $[n]$.

It will often be useful to consider an independence system without the empty set; for this reason, we define $L_* = L \setminus \{\emptyset\}$.

Independence systems are closed under taking unions in the sense that whenever $L \subseteq 2^J$ and $L' \subseteq 2^{J'}$ are independence systems, so is the set of individual unions

$$L \otimes L' := \{I \cup I' \mid I \in L, I' \in L'\} \subseteq 2^{J \cup J'}.$$

Abusing notation, we will write L^k for the independence system emerging from the independence system $L \subseteq 2^{[n]}$ by using this construction on the union of k distinct copies of $[n]$, and likewise L_*^k for the construction with L_* , where we stress that $(L^k)_* \neq L_*^k$.

For later, we will also introduce the matrix $\Omega_L = \{\omega_{ij}\}_{i,j \in [n]} \in \{0, 1\}^{n \times n}$, which we define pointwise as

$$\omega_{ij} = \begin{cases} 1 & \text{if } \{i, j\} \notin L, \\ 0 & \text{else.} \end{cases}$$

We can consider the matrix Ω_L as the adjacency matrix of the graph that represents forbidden pairs in L .

Order Theory

A (non-strict) *partial order* on a set C is a binary relation \preceq on C that satisfies

- (i) $a \preceq a$ (reflexivity),
- (ii) $a \preceq b$ and $b \preceq a$ imply $a = b$ (antisymmetry),
- (iii) $a \preceq b$ and $b \preceq c$ imply $a \preceq c$ (transitivity)

for all $a, b, c \in C$. A partial order \preceq on C induces a binary relation \triangleleft called *strict partial order* for which $a \triangleleft b$ is true if and only if $a \preceq b$ and $a \neq b$. A set C together with a partial order \preceq is called a *partial ordered set*, or *poset* for short. A partial order \preceq on C is called a *total order*, if for each pair $a, b \in C$, either $a \preceq b$ or $b \preceq a$.

For a subset S of a poset C , any element $u \in C$ that satisfies $s \preceq u$ for all $s \in S$ is called an *upperbound* of S . Furthermore, if u^* is an upperbound of $S \subseteq C$ with the property that $u^* \preceq u$ for all upperbounds u of S , then u^* is called the *join* (or *supremum*) of S . If (C, \preceq) is a poset, then the upperbound and join of S in (C, \succeq) are called *lowerbound* and *meet* (or *infimum*) in (C, \preceq) , respectively. The meet and join of a set $S \subseteq C$ are unique when they exist. A poset C for which every pair $a, b \in C$ has a unique meet $a \vee b$ and a unique join $a \wedge b$ is called a *lattice*.

A total order is called a *well-order* if every nonempty subset contains a lowerbound of itself.

2.3 Computational Complexity

We will only give an informal overview of some basic notions here. For a formal introduction, consider the book [AB09].

In the context of this thesis, an algorithm will be treated as a finite list of *elementary instructions* which work on some input data and produce some output data. We assume that each elementary instruction takes the same constant time to be carried out by a computer and that arithmetic operations on rational numbers can be treated as elementary instructions.

We use the term *problem* to describe an assignment of values to a subset of possible inputs to an algorithm. This way, a problem Q defines a partial function f_Q on the sets of all inputs. The inputs for which f_Q is well-defined are called *valid* (for Q). An algorithm A is *correct* (for Q), or rather *solves* Q , if it returns the value of f_Q for all valid inputs, that is, if A computes f_Q . The problem Q is called a *decision-problem* if f_Q is 0/1 valued, where we can interpret 0 as *no* and 1 as *yes*.

The input and output of the algorithm consists of structural and numerical data and we consider the number of numbers in the input as its *input size*. For a problem Q , we can assign a runtime function $t_A : \mathbb{N} \rightarrow \mathbb{N}$ to every algorithm A by defining $t_A(n)$

as the maximal number of elementary instructions that A uses on any valid input of input size n .

The function t_A is well-defined if it only attains finite values, that is, if A *terminates* for any valid input of Q . Note that t_A only predicts the performance in the worst case and may not be indicative for the average case.

For two functions $f : \mathbb{N} \rightarrow \mathbb{N}$ and $g : \mathbb{N} \rightarrow \mathbb{N}$, the *big O notation* is defined as the relation

$$f = O(g) \Leftrightarrow \exists c > 0 \exists n_0 \in \mathbb{N} \forall n > n_0 : f(n) \leq c \cdot g(n).$$

In other words, $f = O(g)$ implies that f is not growing faster than g asymptotically.

A problem Q is said to be *efficiently solvable* if there is a correct algorithm A that solves Q with $t_A \in O(n^k)$ for some fixed $k \in \mathbb{N}$, and we call A efficient in this case. The *complexity-class* P contains all efficiently solvable decision problems. A decision problem Q belongs to the class NP if there is an algorithm A such that for every valid input x with $f_Q(x) = 1$, there is a *certificate* y_x with the property that A can efficiently verify $f_Q(x) = 1$ with input (x, y_x) .

While $P \subseteq NP$ is known, deciding whether $P = NP$ or $P \neq NP$ is one of the major open problems in theoretical computer science. For this thesis, we will assume the following.

Conjecture 2.1:

Finding a solution is harder than verifying it, or, in other words, $P \neq NP$.

A problem Q is called *NP-hard* when an efficient algorithm for Q can be converted into an efficient algorithm for every problem that belongs to the class NP . In particular, if any NP -hard problem is efficiently solvable, then $P=NP$.

2.4 Convex Analysis

A comprehensive treatment of convex analysis can be found in [Roc09].

A set $C \subseteq \mathbb{R}^d$ is called *convex* if for any pair $\mathbf{x}, \mathbf{y} \in C$, their connecting line segment

$$\{\mu \cdot \mathbf{x} + (1 - \mu) \cdot \mathbf{y} \mid \mu \in [0, 1]\}$$

is contained in C as well. Given any set $C \subseteq \mathbb{R}^d$, its *convex hull* $\text{conv}(C)$ is the smallest convex set that contains C (with regards to set inclusion). A *hyperplane* $H \subseteq \mathbb{R}^d$ is a set of the form $\{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{a} \rangle = b\}$ for $\mathbf{a} \in \mathbb{R}^d$ and $b \in \mathbb{R}$. The two associated closed *halfspaces* are $H_- := \{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{a} \rangle \leq b\}$ and $H_+ := \{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{a} \rangle \geq b\}$. A hyperplane *supports* a convex set C if $H \cap C \neq \emptyset$ and either $C \subseteq H_+$ or $C \subseteq H_-$, and the intersection $H \cap C$ is called a *face* of C .

A real valued function $f : C \rightarrow \mathbb{R}$ defined on a convex set $C \subseteq \mathbb{R}^n$ is called *convex* if for any pair $\mathbf{x}, \mathbf{y} \in C$, its graph is below their connecting line segment, or, more

formally, if

$$f(\mu \cdot \mathbf{x} + (1 - \mu) \cdot \mathbf{y}) \leq \mu \cdot f(\mathbf{x}) + (1 - \mu) \cdot f(\mathbf{y})$$

for all choices of $\mu \in [0, 1]$. For a convex function $f : C \rightarrow \mathbb{R}$, its *subdifferential* $\partial f(\mathbf{x}_0)$ at a given point \mathbf{x}_0 parametrizes the supporting hyperplanes of the epigraph of f that pass through $(\mathbf{x}_0, f(\mathbf{x}_0))$. More formally, we have

$$\partial f(\mathbf{x}_0) := \{\mathbf{y} \in \mathbb{R}^n \mid f(\mathbf{x}) - f(\mathbf{x}_0) \geq \langle \mathbf{y}, \mathbf{x} - \mathbf{x}_0 \rangle \forall \mathbf{x} \in C\}.$$

The elements of the subdifferential are called *subgradients* and generalize the gradient for differentiable functions.

For a closed convex set $C \subseteq \mathbb{R}^d$ and a point $\mathbf{x} \in \mathbb{R}^d$, we can define their *Euclidean distance* as

$$\text{dist}(\mathbf{x}, C) := \min \{\|\mathbf{x} - \mathbf{y}\|_2 \mid \mathbf{y} \in C\},$$

which is a convex function in \mathbf{x} . The minimizer of $\text{dist}(\mathbf{x}, C)$ is unique and called the *orthogonal projection* $\pi_C(\mathbf{x})$.

Lemma 2.2:

Let $C \subseteq \mathbb{R}^d$ be closed and convex, then

$$\frac{\mathbf{x} - \pi_C(\mathbf{x})}{\text{dist}(\mathbf{x}, C)} \in \partial \text{dist}(\mathbf{x}, C).$$

Convex Cones

A set $\mathcal{K} \subseteq \mathbb{R}^d$ is called a *cone* if it is invariant under positive rescaling, or, more formally, if $\mathbb{R}_+ \cdot \mathcal{K} \subseteq \mathcal{K}$. A cone is convex when it is closed under addition as well, so when $\mathcal{K} + \mathcal{K} \subseteq \mathcal{K}$. A convex cone \mathcal{K} is called *pointed* if it has the property $\mathcal{K} \cap -\mathcal{K} = \{\mathbf{0}\}$ and *proper* if it is pointed, closed and has nonempty interior. The *dual cone* \mathcal{K}^* of a convex cone \mathcal{K} is defined as

$$\mathcal{K}^* := \{\mathbf{y} \in \mathbb{R}^n \mid \langle \mathbf{x}, \mathbf{y} \rangle \geq 0 \forall \mathbf{x} \in \mathcal{K}\},$$

and a cone \mathcal{K} for which $\mathcal{K} = \mathcal{K}^*$ is called *self-dual*.

A proper cone $\mathcal{K} \subseteq \mathbb{R}^d$ induces a corresponding *conic (partial) order* $\leq_{\mathcal{K}}$ on \mathbb{R}^d by setting

$$\mathbf{y} \leq_{\mathcal{K}} \mathbf{x} \quad \Leftrightarrow \quad \mathbf{x} - \mathbf{y} \in \mathcal{K}.$$

In particular, membership $\mathbf{x} \in \mathcal{K}$ is equivalent to stating $\mathbf{x} \geq_{\mathcal{K}} \mathbf{0}$.

Example 2.3:

The set \mathbb{R}^n can be considered as a convex cone with dual cone $(\mathbb{R}^n)^* = \{\mathbf{0}\}$, but is not pointed and thus not proper. Well-known proper, self-dual cones and their conic orders

include the *non-negative orthant* (\mathbb{R}_+^n, \leq) , the cone of symmetrical positive semidefinite $n \times n$ matrices (\mathcal{S}_+^n, \leq) denoted below and the *Lorenz-cone* $(\mathcal{L}_n, \leq_{\mathcal{L}})$ defined as

$$\mathcal{L}_n = \{(x_0, \tilde{\mathbf{x}}) \in \mathbb{R}_+ \times \mathbb{R}^n \mid \|\tilde{\mathbf{x}}\|_2 \leq x_0\}.$$

Furthermore, the poset (\mathbb{R}_+^n, \leq) is a lattice.

Positives Semidefinite Matrices

A quadratic, symmetrical matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is called *positive semidefinite*, or psd. for short, if the corresponding quadratic form $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is non-negative for all $\mathbf{x} \in \mathbb{R}^n$. Equivalently, there exists a *Cholesky decomposition* of \mathbf{A} , which means there is a matrix $\mathbf{V} \in \mathbb{R}^{n \times \text{rank}(\mathbf{A})}$ such that $\mathbf{A} = \mathbf{V}\mathbf{V}^\top$. A direct consequence is the following lemma.

Lemma 2.4:

Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ and $\mathbf{C} = \mathbf{B}\mathbf{A}\mathbf{B}^\top$ where \mathbf{B} is invertible. Then $\mathbf{A} \geq \mathbf{0} \Leftrightarrow \mathbf{C} \geq \mathbf{0}$.

Furthermore, every psd. matrix \mathbf{A} has a unique psd. square root $\mathbf{A}^{\frac{1}{2}}$ that satisfies $(\mathbf{A}^{\frac{1}{2}})^2 = \mathbf{A}$. The *Schur complement* of a psd. matrix

$$\mathbf{M} = \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{C} \end{pmatrix} \geq \mathbf{0} \quad (2.3)$$

with invertible \mathbf{A} is defined as $\mathbf{M}/\mathbf{A} = \mathbf{C} - \mathbf{B}^\top \mathbf{A}^{-1} \mathbf{B}$. The following is a central argument.

Lemma 2.5 (Schur complement Lemma):

Let \mathbf{M} be as in (2.3). Then $\mathbf{M} \geq \mathbf{0}$ if and only if $\mathbf{M}/\mathbf{A} \geq \mathbf{0}$. Furthermore, \mathbf{M} is invertible if and only if \mathbf{M}/\mathbf{A} is invertible.

Polytopes

The following is a minimal introduction to polytopes, and we refer the reader to [Zie95] for a more elaborate treatment of the topic.

A convex set $P \subseteq \mathbb{R}^d$ that can be written as the convex hull of a finite set of points $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_m)$ is called a (*convex*) *polytope* and has the so-called *inner description*

$$P = \{\mathbf{x} \in \mathbb{R}^d \mid \mathbf{x} = \mathbf{V}\boldsymbol{\lambda}, \boldsymbol{\lambda} \geq \mathbf{0}, \langle \boldsymbol{\lambda}, \mathbf{e} \rangle = 1\}.$$

Example 2.6:

The set

$$\Delta^d := \{\boldsymbol{\lambda} \in \mathbb{R}_+^d \mid \langle \boldsymbol{\lambda}, \mathbf{e} \rangle = 1\} = \text{conv}(\{\mathbf{e}_1, \dots, \mathbf{e}_d\})$$

used to parametrize a polytope in its inner description is called (*standard*) *simplex* and is itself a polytope.

Alternatively, each polytope can be described by a finite intersection of halfspaces, which leads to its so-called *outer description*

$$P = \{\mathbf{x} \in \mathbb{R}^d \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}\}.$$

The faces of a polytope $P \subseteq \mathbb{R}^d$ are polytopes themselves, and its zero-dimensional faces are called its *vertices*, denoted by $\text{vert}(P)$. For any finite set $S \subseteq \mathbb{R}^d$, it holds that $\text{vert}(\text{conv}(S)) \subseteq S$, and $\text{conv}(\text{vert}(P)) = P$ holds for all polytopes $P \subseteq \mathbb{R}^d$.

Example 2.7:

The k -truncated d -dimensional *unit-cube* C_k^d has the outer description

$$C_k^d = \{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{e} \rangle \leq k, 0 \leq x_i \leq 1 \forall i \in [d]\}$$

and its inner description can be constructed from its vertices

$$\text{vert}(C_k^d) = \{\mathbf{x} \in \{0, 1\}^d \mid \langle \mathbf{x}, \mathbf{e} \rangle \leq k\}.$$

The shorthand $C^d := C_d^d$ is used for the d -dimensional unit-cube itself.

One major result about polytopes is the following theorem.

Theorem 2.8 (Carathéodory's theorem, [Die05, Prop. 1.15]):

Let $V \subseteq \mathbb{R}^d$ be a finite set of points and consider a point $\mathbf{x} \in \mathbb{R}^d$ in the polytope $\text{conv}(V)$. Then $\mathbf{x} \in \text{conv}(V')$ for a subset $V' \subseteq V$ of size at most $|V'| \leq \dim(\text{conv}(V)) + 1$.

2.5 Algebra

For a good introduction to computer algebra and the necessary foundations, we recommend the book [CLO15].

Polynomials, Ideals, Varieties

Given vectors $\boldsymbol{\alpha} \in \mathbb{N}^d$ and $\mathbf{x} \in \mathbb{R}^d$, the product $\mathbf{x}^\alpha = \prod_{i \in [d]} x_i^{\alpha_i}$ is called a *monomial* whose *total degree* is defined as $\deg(\mathbf{x}^\alpha) := \langle \mathbf{e}, \boldsymbol{\alpha} \rangle$. Let $\mathbb{R}[\mathbf{x}]$ denote the set of multivariate polynomials in $\mathbf{x} = (x_1, \dots, x_d)$, where we extend

$$\deg(f) := \max \{\deg(\mathbf{x}^\alpha) \mid f_\alpha \neq 0\}$$

for any element $f(\mathbf{x}) = \sum_{\alpha} f_{\alpha} \mathbf{x}^{\alpha}$ of $\mathbb{R}[\mathbf{x}]$.

The vector space of multivariate polynomials of degree at most t is

$$\mathbb{R}_t[\mathbf{x}] := \{f \in \mathbb{R}[\mathbf{x}] \mid \deg(f) \leq t\},$$

where $z_d(t) := \dim(\mathbb{R}_t[\mathbf{x}]) = \binom{d+t}{d}$. By defining

$$\mathbb{N}_t^d := \{\boldsymbol{\alpha} \in \mathbb{N}^d \mid \langle \boldsymbol{\alpha}, \mathbf{e} \rangle \leq t\},$$

we see that each polynomial $f \in \mathbb{R}_t[\mathbf{x}]$ can be written as $f(\mathbf{x}) = \sum_{\boldsymbol{\alpha} \in \mathbb{N}_t^d} f_{\boldsymbol{\alpha}} \mathbf{x}^{\boldsymbol{\alpha}}$ and we may identify $\mathbb{R}_t[\mathbf{x}]$ with $\mathbb{R}^{\mathbb{N}_t^d} \cong \mathbb{R}^{z_d(t)}$ by identifying the polynomial f with its vector of coefficients \mathbf{f} . In this context, we will also write $\mathbf{f} \in \mathbb{R}^{z_d(t)}$ and denote the canonical monomial base $(\mathbf{x}^{\boldsymbol{\alpha}})_{\boldsymbol{\alpha} \in \mathbb{N}_t^d}$ through the *Veronese map* $v_t(\mathbf{x})$, such that

$$f(\mathbf{x}) = \langle \mathbf{f}, v_t(\mathbf{x}) \rangle.$$

Given $\mathbf{y} = (y_{\boldsymbol{\alpha}})_{\boldsymbol{\alpha} \in \mathbb{N}^d} \in \mathbb{R}^{\mathbb{N}^d}$, we can use this identification to define the *Riesz functional* $L_{\mathbf{y}} : \mathbb{R}[\mathbf{x}] \rightarrow \mathbb{R}$ as

$$f \mapsto L_{\mathbf{y}}(f) = \langle \mathbf{f}, \mathbf{y} \rangle,$$

which can be understood as the linearization of a polynomial.

A subset $\mathcal{I} \subseteq \mathbb{R}[\mathbf{x}]$ is called an *ideal* if it is closed under addition ($\mathcal{I} + \mathcal{I} \subseteq \mathcal{I}$) and under multiplication of the whole multivariate polynomial ring ($\mathbb{R}[\mathbf{x}] \cdot \mathcal{I} \subseteq \mathcal{I}$). Any set of polynomials $\mathcal{H} \subseteq \mathbb{R}[\mathbf{x}]$ generates an ideal $\langle \mathcal{H} \rangle := \sum_{h \in \mathcal{H}} \mathbb{R}[\mathbf{x}] \cdot h$. An ideal \mathcal{I} is called *finitely generated* if it can be generated by a finite set $\mathcal{H} \subseteq \mathcal{I}$, in which case \mathcal{H} is called a *basis* of \mathcal{I} . By *Hilbert's basis theorem* [CLO15, Ch. 2, §5, Thm. 4], every ideal $\mathcal{I} \subseteq \mathbb{R}[\mathbf{x}]$ is finitely generated. The set

$$\mathcal{V}_{\mathbb{R}}(\mathcal{I}) := \{\mathbf{x} \in \mathbb{R}^d \mid h(\mathbf{x}) = 0 \forall h \in \mathcal{I}\}$$

is the set of common real zeros of all polynomials in \mathcal{I} and called the *real variety* of \mathcal{I} . Given a finite subset $\mathcal{G} \subseteq \mathbb{R}[\mathbf{x}]$, the set

$$\{\mathbf{x} \in \mathbb{R}^d \mid g(\mathbf{x}) \geq 0 \forall g \in \mathcal{G}\}$$

is called (*closed*) *basic semialgebraic*. In particular, real varieties are basic semialgebraic.

Gröbner bases

A *monomial order* \triangleleft is a strict well-order on \mathbb{N}^d that is translation invariant, or more formally, which implies for all choices of $\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma} \in \mathbb{N}^d$ that $\boldsymbol{\alpha} + \boldsymbol{\gamma} \triangleleft \boldsymbol{\beta} + \boldsymbol{\gamma}$ follows from $\boldsymbol{\alpha} \triangleleft \boldsymbol{\beta}$. Given a monomial order \triangleleft , we also write $\boldsymbol{\alpha} \preceq \boldsymbol{\beta}$ if either $\boldsymbol{\alpha} \triangleleft \boldsymbol{\beta}$ or $\boldsymbol{\alpha} = \boldsymbol{\beta}$.

Example 2.9:

The *lexicographic order* $<_{lex}$ on \mathbb{R}^d is defined as

$$\boldsymbol{\alpha} <_{lex} \boldsymbol{\beta} \Leftrightarrow \alpha_m < \beta_m \text{ for } m = \min \{i \in [d] \mid \alpha_i \neq \beta_i\}$$

and the *graded lexicographic order* $<_{grlex}$ on \mathbb{R}^d is defined as

$$\boldsymbol{\alpha} <_{grlex} \boldsymbol{\beta} \Leftrightarrow \langle \boldsymbol{\alpha}, \mathbf{e} \rangle < \langle \boldsymbol{\beta}, \mathbf{e} \rangle \text{ or } \langle \boldsymbol{\alpha}, \mathbf{e} \rangle = \langle \boldsymbol{\beta}, \mathbf{e} \rangle, \boldsymbol{\alpha} <_{lex} \boldsymbol{\beta}.$$

Both of these orders are monomial orders.

Let $p = \sum_{\alpha} p_{\alpha} \mathbf{x}^{\alpha} \in \mathbb{R}[\mathbf{x}] \setminus \{0\}$ and fix a monomial order \triangleleft to define its *leading term* $LT_{\triangleleft}(p) = p_{\alpha_*} \mathbf{x}^{\alpha_*}$, where α_* is the maximal α w.r.t. \triangleleft such that $p_{\alpha} \neq 0$. More generally, the leading terms $LT_{\triangleleft}(I)$ of an ideal $I \neq \{0\}$ are defined as

$$LT_{\triangleleft}(I) := \{LT_{\triangleleft}(p) \mid p \in I \setminus \{0\}\}.$$

A *Gröbner basis* (or standard basis) of an ideal I is any basis $\mathcal{G} = \{g_1, \dots, g_m\}$ of I with the property that $LT_{\triangleleft}(I) = \langle LT_{\triangleleft}(g_1), \dots, LT_{\triangleleft}(g_m) \rangle$. Additionally, \mathcal{G} is called *reduced* if, for each $i \in [m]$, $LT_{\triangleleft}(g_i)$ is a monomial and no monomial appearing in g_i lies in $\langle LT_{\triangleleft}(\mathcal{G} \setminus \{g_i\}) \rangle$. For every ideal $\{0\} \neq I \subseteq \mathbb{R}[\mathbf{x}]$ and every monomial order \triangleleft , there exists a unique reduced Gröbner basis.

Quotients of Polynomial Rings

Any ideal $I \subseteq \mathbb{R}[\mathbf{x}]$ defines an equivalence relation on $\mathbb{R}[\mathbf{x}]$ by setting

$$g \equiv f \pmod{I} \Leftrightarrow g - f \in I.$$

The set of its equivalence classes $\mathbb{R}[\mathbf{x}]/I = \{[g] \mid g \in \mathbb{R}[\mathbf{x}]\}$ is called the *quotient* of $\mathbb{R}[\mathbf{x}]$ modulo I , and $[g] = \{f \in \mathbb{R}[\mathbf{x}] \mid g \equiv f \pmod{I}\} = g + I$ is a ring itself with the operations $[f] + [g] = [f + g]$ and $[f] \cdot [g] = [f \cdot g]$.

Theorem 2.10 ([CLO15, Ch. 5, §3, Prop. 4]):

Let $I \subseteq \mathbb{R}[\mathbf{x}]$ be an ideal. Then

$$\mathbb{R}[\mathbf{x}]/I \cong \text{lin}(\{\mathbf{x}^{\alpha} \mid \mathbf{x}^{\alpha} \notin LT_{\triangleleft}(I)\}) \quad (2.4)$$

as \mathbb{R} -vector spaces with isomorphism

$$[f] \mapsto \bar{f}.$$

In particular, we can write $[f] = \bar{f} + I$ for all $f \in \mathbb{R}[\mathbf{x}]$.

Remark 2.11:

Given a reduced Gröbner basis for I and \triangleleft , the computation of the reduction map $[f] \mapsto \bar{f}$ from Theorem (2.10) can be carried out efficiently algorithmically, as outlined in [CLO15]. In particular, addition $[f] + [g]$ and multiplication $[f] \cdot [g]$ translate into $\bar{f} + \bar{g}$ and $\bar{f} \cdot \bar{g}$ respectively, where the latter is computationally much more involved than the former.

As a generalization of the degree constrained vector space $\mathbb{R}_t[\mathbf{x}]$, we also define

$$\mathbb{R}_t[\mathbf{x}]/I \cong \text{lin}(\{\mathbf{x}^{\alpha} \mid \mathbf{x}^{\alpha} \notin LT_{\triangleleft}(I), \langle \alpha, \mathbf{e} \rangle \leq t\}).$$

We close this section with examples for various ideals that will be encountered in the rest of the thesis.

Example 2.12 (Binary Vectors):

Let

$$\mathcal{I}(C^d) = \langle x_1^2 - x_1, \dots, x_d^2 - x_d \rangle$$

be the ideal whose variety are the vertices of the d -dimensional unit cube C^d , e.g. the set of d -dimensional binary vectors. A reduced Gröbner basis for $\mathcal{I}(C^d)$ and $<_{grlex}$ is given by $\{x_i^2 \mid i \in [d]\}$, and by Theorem 2.10, we get the \mathbb{R} -vector space isomorphisms

$$\mathbb{R}[x]/\mathcal{I}(C^d) \cong \text{lin} \left(\{x^\alpha \mid \alpha \in \{0, 1\}^d\} \right) \cong \text{lin} \left(\{x_I \mid I \in 2^{[d]}\} \right)$$

by using $(\{0, 1\}^d, \vee) \cong (2^{[d]}, \cup)$ via $\alpha \leftrightarrow \text{supp}(\alpha)$. Then multiplication in $\mathbb{R}[x]/\mathcal{I}(C^d)$ translates into the multiplication $x_I \cdot x_J = x_{I \cup J}$ for all $I, J \subseteq [d]$.

Example 2.13 (Independence Systems):

For an independence system $L \subseteq 2^{[d]}$, let

$$\mathcal{I}(L) := \langle x^\alpha \mid \alpha \in \mathbb{N}^d, \text{supp}(\alpha) \notin L \rangle$$

be its Stanley-Reisner ideal, whose variety is

$$\mathcal{V}_{\mathbb{R}}(L) := \mathcal{V}_{\mathbb{R}}(\mathcal{I}(L)) = \{x \in \mathbb{R}^d \mid \text{supp}(x) \in L\},$$

the union of linear subspaces of \mathbb{R}^d whose support is in L . A reduced Gröbner basis for $\mathcal{I}(L)$ and $<_{grlex}$ is given by

$$\{x^\alpha \mid \alpha \in \{0, 1\}^d, \text{supp}(\alpha) \text{ is minimally dependent in } L\}.$$

By Theorem 2.10, we get the \mathbb{R} -vector space isomorphism

$$\mathbb{R}[x]/\mathcal{I}(L) \cong \text{lin} \left(\{x^\alpha \mid \alpha \in \mathbb{N}^d, \text{supp}(\alpha) \in L\} \right)$$

and multiplication in $\mathbb{R}[x]/\mathcal{I}(L)$ becomes

$$x^\alpha \cdot x^\beta = \begin{cases} x^{\alpha+\beta} & \text{if } \text{supp}(\alpha + \beta) \in L, \\ 0 & \text{else.} \end{cases}$$

The sum of both ideals then yields the intersection of the corresponding varieties, the L -constrained cube given by

$$\mathcal{V}_{\mathbb{R}}(\mathcal{I}(C^d) + \mathcal{I}(L)) = \{0, 1\}^d \cap \mathcal{V}_{\mathbb{R}}(L) = \{x \in \{0, 1\}^d \mid \text{supp}(x) \in L\},$$

and it can be shown that

$$\mathbb{R}[x]/\mathcal{I}(L) \cong \text{lin} \left(\{x^\alpha \mid \alpha \in \{0, 1\}^d \cap \mathcal{V}_{\mathbb{R}}(L)\} \right) \cong \text{lin}(\{x_I \mid I \in L\}),$$

with multiplication

$$x_I \cdot x_J = \begin{cases} x_{I \cup J} & \text{if } I \cup J \in L, \\ 0 & \text{else.} \end{cases}$$

Example 2.14 (Affine Hyperplane):

Let $h = 1 - \langle \mathbf{x}, \mathbf{e} \rangle$ be the linear polynomial whose variety is the hyperplane in \mathbb{R}^d containing Δ^d and let $\mathcal{I}(\Delta^d) = \langle 1 - \langle \mathbf{x}, \mathbf{e} \rangle \rangle \subseteq \mathbb{R}[\mathbf{x}]$ be the corresponding ideal. For the monomial order $<_{grlex}$, a reduced Gröbner basis is h itself, and since $LT_{grlex}(h) = x_1$, it follows that $\mathbb{R}[\mathbf{x}]/\mathcal{I}(\Delta^d) \cong \mathbb{R}[x_2, \dots, x_d]$, which can be realized by just eliminating all occurrences of x_1 via $x_1 = 1 - x_2 - \dots - x_d$.

While $\mathbb{R}[\mathbf{x}]/\mathcal{I}(\Delta^d) \cong \mathbb{R}[x_2, \dots, x_d]$ is intuitive, it creates an artificial asymmetry between the otherwise symmetrical variables in \mathbf{x} . The following result can be used as symmetrical alternative instead.

Theorem 2.15:

We have $\mathbb{R}_t[\mathbf{x}]/\mathcal{I}(\Delta^d) \cong \mathbb{R}_t^h[\mathbf{x}] := \{f \in \mathbb{R}[\mathbf{x}] \mid \deg(f) = t\}$ via the homogenization map

$$f(\mathbf{x}) \mapsto f^{h_t}(\mathbf{x}) := \langle \mathbf{x}, \mathbf{e} \rangle^t \cdot f\left(\frac{\mathbf{x}}{\langle \mathbf{x}, \mathbf{e} \rangle}\right). \quad (2.5)$$

Proof. By definition, $f^{h_t} \equiv f \pmod{\mathcal{I}(\Delta^d)}$ for any $f \in \mathbb{R}_t[\mathbf{x}]$, so the map is well-defined, and all monomials occurring in f^{h_t} have degree t , so the image of $\mathbb{R}_t[\mathbf{x}]/\mathcal{I}(\Delta^d)$ is contained in $\mathbb{R}_t^h[\mathbf{x}]$. It suffices now to show that both vector spaces have the same dimension, and by Example 2.14,

$$\dim(\mathbb{R}_t[\mathbf{x}]/\mathcal{I}(\Delta^d)) = \dim(\mathbb{R}_t[x_2, \dots, x_d]) = z_{d-1}(t) = \binom{d-1+t}{t} = \dim(\mathbb{R}_t^h[\mathbf{x}]),$$

where the last equation follows from [Sta11, p. 18]. \square

2.6 Method of Moments

This section gives a basic description of *method of moments* (MM), and is based mostly on the book [Las15], with some minor changes of notation.

Moments

Moment matrices

For $t \in \mathbb{N}$ and $\mathbf{y} \in \mathbb{R}^{\mathbb{N}_{2t}^d}$, the matrix $\mathbf{M}_t(\mathbf{y})$ indexed by \mathbb{N}_t^d is defined entrywise by

$$(\mathbf{M}_t(\mathbf{y}))_{\alpha, \beta} := L_{\mathbf{y}}(\mathbf{x}^\alpha \cdot \mathbf{x}^\beta) = y_{\alpha+\beta} \quad \forall \alpha, \beta \in \mathbb{N}_t^d$$

and is called the *moment matrix* of order t of \mathbf{y} .

More generally, let $f \in \mathbb{R}[\mathbf{x}]$ be a multivariate polynomial and define

$$\text{td}(f) := \left\lfloor \frac{\deg(f)}{2} \right\rfloor.$$

Then for $t \geq \text{td}(f)$, the matrix $\mathbf{M}_t(f, \mathbf{y})$ indexed by $\mathbb{N}_{t-\text{td}(f)}^d$ and entrywise defined as

$$(\mathbf{M}_t(f, \mathbf{y}))_{\alpha, \beta} := L_{\mathbf{y}}(\mathbf{x}^{\alpha} \cdot \mathbf{x}^{\beta} \cdot f)$$

is called the *localizing moment matrix* of order t of \mathbf{y} with respect to f . Note that each entry of $\mathbf{M}_t(f, \mathbf{y})$ is a linear expression in \mathbf{y} and that we recover $\mathbf{M}_t(\mathbf{y}) = \mathbf{M}_t(1, \mathbf{y})$ as a special case.

Example 2.16:

Consider the bivariate polynomial ring $\mathbb{R}[\mathbf{x}]$ with $\mathbf{x} = (x_1, x_2)$ and the polynomials $h_i(\mathbf{x}) = x_i^2 - x_i$ for $i = 1, 2$. Then, for $t = 2$ we have

$$v_2(\mathbf{x})^{\top} = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$$

and consequently

$$v_2(\mathbf{x})v_2(\mathbf{x})^{\top} = \begin{pmatrix} 1 & x_1 & x_2 & x_1^2 & x_1x_2 & x_2^2 \\ x_1 & x_1^2 & x_1x_2 & x_1^3 & x_1^2x_2 & x_1x_2^2 \\ x_2 & x_1x_2 & x_2^2 & x_1^2x_2 & x_1x_2^2 & x_2^3 \\ x_1^2 & x_1^3 & x_1^2x_2 & x_1^4 & x_1^3x_2 & x_1^2x_2^2 \\ x_1x_2 & x_1^2x_2 & x_1x_2^2 & x_1^3x_2 & x_1^2x_2^2 & x_1x_2^3 \\ x_2^2 & x_1x_2^2 & x_2^3 & x_1^2x_2^2 & x_1x_2^3 & x_2^4 \end{pmatrix}.$$

Applying $L_{\mathbf{y}}$ componentwise on this matrix then shows

$$\mathbf{M}_2(\mathbf{y}) = \begin{pmatrix} \mathcal{Y}(0,0) & \mathcal{Y}(1,0) & \mathcal{Y}(0,1) & \mathcal{Y}(2,0) & \mathcal{Y}(1,1) & \mathcal{Y}(0,2) \\ \mathcal{Y}(1,0) & \mathcal{Y}(2,0) & \mathcal{Y}(1,1) & \mathcal{Y}(3,0) & \mathcal{Y}(2,1) & \mathcal{Y}(1,2) \\ \mathcal{Y}(0,1) & \mathcal{Y}(1,1) & \mathcal{Y}(0,2) & \mathcal{Y}(2,1) & \mathcal{Y}(1,2) & \mathcal{Y}(0,3) \\ \mathcal{Y}(2,0) & \mathcal{Y}(3,0) & \mathcal{Y}(2,1) & \mathcal{Y}(4,0) & \mathcal{Y}(3,1) & \mathcal{Y}(2,2) \\ \mathcal{Y}(1,1) & \mathcal{Y}(2,1) & \mathcal{Y}(1,2) & \mathcal{Y}(3,1) & \mathcal{Y}(2,2) & \mathcal{Y}(1,3) \\ \mathcal{Y}(0,2) & \mathcal{Y}(1,2) & \mathcal{Y}(0,3) & \mathcal{Y}(2,2) & \mathcal{Y}(1,3) & \mathcal{Y}(0,4) \end{pmatrix},$$

which has entries $\mathbf{y} \in \mathbb{R}^{\mathbb{N}_4^2} \cong \mathbb{R}^{15}$ and thus lies in a 15-dimensional subspace. Note that the blocks in both matrices contain entries corresponding to constant total degree. Since $\text{td}(h_i) = 1$, we also get the componentwise maps

$$L_{\mathbf{y}} : h_i(\mathbf{x})v_1(\mathbf{x})v_1(\mathbf{x})^{\top} \mapsto \mathbf{M}_2(h_i, \mathbf{y})$$

which is explicitly given for the case of $i = 1$ as

$$h_1(\mathbf{x})v_1(\mathbf{x})v_1(\mathbf{x})^{\top} = \begin{pmatrix} x_1^2 - x_1 & x_1^3 - x_1^2 & x_1^2x_2 - x_1x_2 \\ x_1^3 - x_1^2 & x_1^4 - x_1^3 & x_1^3x_2 - x_1^2x_2^2 \\ x_1^2x_2 - x_1x_2 & x_1^3x_2 - x_1^2x_2^2 & x_1^2x_2^2 - x_1x_2^3 \end{pmatrix}$$

and

$$\mathbf{M}_2(h_1, \mathbf{y}) = \begin{pmatrix} \mathcal{Y}(2,0) - \mathcal{Y}(1,0) & \mathcal{Y}(3,0) - \mathcal{Y}(2,0) & \mathcal{Y}(2,1) - \mathcal{Y}(1,1) \\ \mathcal{Y}(3,2) - \mathcal{Y}(2,0) & \mathcal{Y}(4,0) - \mathcal{Y}(3,0) & \mathcal{Y}(3,1) - \mathcal{Y}(2,1) \\ \mathcal{Y}(2,1) - \mathcal{Y}(1,1) & \mathcal{Y}(3,1) - \mathcal{Y}(2,1) & \mathcal{Y}(2,2) - \mathcal{Y}(1,2) \end{pmatrix}.$$

Measures and moments

Let $\mathcal{N}(C) \subseteq \mathbb{R}[\mathbf{x}]$ be the convex cone of polynomials that are nonnegative on $C \subseteq \mathbb{R}^d$ and denote its dual cone by

$$\mathcal{N}^*(C) := \left\{ \mathbf{y} \in \mathbb{R}^{\mathbb{N}^d} \mid L_{\mathbf{y}}(f) \geq 0 \ \forall f \in \mathcal{N}(C) \right\}.$$

Denote by $\mathcal{M}_+(C)$ the space of finite (nonnegative) Borel measures supported on C and by $\mathcal{P}(C)$ the subset of probability measures on C . We can recover the cone of the corresponding moments

$$\left\{ \mathbf{y} \in \mathbb{R}^{\mathbb{N}^d} \mid \exists \mu \in \mathcal{M}_+(C): y_{\alpha} = \int_C \mathbf{x}^{\alpha} \cdot d\mu \quad \forall \alpha \in \mathbb{N}^d \right\} \subseteq \mathcal{N}^*(C), \quad (2.6)$$

where equality holds if C is compact [Las15, Lemma 4.7].

The Method

Reformulation of Optimization Problems

Let $C \subseteq \mathbb{R}^d$ be a compact set and $f(\mathbf{x}) = \sum_{\alpha \in \mathbb{N}_t^d} f_{\alpha} \mathbf{x}^{\alpha}$ be a real-valued multivariate polynomial, then

$$\inf_{\mathbf{x} \in C} f(\mathbf{x}) = \inf_{\mu \in \mathcal{P}(C)} \int_C f \cdot d\mu \quad (2.7)$$

can be reduced to a convex linear programming problem. Indeed, we have that

$$\int_C f \cdot d\mu = \int_C \sum_{\alpha \in \mathbb{N}_t^d} f_{\alpha} \mathbf{x}^{\alpha} \cdot d\mu = \sum_{\alpha \in \mathbb{N}_t^d} f_{\alpha} \int_C \mathbf{x}^{\alpha} \cdot d\mu = L_{\mathbf{y}}(f),$$

where

$$y_{\alpha} = \int_C \mathbf{x}^{\alpha} \cdot d\mu$$

is the moment of order α .

Consequently,

$$\inf \{ L_{\mathbf{y}}(f) \mid y_0 = 1, \mathbf{y} \in \mathcal{N}^*(C) \} \quad (2.8)$$

is a relaxation of problem (2.7) with the benefit of being a reformulation whenever equality holds in (2.6).

Note that the constraint $y_0 = 1$ normalizes $\mathbf{y} \in \mathcal{N}^*(C)$ to represent a probability measure in $\mathcal{P}(C) \subseteq \mathcal{M}_+(C)$.

Although problem (2.8) is a convex linear programming problem, the characterization of $\mathbf{y} \in \mathcal{N}^*(C)$ (known as C -moment problem in the literature) may be notoriously hard for general C .

However, for compact *semi-algebraic* C given as

$$C = \{ \mathbf{x} \in \mathbb{R}^d \mid g(\mathbf{x}) \geq 0 \quad \forall g \in \mathcal{G} \} \quad (2.9)$$

for some *finite* set of polynomials $\mathcal{G} \subseteq \mathbb{R}[\mathbf{x}]$, an explicit characterization of $\mathcal{N}^*(C)$ is available. Since C is assumed to be compact, we will assume without loss of generality that

$$0 \leq R^2 - \|\mathbf{x}\|^2 \in \mathcal{G},$$

where R is a sufficiently large positive constant.

Remark 2.17:

In fact, we would only need any function u in the quadratic module generated by \mathcal{G} to have a compact superlevel set $\{ \mathbf{x} \in \mathbb{R}^d \mid u(\mathbf{x}) \leq 0 \}$ for the following.

This representation allows the application of a theorem on positivity by Putinar [Las15, Theorem 2.15], which leads to

$$\begin{aligned} \mathcal{N}^*(C) &= \left\{ \mathbf{y} \in \mathbb{R}^{\mathbb{N}^d} \mid \mathbf{M}_t(\mathbf{y}) \geq \mathbf{0}, \mathbf{M}_t(g, \mathbf{y}) \geq \mathbf{0} \quad \forall g \in \mathcal{G}, \forall t \in \mathbb{N} \right\} \\ &=: \mathcal{N}_{\geq}^*(\mathcal{G}). \end{aligned}$$

In particular, problem (2.8) is equivalent to

$$\inf \{ L_{\mathbf{y}}(f) \mid y_0 = 1, \mathbf{y} \in \mathcal{N}_{\geq}^*(\mathcal{G}) \}.$$

To summarize, if f is a polynomial and C a compact semi-algebraic set, then problem (2.7) is equivalent to a convex linear programming problem with an infinite number of linear constraints on an infinite number of decision variables.

Semidefinite Relaxations

Now, for $t \geq \text{td}(f)$, consider the finite-dimensional truncations

$$\rho_t = \inf \{ L_{\mathbf{y}}(f) \mid y_0 = 1, \mathbf{y} \in \mathcal{N}_t^*(\mathcal{G}) \} \quad (2.10)$$

of problem (2.8) where

$$\mathcal{N}_t^*(\mathcal{G}) := \left\{ \mathbf{y} \in \mathbb{R}^{\mathbb{N}_{2t}^d} \mid \begin{array}{l} \mathbf{M}_t(\mathbf{y}) \geq \mathbf{0}, \\ \mathbf{M}_t(g, \mathbf{y}) \geq \mathbf{0} \quad \forall g \in \mathcal{G} : t \geq \text{td}(g) \end{array} \right\}.$$

By construction, $\{ \mathcal{N}_t^*(\mathcal{G}) \}_{t \in \mathbb{N}}$ generates a hierarchy of relaxations for $\mathcal{N}^*(C)$ in problem (2.8), where each $\mathcal{N}_t^*(\mathcal{G})$ is concerned with moment and localizing matrices of fixed size t . The lowerbounds ρ_t monotonically converge toward the optimal value of (2.7) [Las15, Theorem 6.2] and finite convergence may take place, which can be efficiently checked [Las15, Theorem 6.6].

Furthermore, in the best case of finite convergence, (2.10) will yield the global optimal value and a convex combination of global optimal solutions as minimizer, which can be efficiently decomposed into optimal solutions [Las15, Sct. 6.1.2].

In the noncompact case, the ρ_t are still monotonically increasing lower bounds of (2.7), but convergence to the optimum is not guaranteed.

Remark 2.18:

In the literature, this construction is known as *Method of Moments* (MM), where it is assumed that $t \geq \max_i \text{td}(g_i)$ in addition to $t \geq \text{td}(f)$ in order to start with a complete description of all the constraints used in the problem. Our slightly different definition is more flexible by enabling us to start with an incomplete set of constraints of low degree while still fitting into the overall hierarchy.

It should be noted that using a value of t that truncates most of the 'relevant' inequalities for the problem is not likely to yield a useful lower bound.

For convenience, we will also introduce a shorthand notation for polynomial equations $h(\mathbf{x}) = 0$ (imposed by having both $h(\mathbf{x}) \geq 0$ and $-h(\mathbf{x}) \geq 0$) by setting

$$\mathcal{N}_t^*(\mathcal{H}, \mathcal{G}) := \left\{ \mathbf{y} \in \mathbb{R}^{\mathbb{N}_{2t}^d} \left| \begin{array}{l} \mathbf{M}_t(\mathbf{y}) \geq \mathbf{0} \\ \mathbf{M}_t(g, \mathbf{y}) \geq \mathbf{0} \quad \forall g \in \mathcal{G} : t \geq \text{td}(g) \\ \mathbf{M}_t(h, \mathbf{y}) = \mathbf{0} \quad \forall h \in \mathcal{H} : t \geq \text{td}(h) \end{array} \right. \right\}. \quad (2.11)$$

If the description of a basic semialgebraic set

$$C = \{ \mathbf{x} \in \mathbb{R}^d \mid h(\mathbf{x}) = 0 \quad \forall h \in \mathcal{H}, \quad g(\mathbf{x}) \geq 0 \quad \forall g \in \mathcal{G} \}$$

is clear from context, we will also abuse notation to write $\mathcal{N}_t^*(C)$ instead of $\mathcal{N}_t^*(\mathcal{H}, \mathcal{G})$.

Chapter 3

Computational Aspects

The goal of this chapter is to outline the computational aspects of MM, that is, how to implement and formulate the relaxations encountered in Section 2.6 in practice.

As a preprocessing step, we explain in Section 3.1 how the results from Section 2.5 can be exploited to reduce the size of the moment matrices when optimizing over a variety. Section 3.2 then describes the software and hardware that we used to carry out all experiments throughout the thesis. Lastly, Section 3.3 gives a new method to solve conic linear problems.

3.1 Method of Moments on a Variety

When using MM on a set K whose description includes polynomial equations, the size of the moment matrices involved can be reduced. To see this, let (2.9) be the intersection

$$C = \{ \mathbf{x} \in \mathbb{R}^d \mid g(\mathbf{x}) \geq 0 \quad \forall g \in \mathcal{G} \} \cap \mathcal{V}_{\mathbb{R}}(\langle \mathcal{H} \rangle)$$

of an basic semialgebraic set and a real variety. From an algebraic point of view, the constraints $\mathbf{M}_t(h, \mathbf{y}) = \mathbf{0}$ in (2.11) are not very natural, since they depend on a specific generating system of their underlying ideal $\mathcal{I} = \langle \mathcal{H} \rangle$ and neglect the overall structure of \mathcal{I} . In addition, we have the issue of redundancy in the form of inherent linear dependent columns in the moment matrices, as shown by the following lemma.

Lemma 3.1:

For $h, g \in \mathbb{R}[\mathbf{x}]$, let $t \in \mathbb{N}$ satisfy $t \geq \deg(h) + \text{td}(g)$. If $\mathbf{y} \in \mathbb{R}^{\mathbb{N}^d}$ satisfies $\mathbf{M}_{t+1}(h, \mathbf{y}) = \mathbf{0}$, then

$$\mathbf{M}_t(g, \mathbf{y}) \cdot \mathbf{h} = \mathbf{0}.$$

Proof. Let $h(\mathbf{x}) = \sum_{\alpha \in \mathbb{N}^d} h_{\alpha} \cdot \mathbf{x}^{\alpha}$ and $g(\mathbf{x}) = \sum_{\alpha \in \mathbb{N}^d} g_{\alpha} \cdot \mathbf{x}^{\alpha}$. Due to linearity of $L_{\mathbf{y}}$, for

each $\boldsymbol{y} \in \mathbb{N}_{t-\text{td}(g)}^d$ it holds that

$$\begin{aligned} (\mathbf{M}_t(g, \boldsymbol{y}) \cdot \mathbf{h})_{\boldsymbol{y}} &= \sum_{\boldsymbol{\alpha} \in \mathbb{N}_{\deg(h)}^d} h_{\boldsymbol{\alpha}} \cdot L_{\boldsymbol{y}}(\mathbf{x}^{\boldsymbol{\alpha}} \cdot \mathbf{x}^{\boldsymbol{y}} \cdot g) = L_{\boldsymbol{y}}(\mathbf{x}^{\boldsymbol{y}} \cdot g \cdot h) \\ &= \sum_{\boldsymbol{\beta} \in \mathbb{N}_{\deg(g)}^d} g_{\boldsymbol{\beta}} \cdot L_{\boldsymbol{y}}(\mathbf{x}^{\boldsymbol{\beta}} \cdot \mathbf{x}^{\boldsymbol{y}} \cdot h). \end{aligned} \quad (3.1)$$

Then the inequality $\deg(\mathbf{x}^{\boldsymbol{\beta}+\boldsymbol{y}}) \leq 2(t+1-\text{td}(h))$ is sharp, depending on the parities of $\deg(g)$ and $\deg(h)$. Consequently, $\mathbf{M}_{t+1}(h, \boldsymbol{y})$ is the smallest moment matrix that sets all entries $L_{\boldsymbol{y}}(\mathbf{x}^{\boldsymbol{\beta}} \cdot \mathbf{x}^{\boldsymbol{y}} \cdot h)$ occurring in (3.1) to 0. \square

Lemma 3.1 needs $\mathbf{M}_{t+1}(h, \boldsymbol{y}) = \mathbf{0}$ rather than $\mathbf{M}_t(h, \boldsymbol{y}) = \mathbf{0}$ when $\deg(g)$ is even while $\deg(h)$ is odd, as shown in the next example.

Example 3.2:

Let $d = 2$, $g(\mathbf{x}) = 1$ and $h(\mathbf{x}) = x_1^3 - x_2^2 \in \mathbb{R}[\mathbf{x}]$. For $t = 3$, we have

$$\mathbf{M}_2(h, \boldsymbol{y}) = \begin{pmatrix} \mathcal{Y}(3,0) - \mathcal{Y}(0,2) & \mathcal{Y}(4,0) - \mathcal{Y}(1,2) & \mathcal{Y}(3,1) - \mathcal{Y}(0,3) \\ \mathcal{Y}(4,0) - \mathcal{Y}(1,2) & \mathcal{Y}(5,0) - \mathcal{Y}(2,2) & \mathcal{Y}(4,1) - \mathcal{Y}(1,3) \\ \mathcal{Y}(3,1) - \mathcal{Y}(0,3) & \mathcal{Y}(4,1) - \mathcal{Y}(1,3) & \mathcal{Y}(3,2) - \mathcal{Y}(0,4) \end{pmatrix},$$

but

$$\mathbf{M}_2(g, \boldsymbol{y}) \cdot \mathbf{h} = \begin{pmatrix} \mathcal{Y}(3,0) - \mathcal{Y}(0,2) \\ \mathcal{Y}(4,0) - \mathcal{Y}(1,2) \\ \mathcal{Y}(3,1) - \mathcal{Y}(0,3) \\ \mathcal{Y}(5,0) - \mathcal{Y}(2,2) \\ \mathcal{Y}(4,1) - \mathcal{Y}(1,3) \\ \mathcal{Y}(3,2) - \mathcal{Y}(0,4) \\ \mathcal{Y}(6,0) - \mathcal{Y}(3,2) \\ \mathcal{Y}(5,1) - \mathcal{Y}(2,3) \\ \mathcal{Y}(4,2) - \mathcal{Y}(1,4) \\ \mathcal{Y}(3,3) - \mathcal{Y}(0,5) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \mathcal{Y}(6,0) - \mathcal{Y}(3,2) \\ \mathcal{Y}(5,1) - \mathcal{Y}(2,3) \\ \mathcal{Y}(4,2) - \mathcal{Y}(1,4) \\ \mathcal{Y}(3,3) - \mathcal{Y}(0,5) \end{pmatrix},$$

is not necessarily $\mathbf{0}$. Note that the non-zero entries all contain linearizations of polynomials containing monomials of degree $2t$.

In particular, we get the following corollary.

Corollary 3.3:

For $h, g \in \mathbb{R}[\mathbf{x}]$, let $t \in \mathbb{N}$ satisfy $t \geq \deg(h) + \text{td}(g)$. Then

$$\mathbf{M}_t(g, \boldsymbol{y}) \cdot \mathbf{h} = \mathbf{0}$$

is a valid constraint for $\mathcal{N}_t^*(\{h\}, \{g\})$ and can be used to reduce the size of $\mathbf{M}_t(g, \boldsymbol{y}) \geq \mathbf{0}$.

Proof. Lemma 3.1 shows that the constraint is valid, and Lemma 2.4 shows that we can simultaneously reduce one column and row of $\mathbf{M}_t(g, \mathbf{y})$ to $\mathbf{0}$, making it obsolete for the psd. constraint. \square

To remedy this, one can instead set up MM in the quotient ring $\mathbb{R}[\mathbf{x}]/\mathcal{I}$ as opposed to $\mathbb{R}[\mathbf{x}]$, working with the projections of the inequalities given by \mathcal{G} onto $\mathbb{R}[\mathbf{x}]/\mathcal{I}$. This approach results in the theory of theta bodies, as explained in [BPT13, Chapter 6]. There, it is shown that compared to the approach outlined so far, working in the quotient ring can yield tighter relaxations with smaller moment matrices, albeit at the cost of doing computations in $\mathbb{R}[\mathbf{x}]/\mathcal{I}$ instead of $\mathbb{R}[\mathbf{x}]$.

Definition 3.4:

Given an ideal $\mathcal{I} \subseteq \mathbb{R}[\mathbf{x}]$, $f \in \mathbb{R}[\mathbf{x}]$ and $t \geq \text{td}(f)$, the matrix $\mathbf{M}_t^{\mathcal{I}}(f, \mathbf{y})$ indexed by

$$\left\{ \boldsymbol{\alpha} \in \mathbb{N}_{t-\text{td}(f)}^d \mid \mathbf{x}^{\boldsymbol{\alpha}} \notin \text{LT}_{\text{grlex}}(\mathcal{I}) \right\}$$

and entrywise defined as

$$(\mathbf{M}_t^{\mathcal{I}}(f, \mathbf{y}))_{\boldsymbol{\alpha}, \boldsymbol{\beta}} := L_{\mathbf{y}} \left(\overline{\mathbf{x}^{\boldsymbol{\alpha}} \cdot \mathbf{x}^{\boldsymbol{\beta}} \cdot f} \right),$$

is called the *reduced localizing moment matrix* of order t of \mathbf{y} with respect to f .

Example 3.5:

Recall Example 2.16, where we considered moment matrices for $t = 2$ and the bivariate polynomials $h_i(\mathbf{x}) = x_i^2 - x_i$, which define a basis for $\mathcal{I}(C^2)$. Instead of using the constraints

$$\mathbf{M}_2(\mathbf{y}) \geq \mathbf{0}, \quad \mathbf{M}_2(h_1, \mathbf{y}) = \mathbf{0} \quad \text{and} \quad \mathbf{M}_2(h_2, \mathbf{y}) = \mathbf{0},$$

to describe C^2 , we can instead consider $v_2(\mathbf{x})v_2(\mathbf{x})^{\top} \bmod \mathcal{I}(C^2)$ to get

$$v_2(\mathbf{x})v_2(\mathbf{x})^{\top} \equiv \left(\begin{array}{c|cc|ccc} 1 & x_1 & x_2 & x_1 & x_1x_2 & x_2 \\ \hline x_1 & x_1 & x_1x_2 & x_1 & x_1x_2 & x_1x_2 \\ x_2 & x_1x_2 & x_2 & x_1x_2 & x_1x_2 & x_2 \\ \hline x_1 & x_1 & x_1x_2 & x_1 & x_1x_2 & x_1x_2 \\ x_1x_2 & x_1x_2 & x_1x_2 & x_1x_2 & x_1x_2 & x_1x_2 \\ x_2 & x_1x_2 & x_2 & x_1x_2 & x_1x_2 & x_2 \end{array} \right) \bmod \mathcal{I}(C^2). \quad (3.2)$$

As one can see, Corollary 3.3 shows that the coefficient vectors of h_1 and h_2 given by

$$\mathbf{h}_1 = (0, -1, 0, 1, 0, 0) \quad \text{and} \quad \mathbf{h}_2 = (0, 0, -1, 0, 0, 1)$$

belong to the kernel of the matrix-polynomial (3.2) independently of the evaluation of the variable \mathbf{x} . Consequently, using a basis of $\mathbb{R}_2[\mathbf{x}]/\mathcal{I}(C^2)$ removes the linear depen-

dent columns and rows arising this way and yields the reduced matrices

$$\left(\begin{array}{c|cc|c} 1 & x_1 & x_2 & x_1x_2 \\ \hline x_1 & x_1 & x_1x_2 & x_1x_2 \\ x_2 & x_1x_2 & x_2 & x_1x_2 \\ \hline x_1x_2 & x_1x_2 & x_1x_2 & x_1x_2 \end{array} \right) \mapsto \mathbf{M}_2^{\mathcal{I}(C^2)}(\mathbf{y}) = \left(\begin{array}{c|cc|c} \mathcal{Y}(0,0) & \mathcal{Y}(1,0) & \mathcal{Y}(0,1) & \mathcal{Y}(1,1) \\ \hline \mathcal{Y}(1,0) & \mathcal{Y}(1,0) & \mathcal{Y}(1,1) & \mathcal{Y}(1,1) \\ \mathcal{Y}(0,1) & \mathcal{Y}(1,1) & \mathcal{Y}(0,1) & \mathcal{Y}(1,1) \\ \hline \mathcal{Y}(1,1) & \mathcal{Y}(1,1) & \mathcal{Y}(1,1) & \mathcal{Y}(1,1) \end{array} \right),$$

where the latter only needs 4 variables out of the total 15 variables contained in $\mathbf{y} \in \mathbb{R}^{\mathbb{N}_4^2}$ used in the original approach.

This motivates the usage of the following theorem.

Theorem 3.6 (Moment Matrix Reduction):

Let $\mathcal{H}, \mathcal{G} \subseteq \mathbb{R}[\mathbf{x}]$, $\langle \mathcal{H} \rangle = \mathcal{I}$ and $t \geq \max \{ \text{td}(h) \mid h \in \mathcal{H} \}$, then for

$$\overline{\mathcal{N}}_t^*(\mathcal{H}, \mathcal{G}) := \left\{ \mathbf{y} \in \mathbb{R}^{\mathbb{N}_{2t}^d} \mid \begin{array}{l} \mathbf{M}_t^{\mathcal{I}}(\mathbf{y}) \geq \mathbf{0} \\ \mathbf{M}_t^{\mathcal{I}}(g, \mathbf{y}) \geq \mathbf{0} \quad \forall g \in \mathcal{G} : t \geq \text{td}(g) \\ \mathbf{M}_t(h, \mathbf{y}) = \mathbf{0} \quad \forall h \in \mathcal{H} \end{array} \right\}$$

we get the inclusion

$$\mathcal{N}^*(\mathcal{H}, \mathcal{G}) \subseteq \overline{\mathcal{N}}_t^*(\mathcal{H}, \mathcal{G}) \subseteq \mathcal{N}_t^*(\mathcal{H}, \mathcal{G}).$$

In particular, the constraints $\mathbf{M}_t(h, \mathbf{y}) = \mathbf{0}$ are only necessary to compute the entries of $\mathbf{y} \in \mathcal{N}_t^*(\mathcal{H}, \mathcal{G})$ that don't appear in the matrix $\mathbf{M}_t^{\mathcal{I}}(\mathbf{y})$.

The reason why $\mathbb{R}[\mathbf{x}]/\mathcal{I}$ and its degree-truncations admit smaller moment matrices in Theorem 3.6 is the reduction given by Theorem 2.10. Replacing the standard monomial basis by the basis given in Theorem 2.10 removes all linear dependencies that are implied by polynomials in \mathcal{I} through Corollary 3.3. However, depending on \mathcal{I} , the computation of a Gröbner basis may not be available in practice. Fortunately, the ideals shown in Examples 2.12, 2.13 and 2.14 already cover the most important cases we will examine. In particular, reduced moment matrices using the basis of $\mathbb{R}[\mathbf{x}]/\mathcal{I}(C^d)$ in Example 2.12 are called *combinatorial moment matrices* and have long been studied in the context of combinatorial optimization [Lau03].

In the case that a Gröbner basis is not available, we can still use Corollary 3.3 iteratively as a preprocessing step to reduce the size of $\mathbf{M}_t(g, \mathbf{y})$ significantly to the point where the columns do not contain the support of any entry in $\mathbf{M}_t(h, \mathbf{y})$.

3.2 Conic Linear Programming

The task of optimizing over the sets $\mathcal{N}_t^*(\mathcal{H}, \mathcal{G})$ or $\overline{\mathcal{N}}_t^*(\mathcal{H}, \mathcal{G})$ belongs to the problem class of *conic linear problems* (or *CLP*), which all can be transformed into the standard form

$$\min \{ \langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{Ax} = \mathbf{b}, \mathbf{x} \in \mathcal{K} \}, \quad (3.3)$$

for some convex cone $\mathcal{K} \subseteq \mathbb{R}^d$.

We will assume that the reader is familiar with the basic concepts of (conic) linear programming and suggest the book [AL12] for an in-depth treatment of both theory and practice. Currently, the field of conic linear programming is dominated by interior point methods, and this class of methods has seen numerous improvements due to being studied extensively in recent years [NT08]. Unfortunately, scaling becomes a major problem with these methods, making them prohibitive for large problems.

To fill this gap, another line of research is concerned with well-scaling first-order methods [EZC10]. For example, [OCPB16] recently proposed a combination of operator splitting and homogeneous self-dual embedding to tackle this problem and was able to beat state-of-the-art interior point methods on large instances. However, while first-order methods may perform faster in general, they usually do so at the expense of accuracy.

Setup

In order to solve any conic linear problems arising in this thesis, we utilized the widely used package SDPT3 [TTT96, TTT03] for MATLAB version 8.1.0.604 (R2013) with an Intel i5 of 3.2 GHz \times 4 and 16 GB of memory. In particular, the package can solve problems given in the standard form (3.3) when $\mathcal{K} \subseteq \mathbb{R}^d$ is a convex cone given as the Cartesian product of any combination of the cones \mathbb{R}^n , \mathbb{R}_+^n , \mathcal{S}_+^n and \mathcal{L}_n in various dimensions n , which is true for all problems that we will encounter.

Remark 3.7:

In this context, the psd. matrix cone \mathcal{S}_+^n is identified with the vector space of its upper triangular entries in $\mathbb{R}^{\binom{n}{2}}$ and thus vectorized.

The package uses an infeasible primal-dual path-following algorithm and thus belongs to the class of interior point methods. While the algorithm can be considered as an oracle for the sake of this thesis, we encourage the interested reader to read up on the underlying ideas in [TTT03].

3.3 The LP-Newton Method for CLPs

As an alternative to the interior point method implemented in SDPT3, and to give a self contained treatment of how to solve (3.3), we introduce the *conic LP-Newton method* (or CLP-Newton method), which is a generalization of the *LP-Newton method* devised in [FHYZ08].

The original method computes the end-point of the intersection of a line and a zonotope, and L and P stand for *line* and (convex) *polyhedron*, respectively. In particular, it was shown that for the case of $\mathcal{K} = \mathbb{R}_+^n$, solving (3.3) can be recast as finding the

endpoint of such an intersection under mild assumptions, and we will show that the same is true for any self-dual cone \mathcal{K} .

While interior point methods either start in or eventually enter the feasible region and then traverse it towards an optimal solution, the CLP-Newton-method starts with an infeasible point and converges towards an optimal feasible solution exclusively from outside of the feasible set. Although both approaches utilize the Newton method in some form, they function very differently, and the CLP-Newton method is conceptually closer to algorithms designed for feasibility problems.

Overview

This section is divided into two parts.

In the first part, we restate the results from [FHYZ08] in the setting of CLPs. For this, we introduce conic zonotopes and construct directly the adapted CLP-Newton method for them in Section (3.3.1). Since the minimum-norm-point algorithm is necessary for the projection step of the algorithm, it is recalled in Section (3.3.2).

In the second part, we show how to make the method explicit for widespread convex cones and evaluate the method based on experiments. In particular, Section (3.3.3) contains conditions for \mathcal{K} to be exploited by the algorithm and considers the cones \mathbb{R}_+^n , \mathcal{L}_n and \mathcal{S}_+^n , which are all covered by SDPT3 as well. Finally, Section (3.3.4) reports experiments on \mathcal{L}_n and \mathcal{S}_+^n , which shows how the algorithm performs in cases that were not considered in the original paper [FHYZ08].

3.3.1 The CLP-Newton Method

Conic Zonotopes

Throughout this section, let $\mathcal{K} \subseteq \mathbb{R}^n$ be a proper self-dual cone as in Section 2.4. Recall that $\leq_{\mathcal{K}}$ defines a partial order on \mathbb{R}^n by defining for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ that

$$\mathbf{y} \leq_{\mathcal{K}} \mathbf{x} \quad \Leftrightarrow \quad \mathbf{x} - \mathbf{y} \in \mathcal{K},$$

where $\mathbf{0} \leq_{\mathcal{K}} \mathbf{x}$ is short for membership $\mathbf{x} \in \mathcal{K}$. Extending this concept, we also define

$$\mathbf{y} <_{\mathcal{K}} \mathbf{x} \quad \Leftrightarrow \quad \mathbf{x} - \mathbf{y} \in \mathcal{K}, \quad \mathbf{x} \neq \mathbf{y}.$$

For any two points $\mathbf{l} \leq_{\mathcal{K}} \mathbf{u} \in \mathbb{R}^n$, we define their \mathcal{K} -box as

$$[\mathbf{l}, \mathbf{u}]_{\mathcal{K}} := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{l} \leq_{\mathcal{K}} \mathbf{x} \leq_{\mathcal{K}} \mathbf{u}\}.$$

Throughout this section, we will assume that $\mathbf{l} <_{\mathcal{K}} \mathbf{u}$, which implies the inclusion of the line segment between \mathbf{l} and \mathbf{u} in their \mathcal{K} -box, showing that

$$\text{conv}(\{\mathbf{l}, \mathbf{u}\}) \subseteq [\mathbf{l}, \mathbf{u}]_{\mathcal{K}}$$

is nonempty and nontrivial. Now let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and define

$$Z = \{z \mid z = \mathbf{A}\mathbf{x}, \mathbf{x} \in [\mathbf{l}, \mathbf{u}]_{\mathcal{K}}\},$$

which we will call a \mathcal{K} -zonotope in \mathbb{R}^m . This generalizes the already established zonotopes [Zie95], which correspond to \mathbb{R}_+^n -zonotopes in our setting.

Remark 3.8:

Linear optimization over Z is equivalent to linear optimization over $[\mathbf{l}, \mathbf{u}]_{\mathcal{K}}$, as

$$\max \{\langle \mathbf{c}, z \rangle \mid z \in Z\} = \max \{\langle \mathbf{A}^\top \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{x} \in [\mathbf{l}, \mathbf{u}]_{\mathcal{K}}\}.$$

CLP reformulation

Consider the CLP

$$\max \{\langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \in [\mathbf{l}, \mathbf{u}]_{\mathcal{K}}\}. \quad (\text{Box-CLP})$$

This is slightly different from the standard form (3.3) for CLPs, but can be cast into the form of (Box-CLP) given appropriate bounds on the feasible region, which are often available or can be easily computed.

Following [FHYZ08], we define an $(m + 1) \times n$ matrix

$$\bar{\mathbf{A}} = \begin{pmatrix} \mathbf{A} \\ \mathbf{c}^\top \end{pmatrix}$$

as well as a \mathcal{K} -zonotope

$$\bar{Z} = \{z \mid z = \bar{\mathbf{A}}\mathbf{x}, \mathbf{x} \in [\mathbf{l}, \mathbf{u}]_{\mathcal{K}}\}$$

and the line

$$LL = \left\{ \begin{pmatrix} \mathbf{b} \\ \gamma \end{pmatrix} \mid \gamma \in \mathbb{R} \right\}.$$

Using this notation, problem (Box-CLP) can be restated as

$$\gamma^* = \max \left\{ \gamma \mid \begin{pmatrix} z \\ \gamma \end{pmatrix} \in LL \cap \bar{Z} \right\}, \quad (\text{Box-CLP}')$$

where $\gamma \in \mathbb{R}$. The rest of this section is dedicated to solving (Box-CLP').

The Method

In order to describe the CLP-Newton method, we need to introduce additional definitions. To start, let

$$\gamma_0 = \max \{ \langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{x} \in [l, \mathbf{u}]_{\mathcal{K}} \}, \quad (3.4)$$

which is a relaxation of (Box-CLP') and as such yields an upperbound of γ^* . We will also need the orthogonal projection π_C on a given convex set $C \subseteq \mathbb{R}^n$ from Section 2.4. Lastly, we introduce the shorthand

$$\bar{\mathbf{b}}(\gamma) = \begin{pmatrix} \mathbf{b} \\ \gamma \end{pmatrix} \quad \forall \gamma \in \mathbb{R}$$

to parametrize LL and simplify notation. The following observation is the basis for the CLP-Newton method.

Consider the continuous, convex scalar function $g : \mathbb{R} \rightarrow \mathbb{R}$ given by

$$\gamma \mapsto g(\gamma) := \|\bar{\mathbf{b}}(\gamma) - \pi_{\bar{Z}}(\bar{\mathbf{b}}(\gamma))\|_2 = \text{dist}(\bar{\mathbf{b}}(\gamma), \bar{Z}).$$

Then $g(\gamma) = 0$ if and only if $\bar{\mathbf{b}}(\gamma) \in \bar{Z}$, and since $\bar{\mathbf{b}}(\gamma)$ parametrizes the line LL , the set of zeros of g is a parametrization of the set $LL \cap \bar{Z}$ via $\gamma \mapsto \bar{\mathbf{b}}(\gamma)$. Moreover, it follows by definition that the zeros of g also coincide with all possible values of the objective function of (Box-CLP'), which necessarily form a closed interval. This leads us to the following result.

Lemma 3.9:

The optimal value γ^* of (Box-CLP') is equal to

$$\max \{ \gamma \in \mathbb{R} \mid g(\gamma) = 0 \}. \quad (\text{MZ})$$

Now, applying the CLP-Newton method to solve (Box-CLP') consists of solving (MZ) with a generalized Newton method instead. The idea is that since g is convex and γ_0 is an upperbound to γ^* , we can start such a method at γ_0 and are guaranteed to converge towards γ^* from above. In general, g will be non-differentiable, and we need to use the subdifferential instead of the usual derivative. For this, the following lemma shows that we can extract a point in the subdifferential of g from a projection onto \bar{Z} .

Lemma 3.10:

For all $\gamma \in \mathbb{R}$ we have

$$\frac{\gamma - \langle \mathbf{e}_{m+1}, \pi_{\bar{Z}}(\bar{\mathbf{b}}(\gamma)) \rangle}{\text{dist}(\bar{\mathbf{b}}(\gamma), \bar{Z})} \in \partial g(\gamma).$$

Proof. The function $g(\gamma) = \text{dist}(\bar{\mathbf{b}}(\gamma), \bar{Z})$ is the composition of an affine map

$$\bar{\mathbf{b}}(\gamma) = \begin{pmatrix} \mathbf{0} \\ 1 \end{pmatrix} \gamma + \begin{pmatrix} \mathbf{b} \\ 0 \end{pmatrix}$$

and a distance function, so we can use the chain rule for subderivatives [Roc09, Ch. A] to yield

$$\partial g(\gamma) = \mathbf{e}_{m+1}^\top \cdot \partial_{\mathbf{x}=\bar{\mathbf{b}}(\gamma)} \text{dist}(\mathbf{x}, \bar{\mathbf{Z}}).$$

By Lemma 2.2, we have

$$\frac{\mathbf{x} - \pi_{\bar{\mathbf{Z}}}(\mathbf{x})}{\text{dist}(\mathbf{x}, \bar{\mathbf{Z}})} \in \partial \text{dist}(\mathbf{x}, \bar{\mathbf{Z}}),$$

and thus

$$\frac{\gamma - \langle \mathbf{e}_{m+1}, \pi_{\bar{\mathbf{Z}}}(\bar{\mathbf{b}}(\gamma)) \rangle}{\text{dist}(\bar{\mathbf{b}}(\gamma), \bar{\mathbf{Z}})} = \frac{\langle \mathbf{e}_{m+1}, \bar{\mathbf{b}}(\gamma) - \pi_{\bar{\mathbf{Z}}}(\bar{\mathbf{b}}(\gamma)) \rangle}{\text{dist}(\bar{\mathbf{b}}(\gamma), \bar{\mathbf{Z}})} \in \partial g(\gamma).$$

□

Now that we can extract a point in the subdifferential from a projection onto $\bar{\mathbf{Z}}$, we are able to state the following algorithm.

Algorithm 3.1: The CLP-Newton Method (CLPN)

Data: Data $\mathbf{A}, \mathbf{b}, \mathbf{c}, \mathbf{l}, \mathbf{u}$ for (Box-CLP'), error tolerance $\varepsilon > 0$.

Result: An approximate solution \mathbf{x}_k for (Box-CLP') or the detection that (Box-CLP') is infeasible.

- 1 Compute $\gamma_0 = \max \{ \langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{x} \in [\mathbf{l}, \mathbf{u}]_{\mathcal{K}} \}$;
 - 2 **for** $k = 1, 2, \dots$ **do**
 - 3 Find \mathbf{x}_k such that $\bar{\mathbf{A}}\mathbf{x}_k = \pi_{\bar{\mathbf{Z}}}(\bar{\mathbf{b}}(\gamma_{k-1}))$;
 - 4 **if** $\| \bar{\mathbf{A}}\mathbf{x}_k - \bar{\mathbf{b}}(\gamma_{k-1}) \|_2 < \varepsilon$ **then**
 - 5 **return** \mathbf{x}_k ;
 - 6 Set $(\mathbf{z}_k^\top, \zeta_k)^\top = \bar{\mathbf{A}}\mathbf{x}_k$;
 - 7 **if** $\zeta_k \geq \gamma_{k-1}$ **then**
 - 8 **return** "Problem (Box-CLP') is infeasible";
 - 9 Compute $\gamma_k = \zeta_k - \| \mathbf{b} - \mathbf{z}_k \|^2 \cdot (\gamma_{k-1} - \zeta_k)^{-1}$;
-

After initializing γ_0 as the starting value, the k -th step of Algorithm 3.1 checks for stopping criteria and then performs a Newton step on the current iterate γ_{k-1} as follows.

First, $g(\gamma_{k-1})$ is computed, and if $g(\gamma_{k-1}) \in [0, \varepsilon)$ for a given precision ε , then γ_{k-1} is accepted as a zero of g , thus terminating the algorithm.

In the next part, it is checked whether the new iterate passed the minimum of g , indicating that g has no zeros, thus proving the problem is infeasible.

If neither of these conditions is fulfilled, a new iterate γ_k can be computed by performing a Newton step. To do this, we choose $h_k \in \partial g(\gamma_{k-1})$ as in Lemma 3.10 in the recursion

$$\gamma_k = \gamma_{k-1} - \frac{g(\gamma_{k-1})}{h_k}$$

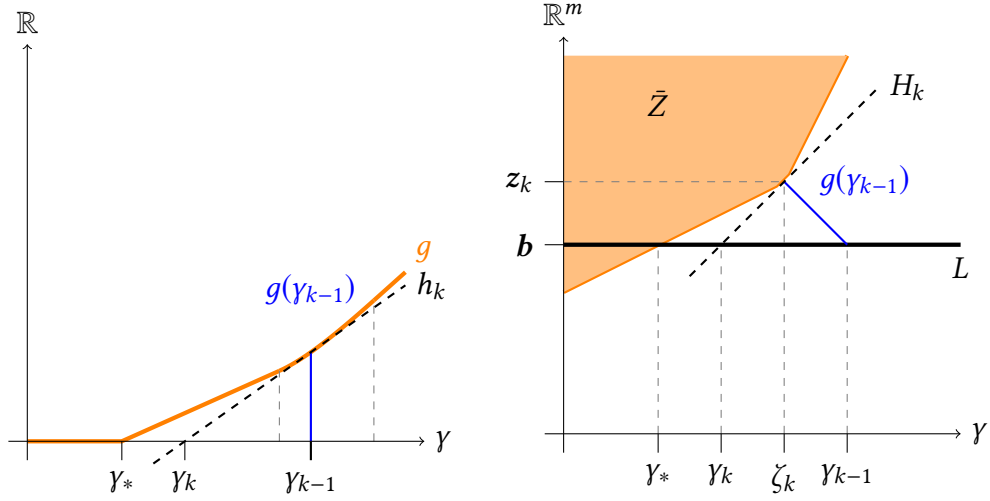


Figure 3.1: Visualization of the update step (3.5). Left: Newton step using g . Right: Geometric deduction from supporting hyperplane H_k .

to get

$$\gamma_k = \gamma_{k-1} - \frac{\text{dist}(\bar{\mathbf{b}}(\gamma_{k-1}), \bar{Z})^2}{\gamma_{k-1} - \zeta_k} = \zeta_k - \frac{\|\mathbf{b} - \mathbf{z}_k\|_2^2}{\gamma_{k-1} - \zeta_k}. \quad (3.5)$$

We can also understand the Newton step geometrically by first noting that

$$H_k = \left\{ \begin{pmatrix} \mathbf{z} \\ \zeta \end{pmatrix} \in \mathbb{R}^{m+1} \mid \left\langle \begin{pmatrix} \mathbf{z} \\ \zeta \end{pmatrix} - \begin{pmatrix} \mathbf{z}_k \\ \zeta_k \end{pmatrix}, \begin{pmatrix} \mathbf{b} \\ \gamma_{k-1} \end{pmatrix} - \begin{pmatrix} \mathbf{z}_k \\ \zeta_k \end{pmatrix} \right\rangle = 0 \right\}$$

is a supporting hyperplane of \bar{Z} at $\pi_{\bar{Z}}(\bar{\mathbf{b}}(\gamma_{k-1}))$. Then any feasible point in $LL \cap \bar{Z}$ is necessarily contained in the halfspace defined by H_k not containing $\bar{\mathbf{b}}(\gamma_{k-1})$, and (3.5) computes the intersection $LL \cap H_k$. Figure 3.1 compares the two approaches to deduce (3.5).

The following was shown in the original paper [FHYZ08] and generalizes to our setting.

Lemma 3.11 ([FHYZ08]):

The following statements hold for all values of k attained in Algorithm 3.1.

- (i) If $\gamma_{k-1} > \zeta_k$, then $\zeta_k > \gamma_k$.
- (ii) If $\gamma_{k-1} < \zeta_k$ or $\gamma_{k-1} = \zeta_k$ and $\mathbf{z}_k \neq \mathbf{b}$, then CLPN correctly assesses infeasibility of (Box-CLP').
- (iii) If $\gamma_{k-1} = \zeta_k$ and $\mathbf{z}_k = \mathbf{b}$, then γ_{k-1} is equal to the optimal value of (Box-CLP').

Remark 3.12:

Since g is convex, Algorithm 3.1 falls into the class of generalized Newton methods, which immediately shows asymptotic convergence in case that (Box-CLP') is feasible and finite termination in case that there is no feasible solution.

3.3.2 The Minimum-Norm-Point Algorithm

Algorithm 3.1 is limited by its routines to compute γ_0 and evaluate $\pi_{\bar{Z}}$, so both operations need to be efficient for practical usage. Fortunately, we can use Algorithm 3.2, an adaption of the minimum-norm-point algorithm taken from [Bac11, Jag13], to reduce the computation of $\pi_{\bar{Z}}$ to a series of conic LPs of the form (3.4).

Algorithm 3.2: Minimum-Norm-Point Algorithm (MNP) for $\text{dist}(y, \bar{Z})$

Data: Data l, u, \bar{A} for (Box-CLP'), $\bar{b} \in \mathbb{R}^{m+1}$, $s_0 \in [l, u]_{\mathcal{K}}$, error tolerance ε .

Result: $\hat{x} \in [l, u]_{\mathcal{K}}$ such that $0 \leq \|\bar{A}\hat{x} - \bar{b}\|_2 - \text{dist}(\bar{b}, \bar{Z}) \leq \varepsilon$.

```

1 Set  $P = \{s_0\}$  and  $k = 0$ ;
2 for  $k = 1, 2, \dots$  do
3   Compute  $x_k = \text{argmin} \{ \|\bar{A}x - \bar{b}\|_2^2 \mid x \in \text{aff}(P) \}$ ;
4   if  $x_k \in \text{conv}(P)$  then
5     Compute  $s_k = \text{argmin} \{ \langle s, \bar{A}^\top(\bar{A}x_k - \bar{b}) \rangle \mid s \in [l, u]_{\mathcal{K}} \}$ ;
6     if  $2\langle \bar{A}(x_k - s_k), \bar{A}x_k - \bar{b} \rangle < \varepsilon$  then
7       return  $x_k$ ;
8     else
9       Set  $P = P \cup \{s_k\}$ ;
10  else
11    Compute  $\hat{\mu} = \max \{ \mu \mid \mu \in [0, 1], x_{k-1} + \mu(x_k - x_{k-1}) \in \bar{Z} \}$ ;
12    Set  $x_k = x_{k-1} + \hat{\mu}(x_k - x_{k-1})$ ;
13    Set  $P$  to the minimal subset  $P' \subseteq P$  such that  $x_k \in \text{conv}(P')$ ;

```

While finite termination of MNP is established by the following theorem, it is an open problem whether the runtime of MNP can be bounded by a polynomial in the input size.

Theorem 3.13 ([Jag13]):

Algorithm 3.2 produces a sequence $\{x_k\}_{k \in \mathbb{N}}$ such that for $k > 1$ and

$$h_k := 2\langle \bar{A}(x_{k-1} - s_k), \bar{A}x_{k-1} - y \rangle,$$

we get

$$0 \leq \|\bar{A}x_k - y\|_2 - \text{dist}(y, \bar{Z}) \leq h_{k+1} \leq \frac{27 \text{diam}(\bar{Z})^2}{4(k+2)}.$$

In particular, the algorithm works correctly and terminates after a finite number of steps.

Remark 3.14:

In [Jag13], it is shown that the preceding theorem also applies to other variants of the Frank-Wolfe algorithm that can approximate $\pi_{\mathcal{Z}}(\bar{\mathbf{b}}(y_k))$. However, preliminary experiments have shown that among the variants we tested, the minimum-norm-point algorithm was the fastest. For another discussion of MNP, consider [Bac11, Sct. 9.2].

In the following, we will formalize a suitable class of problems for MNP, based on the difficulty of solving the related conic linear optimization problem of the form (3.4).

Definition 3.15:

$\mathcal{K} \subseteq \mathbb{R}^n$ is called *suitable* for CLPN if the problem

$$\max \{ \langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{x} \in [\mathbf{l}, \mathbf{u}]_{\mathcal{K}} \subseteq \mathbb{R}^n \} \quad (3.6)$$

can be solved in time $\mathcal{O}(n^2)$.

This property ensures that the bottleneck of computing MNP is the evaluation of projection $\mathbf{x}_k = \pi_{\text{aff}(P)}(\bar{\mathbf{b}})$, corresponding to solving a linear system of dimension at most $n \times n$. At the same time, it guarantees an efficient way to compute the starting point y_0 .

Remark 3.16:

The set of cones suitable for CLPN is closed under taking Cartesian products. In particular, if $\mathcal{K} = \otimes_{i \in I} \mathcal{K}_i$, then the corresponding \mathcal{K} -zonotope can be decomposed into several smaller \mathcal{K}_i -zonotopes, making (3.6) solvable in parallel for each \mathcal{K}_i and potentially accelerating the algorithm by a huge margin.

3.3.3 Linear Optimization on \mathcal{K} -Zonotopes

In the following, we will investigate the structure of (3.6) to get a better understanding when \mathcal{K} is suitable for CLPN, and apply these insights on three well-studied cones from the literature.

Necessary conditions

Since (3.6) involves linear optimization over a convex set, the optimal solutions will necessarily be extreme points, as summarized in the following result.

Lemma 3.17 (Extreme points of $[\mathbf{l}, \mathbf{u}]_{\mathcal{K}}$):

Problem (3.6) is equivalent to

$$\max \{ \langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{x} \in \{ \mathbf{l}, \mathbf{u} \} \cup (\text{bd}(\mathbf{l} + \mathcal{K}) \cap \text{bd}(\mathbf{u} - \mathcal{K})) \}.$$

Proof. The set of extreme points of $[\mathbf{l}, \mathbf{u}]_{\mathcal{K}}$ is

$$\text{bd}([\mathbf{l}, \mathbf{u}]_{\mathcal{K}}) = (\text{bd}(\mathbf{l} + \mathcal{K}) \cap (\mathbf{u} - \mathcal{K})) \cup ((\mathbf{l} + \mathcal{K}) \cap \text{bd}(\mathbf{u} - \mathcal{K})).$$

Without loss of generality, we can assume for an optimal solution \mathbf{x}^* of (3.6) that

$$\mathbf{x}^* \in \text{bd}(\mathbf{l} + \mathcal{K}) \cap \text{int}(\mathbf{u} - \mathcal{K}) \setminus \{\mathbf{l}\},$$

so that we can write $\mathbf{x}^* = \mathbf{l} + \mathbf{y}$ with $\mathbf{y} \in \mathcal{K} \setminus \{\mathbf{0}\}$. Then for small $\varepsilon > 0$, we maintain

$$\mathbf{l} + (1 \pm \varepsilon)\mathbf{y} \in \text{bd}(\mathbf{l} + \mathcal{K}) \cap \text{int}(\mathbf{u} - \mathcal{K}) \setminus \{\mathbf{l}\},$$

and the optimality of \mathbf{x}^* implies $\langle \mathbf{c}, \mathbf{y} \rangle = 0$. But then $\langle \mathbf{c}, \mathbf{x}^* \rangle = \langle \mathbf{c}, \mathbf{l} \rangle$ and we can choose \mathbf{l} as maximizer instead. \square

The related dual problem of (3.6) also allows us some insights.

Lemma 3.18:

Let $\mathbf{c} = \mathbf{c}_+ + \mathbf{c}_-$ be the Moreau decomposition where

$$\mathbf{c}_+ = \pi_{\mathcal{K}}(\mathbf{c}) \quad \text{and} \quad \mathbf{c}_- = \pi_{-\mathcal{K}^*}(\mathbf{c}) = -\pi_{\mathcal{K}}(-\mathbf{c}),$$

since \mathcal{K} is self-dual. Then the dual of (3.6) is equivalent to

$$\min \{ \langle \mathbf{l} - \mathbf{u}, \mathbf{y} \rangle \mid \mathbf{y} \geq_{\mathcal{K}} -\mathbf{c}_+, \mathbf{y} \geq_{\mathcal{K}} \mathbf{c}_- \}. \quad (3.7)$$

Proof. The dual problem reads

$$\min \{ \langle \mathbf{l}, \mathbf{y}_2 \rangle - \langle \mathbf{u}, \mathbf{c}_1 \rangle \mid \mathbf{y}_2 - \mathbf{y}_1 = \mathbf{c}, \mathbf{y}_1, \mathbf{y}_2 \in \mathcal{K} \}, \quad (3.8)$$

since \mathcal{K} is self-dual. We can reparametrize $\mathbf{y}_1 = \mathbf{y} - \mathbf{c}_-$ and $\mathbf{y}_2 = \mathbf{y} + \mathbf{c}_+$ to satisfy the equality constraint and get

$$\min \{ \langle \mathbf{l} - \mathbf{u}, \mathbf{y} \rangle + \langle \mathbf{l}, \mathbf{c}_+ \rangle + \langle \mathbf{u}, \mathbf{c}_- \rangle \mid \mathbf{y} \geq_{\mathcal{K}} -\mathbf{c}_+, \mathbf{y} \geq_{\mathcal{K}} \mathbf{c}_- \}.$$

Since $\langle \mathbf{l}, \mathbf{c}_+ \rangle + \langle \mathbf{u}, \mathbf{c}_- \rangle$ is constant, the result follows. \square

While Lemma 3.18 might seem rather uninteresting on its own, it has implications for cones whose cone orders $\leq_{\mathcal{K}}$ have special structural properties.

Remark 3.19:

The optimal solution of the dual (3.7) is a least upperbound on the set $\{\mathbf{c}_-, -\mathbf{c}_+\}$ in the partial ordered set $(\mathbb{R}^n, \leq_{\mathcal{K}})$. If $(\mathbb{R}^n, \leq_{\mathcal{K}})$ is a lattice, then the solution of the dual can be recovered from the join $\mathbf{c}_- \vee -\mathbf{c}_+$ in $(\mathbb{R}^n, \leq_{\mathcal{K}})$.

The nonnegative orthant $\mathcal{K} = \mathbb{R}_+^n$

Since the original paper [FHYZ08] treats this case, it is naturally suited for CLPN. In contrast to our general setting, the cone \mathbb{R}_+^n satisfies the useful property that the minimum-norm-point algorithm converges to the optimal solution in a finite number of iterations [Wol]. As such, the LP-Newton method converges in a finite number of steps [FHYZ08], and an important open problem is to decide whether this number can be bounded by a polynomial in the input size.

To see that (3.6) can easily be solved in linear time, note that a solution \mathbf{x}^* is given by greedily choosing the largest increase of the objective function componentwise by setting

$$x_i^* = \begin{cases} l_i & \text{if } c_i < 0, \\ u_i & \text{else.} \end{cases} \quad (3.9)$$

Another way to see the optimality of the greedy algorithm is by noting that \leq satisfies the lattice property. Then Remark 3.19 shows $\mathbf{y} = \mathbf{0}$ in (3.7), since

$$y_i = (\mathbf{c}_- \vee -\mathbf{c}_+)_i = \max\{0, -|c_i|\} = 0,$$

which confirms (3.9) through the dual variables $\mathbf{y}_1 = -\pi_{\mathbb{R}_+^n}(\mathbf{c})$ and $\mathbf{y}_2 = \pi_{\mathbb{R}_+^n}(\mathbf{c})$ in (3.8).

The Lorentz-cone $\mathcal{K} = \mathcal{L}_n$

Recall from Section 2.4 that the Lorentz-cone is given as

$$\mathcal{L}_n = \{(x_0, \tilde{\mathbf{x}}) \in \mathbb{R}_+ \times \mathbb{R}^n \mid \|\tilde{\mathbf{x}}\|_2 \leq x_0\}.$$

A useful property of the Lorentz-cone is the explicit description of its boundary

$$\text{bd}(\mathcal{L}_n) = \{(x_0, \tilde{\mathbf{x}}) \in \mathbb{R}_+ \times \mathbb{R}^n \mid \|\tilde{\mathbf{x}}\|_2 = x_0\},$$

which we can exploit by applying Lemma 3.17.

To simplify notation, we will denote the vector $\mathbf{a} = (a_0, a_1, \dots, a_n) \in \mathbb{R}^{n+1}$ by writing $\mathbf{a}^\top = (a_0, \tilde{\mathbf{a}}^\top)$, where $\tilde{\mathbf{a}} = (a_1, \dots, a_n) \in \mathbb{R}^n$ and $a_0 \in \mathbb{R}$. In addition, we will assign a special set $\mathcal{E}(\mathbf{w})$ to each $\mathbf{w} \in \mathcal{L}_n$, which we formally define as

$$\mathcal{E}(\mathbf{w}) = \{\mathbf{x} \in \mathbb{R}^{n+1} \mid \|\tilde{\mathbf{x}}\|_2^2 = x_0^2, \quad \|\tilde{\mathbf{w}} - \tilde{\mathbf{x}}\|_2^2 = (w_0 - x_0)^2\}. \quad (3.10)$$

Using this notation, we can show the following representation result.

Lemma 3.20:

For any $\mathbf{w} = (w_0, \tilde{\mathbf{w}}) \in \text{int}(\mathcal{L}_n)$, define the parameters

$$\bar{\mathbf{w}} := \frac{1}{w_0} \tilde{\mathbf{w}}, \quad \bar{w}_0 := \frac{w_0^2 - \|\tilde{\mathbf{w}}\|_2^2}{2w_0}, \quad \mathbf{Q} := \mathbf{I}_n - \bar{\mathbf{w}}\bar{\mathbf{w}}^\top \quad \text{and} \quad \gamma := \sqrt{\frac{1}{2}w_0\bar{w}_0}.$$

Then

$$\mathcal{E}(\mathbf{w}) = \left\{ \mathbf{x} \in \mathbb{R}^{n+1} \mid x_0 = \langle \tilde{\mathbf{x}}, \bar{\mathbf{w}} \rangle + \bar{w}_0, \left\| \mathbf{Q}^{\frac{1}{2}} \left(\tilde{\mathbf{x}} - \frac{1}{2} \tilde{\mathbf{w}} \right) \right\|_2^2 = \gamma^2 \right\}$$

and in particular, $\mathcal{E}(\mathbf{w})$ is an n -dimensional ellipsoid.

Proof. Subtracting the equations in (3.10) immediately shows containment in the hyperplane

$$H = \left\{ \mathbf{x} \in \mathbb{R}^{n+1} \mid x_0 = \langle \tilde{\mathbf{x}}, \bar{\mathbf{w}} \rangle + \bar{w}_0 \right\}, \quad (3.11)$$

where $\bar{\mathbf{w}}$ and \bar{w}_0 are well defined since $\mathbf{w} \in \text{int}(\mathcal{L}_n)$.

Using (3.11) in either equation in (3.10) on x_0 yields an equation of the form

$$0 = \tilde{\mathbf{x}}^\top \mathbf{Q} \tilde{\mathbf{x}} - 2 \langle \bar{w}_0 \bar{\mathbf{w}}, \tilde{\mathbf{x}} \rangle - \bar{w}_0^2$$

where \mathbf{Q} is positive definite since $\|\bar{\mathbf{w}}\|_2 < 1$. Completing the square yields the equivalent condition

$$\left\| \mathbf{Q}^{\frac{1}{2}} \left(\tilde{\mathbf{x}} - \bar{w}_0 \mathbf{Q}^{-1} \bar{\mathbf{w}} \right) \right\|_2^2 = \bar{w}_0^2 + \bar{w}_0^2 \cdot \bar{\mathbf{w}}^\top \mathbf{Q}^{-1} \bar{\mathbf{w}},$$

which defines an n -dimensional ellipsoid.

Using the Sherman-Morrison formula (2.1) we can simplify

$$\bar{w}_0 \mathbf{Q}^{-1} \bar{\mathbf{w}} = \frac{1}{2} \tilde{\mathbf{w}}$$

and

$$\bar{w}_0^2 + \bar{w}_0^2 \cdot \bar{\mathbf{w}}^\top \mathbf{Q}^{-1} \bar{\mathbf{w}} = \frac{1}{2} w_0 \bar{w}_0 = \gamma^2.$$

□

Lemma 3.21:

For any $\mathbf{w} \in \mathcal{L}_n$, the linear optimization problem

$$\max \{ \langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{x} \in \mathcal{E}(\mathbf{w}) \} \quad (3.12)$$

can be solved in $\mathcal{O}(n)$.

Proof. We will distinguish the cases $\mathbf{w} \in \text{bd}(\mathcal{L}_n)$ and $\mathbf{w} \in \text{int}(\mathcal{L}_n)$, which can be checked in $\mathcal{O}(n)$.

For $\mathbf{w} \in \text{bd}(\mathcal{L}_n)$, we claim that $\mathcal{E}(\mathbf{w}) = \text{conv}(\{\mathbf{0}, \mathbf{w}\})$. By assumption, $\|\tilde{\mathbf{w}}\|_2 = w_0$ and $\|\tilde{\mathbf{x}}\|_2 = x_0$, so

$$\|\tilde{\mathbf{w}} - \tilde{\mathbf{x}}\|_2 = \|\tilde{\mathbf{w}}\|_2 - \|\tilde{\mathbf{x}}\|_2 \quad \forall \mathbf{x} \in \mathcal{E}(\mathbf{w}),$$

which is only possible if $\tilde{\mathbf{x}}$ is a multiple of $\tilde{\mathbf{w}}$. Consequentially, x_0 is the same multiple of w_0 , which shows the claim. Then (3.12) is equal to either 0 or $\langle \mathbf{c}, \mathbf{w} \rangle$.

For $\mathbf{w} \in \text{int}(\mathcal{L}_n)$, we can rewrite the result of Lemma 3.20 as

$$\mathcal{E}(\mathbf{x}) = \left\{ \mathbf{x} \in \mathbb{R}^{n+1} \mid x_0 = \langle \tilde{\mathbf{x}}, \bar{\mathbf{w}} \rangle + \bar{w}_0, \tilde{\mathbf{x}} = \mathbf{Q}^{-\frac{1}{2}} \mathbf{y} + \frac{1}{2} \tilde{\mathbf{w}}, \|\mathbf{y}\|_2^2 = \gamma^2 \right\}.$$

This yields a reparametrization of (3.12) in terms of \mathbf{y} with objective value

$$\langle \mathbf{c}, \mathbf{x} \rangle = \langle \tilde{\mathbf{c}} + \frac{c_0}{w_0} \tilde{\mathbf{w}}, \tilde{\mathbf{x}} \rangle + c_0 \bar{w}_0 \equiv \langle \tilde{\mathbf{c}} + \frac{c_0}{w_0} \tilde{\mathbf{w}}, \mathbf{Q}^{-\frac{1}{2}} \mathbf{y} + \frac{1}{2} \tilde{\mathbf{w}} \rangle \equiv \langle \mathbf{Q}^{-\frac{1}{2}} (\tilde{\mathbf{c}} + \frac{c_0}{w_0} \tilde{\mathbf{w}}), \mathbf{y} \rangle,$$

where \equiv denotes equality up to a constant difference. We need to distinguish two cases:

If $\mathbf{Q}^{-\frac{1}{2}} (\tilde{\mathbf{c}} + \frac{c_0}{w_0} \tilde{\mathbf{w}}) = \mathbf{0}$, then every solution is optimal, and we set

$$\mathbf{y}^* = \gamma \cdot \frac{\mathbf{e}}{\sqrt{n}}.$$

Otherwise, we set

$$\mathbf{y}^* = \gamma \cdot \frac{\mathbf{Q}^{-\frac{1}{2}} (\tilde{\mathbf{c}} + \frac{c_0}{w_0} \tilde{\mathbf{w}})}{\|\mathbf{Q}^{-\frac{1}{2}} (\tilde{\mathbf{c}} + \frac{c_0}{w_0} \tilde{\mathbf{w}})\|_2},$$

since the optimal solution of a linear optimization problem over a scaled Euclidean ball is parallel to the objective function.

In both cases, we recover the optimal solution of (3.12) by using the parametrization

$$\tilde{\mathbf{x}}^* = \mathbf{Q}^{-\frac{1}{2}} \mathbf{y}^* + \frac{1}{2} \tilde{\mathbf{w}}, \quad x_0^* = \langle \tilde{\mathbf{x}}^*, \bar{\mathbf{w}} \rangle + \bar{w}_0.$$

To see that \mathbf{x}^* can be computed in linear time, note that this is equivalent to carrying out a multiplication by $\mathbf{Q}^{-\frac{1}{2}}$ in time $O(n)$. This can be achieved by using a specific formula for $\mathbf{Q}^{-\frac{1}{2}}$, as we will show now. For this, we first compute a set of parameters depending on $\tilde{\mathbf{w}}$.

If $\tilde{\mathbf{w}} \neq \mathbf{0}$, we set

$$\alpha := \frac{1 - \frac{2\gamma}{w_0}}{\|\tilde{\mathbf{w}}\|_2^2}, \quad \beta := \frac{\alpha}{1 - \alpha \|\tilde{\mathbf{w}}\|_2^2} = \frac{\frac{w_0}{2\gamma} - 1}{\|\tilde{\mathbf{w}}\|_2^2},$$

and $\alpha = \beta = 0$ otherwise. Using $w_0 > \|\tilde{\mathbf{w}}\|_2$, one can verify that $\alpha, \beta \geq 0$ and a straightforward computation shows the identities

$$\mathbf{Q}^{\frac{1}{2}} = \mathbf{I}_n - \alpha \tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top, \quad \mathbf{Q}^{-\frac{1}{2}} = \mathbf{I}_n + \beta \tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top, \quad (3.13)$$

where one identity can be reduced to the other by the Sherman-Morrison formula (2.1). Then (3.13) shows that multiplication by $\mathbf{Q}^{-\frac{1}{2}}$ amounts to

$$\mathbf{Q}^{-\frac{1}{2}} \tilde{\mathbf{x}} = \tilde{\mathbf{x}} + \beta \langle \tilde{\mathbf{x}}, \tilde{\mathbf{w}} \rangle \tilde{\mathbf{w}},$$

where the right-hand side can be computed in $O(n)$. □

Theorem 3.22:

The cone \mathcal{L}_n is suitable for CLPN.

Proof. Through translation we can assume that $\mathbf{l} = \mathbf{0}$ and focus on the case

$$\max \{ \langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{x} \in [\mathbf{0}, \mathbf{w}]_{\mathcal{L}_n} \subseteq \mathbb{R}^{n+1} \} \quad (3.14)$$

where $\mathbf{w} = \mathbf{u} - \mathbf{l} \in \mathcal{L}_n \setminus \{\mathbf{0}\}$ and consequently $w_0 = u_0 - l_0 > 0$. By Lemma 3.17, it suffices to compute

$$\max \{ \langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{x} \in \text{bd}(\mathcal{L}_n) \cap \text{bd}(\mathbf{w} - \mathcal{L}_n) =: \mathcal{E}'(\mathbf{w}) \}$$

and compare this value to $\langle \mathbf{c}, \mathbf{0} \rangle = 0$ and $\langle \mathbf{c}, \mathbf{w} \rangle$. Thus, we have

$$\begin{aligned} \mathcal{E}'(\mathbf{w}) &= \{ \mathbf{x} \in \mathbb{R}^{n+1} \mid \|\tilde{\mathbf{x}}\|_2 = x_0, \|\tilde{\mathbf{w}} - \tilde{\mathbf{x}}\|_2 = w_0 - x_0 \} \\ &= \{ \mathbf{x} \in \mathcal{E}(\mathbf{w}) \mid x_0 \in [0, w_0] \} \end{aligned}$$

by (3.10) and claim that $\mathcal{E}'(\mathbf{w}) = \mathcal{E}(\mathbf{w})$.

To see this, we can use Lemma 3.21 with objective $\mathbf{c}^\top = (\pm 1, \mathbf{0}^\top)$ to get

$$\max \{ \pm x_0 \mid \mathbf{x} \in \mathcal{E}(\mathbf{w}) \} = \frac{1}{2}w_0 \pm \frac{1}{2}\|\tilde{\mathbf{w}}\|_2 \in [0, w_0],$$

where the bounds follow from $\mathbf{w} \in \mathcal{L}_n$. Then (3.14) is equivalent to

$$\max \{ \langle \mathbf{c}, \mathbf{x} \rangle \mid \mathbf{x} \in \mathcal{E}(\mathbf{w}) \},$$

which can be done in linear time according to Lemma 3.21. \square

As a summary of the preceding results, we close this subsection with Algorithm 3.3, an explicit linear time algorithm for solving (3.14).

Algorithm 3.3: Explicit solution for (3.14)

Data: Data $\mathbf{l}, \mathbf{u}, \mathbf{c}$ for problem (3.14), \bar{w}_0, β, γ as in Theorem 3.22.

Result: Solution $\mathbf{x}^* \in [\mathbf{l}, \mathbf{u}]_{\mathcal{K}}$ to (3.14).

- 1 Initialize $\mathbf{w} = \mathbf{u} - \mathbf{l}$, $\mathbf{x}^* = \mathbf{0}$;
 - 2 **if** $\|\tilde{\mathbf{w}}\|_2^2 < w_0^2$ **then**
 - 3 $\tilde{\mathbf{y}}^* = \tilde{\mathbf{c}} + \frac{c_0}{w_0}\tilde{\mathbf{w}}$;
 - 4 $\tilde{\mathbf{y}}^* = \tilde{\mathbf{y}}^* + \beta\langle \tilde{\mathbf{y}}^*, \tilde{\mathbf{w}} \rangle \tilde{\mathbf{w}}$;
 - 5 **if** $\tilde{\mathbf{y}}^* = \mathbf{0}$ **then**
 - 6 $\tilde{\mathbf{y}}^* = \gamma \cdot \frac{\mathbf{e}}{\sqrt{n}}$;
 - 7 **else**
 - 8 $\tilde{\mathbf{y}}^* = \gamma \cdot \frac{\tilde{\mathbf{y}}^*}{\|\tilde{\mathbf{y}}^*\|_2}$;
 - 9 $\tilde{\mathbf{x}}^* = \tilde{\mathbf{y}}^* + (\beta\langle \tilde{\mathbf{y}}^*, \tilde{\mathbf{w}} \rangle + \frac{1}{2})\tilde{\mathbf{w}}$;
 - 10 $x_0^* = \frac{1}{w_0}\langle \tilde{\mathbf{x}}^*, \tilde{\mathbf{w}} \rangle + \bar{w}_0$;
 - 11 $\mathbf{x}^* = \mathbf{x}^* + \mathbf{l}$;
 - 12 **return** $\text{argmax} \{ \langle \mathbf{y}, \mathbf{c} \rangle \mid \mathbf{y} \in \{\mathbf{l}, \mathbf{x}^*, \mathbf{u}\} \}$;
-

Remark 3.23:

The parameters \bar{w}_0 , β and γ only depend on $\mathbf{w} = \mathbf{u} - \mathbf{l}$. When optimizing multiple times over $[\mathbf{l}, \mathbf{u}]_{\mathcal{K}}$ with different objective functions, like in our setting, these parameters can be stored and need to be computed only once.

The positive semidefinite cone $\mathcal{K} = \mathcal{S}_+^n$

Let \mathcal{S}_+^n denote the cone of symmetrical $n \times n$ matrices that are positive semidefinite and let \leq be the corresponding conic order. For statements regarding complexity, note that we can embed $\mathcal{S}_+^n \subseteq \mathbb{R}^N$ for $N = \binom{n}{2} \in \mathcal{O}(n^2)$.

Then (3.6) reads

$$\max \{ \langle \mathbf{C}, \mathbf{X} \rangle \mid \mathbf{L} \leq \mathbf{X} \leq \mathbf{U} \}.$$

In order to solve this problem, we will transform it into a more suitable form. Just as with the case of the Lorentz cone, we first use the substitution $\mathbf{Y} = \mathbf{X} - \mathbf{L}$ to get the equivalent problem

$$\max \{ \langle \mathbf{C}, \mathbf{Y} \rangle \mid \mathbf{0} \leq \mathbf{Y} \leq \mathbf{W} \},$$

where $\mathbf{W} := \mathbf{U} - \mathbf{L}$ and the constant $\langle \mathbf{C}, \mathbf{L} \rangle$ was dropped from the objective. To simplify the argument, we will assume without loss of generality that $\mathbf{W} \in \text{int}(\mathcal{S}_+^n)$, since singular \mathbf{W} can be reduced to this case by changing to a suitable subspace. Using the Cholesky decomposition $\mathbf{W} = \mathbf{V}\mathbf{V}^\top$, we can rewrite $\mathbf{Y} = \mathbf{V}\mathbf{Z}\mathbf{V}^\top$ to get the equivalent problem

$$\max \{ \langle \mathbf{C}', \mathbf{Z} \rangle \mid \mathbf{0} \leq \mathbf{Z} \leq \mathbf{I}_n \}, \quad (3.15)$$

where $\mathbf{C}' = \mathbf{V}^\top \mathbf{C} \mathbf{V}$.

By applying these preprocessing steps, we reduced the original constraints to box-constraints on the eigenvalues of \mathbf{Z} , which allow us to solve the problem explicitly.

Theorem 3.24:

Let $\mathbf{C}' = \mathbf{B}^\top \mathbf{D} \mathbf{B}$ be the eigenvalue decomposition of \mathbf{C}' . Then the solution to (3.15) is given by

$$\mathbf{Z} = \mathbf{B}^\top \mathbf{\Lambda}^* \mathbf{B},$$

where $\mathbf{\Lambda}^*$ is a diagonal matrix with $\text{diag}(\mathbf{\Lambda}^*) = \boldsymbol{\lambda}^*$ and $\boldsymbol{\lambda}^*$ is the solution of

$$\max \{ \langle \text{diag}(\mathbf{D}), \boldsymbol{\lambda} \rangle \mid \boldsymbol{\lambda} \in [0, 1]^n \}. \quad (3.16)$$

Proof. Let $\mathbf{d} = \text{diag}(\mathbf{D})$ and let $\boldsymbol{\lambda}$ denote the eigenvalues of \mathbf{Z} . Then the Hoffman-Wielandt inequality [Ren10, Sct. 18.3.7] states

$$\langle \mathbf{C}', \mathbf{Z} \rangle \leq \langle \mathbf{d}, \mathbf{P}\boldsymbol{\lambda} \rangle,$$

where \mathbf{P} is a permutation that assigns the i -th biggest entry of $\boldsymbol{\lambda}$ to the i -th biggest entry of \mathbf{d} for all $i \in [n]$. Since $\boldsymbol{\lambda} \in [0, 1]^n$, the right-hand side is maximal when $\boldsymbol{\lambda}$ is the solution $\boldsymbol{\lambda}^*$ of (3.16) and \mathbf{P} the identity. Then one can verify that the left hand-side also attains this upperbound by choosing $\mathbf{Z} = \mathbf{B}^\top \text{Diag}(\boldsymbol{\lambda}^*) \mathbf{B}$. \square

Now that we showed that solving (3.15) reduces to the computation (3.9), it is important to note that the computational burden actually lies in the preprocessing. From the point of view of complexity, we have the following.

Theorem 3.25:

The cone \mathcal{S}_+^n is suitable for CLPN.

Proof. The complexity of computing the eigenvalues of a $n \times n$ matrix as well as matrix multiplication is contained in $\mathcal{O}(n^3)$. Since $\mathcal{S}_+^n \subseteq \mathbb{R}^N$ with $n \in \mathcal{O}(N^{1/2})$, we get an algorithm in $\mathcal{O}(n^3) \subseteq \mathcal{O}(N^{3/2}) \subseteq \mathcal{O}(N^2)$. \square

Remark 3.26:

The preceding theorem is noteworthy in terms of Remark 3.19, since the conic order \preceq induced by \mathcal{S}_+^n is explicitly known to *not* satisfy the lattice property. In particular, if we have a proper interval $\mathbf{L} < \mathbf{U}$, then Slater's condition holds and we expect strong duality to hold in Lemma 3.18, so that the preceding theorem yields an oracle for elements of the set of least upperbounds of $\{\mathbf{C}_-, -\mathbf{C}_+\}$ in (\mathbb{R}^N, \preceq) .

3.3.4 Experiments

In this section, we show some experiments done with a simple implementation of both MNP and CLPN. As a reference, we used the widespread SDPT3 package [TTT03].

Data generation

SOCP

Based on parameter tuples $\frac{n}{m}$, we generated random instances for $\mathcal{K} = \mathcal{L}_n$. We set $\mathbf{l} = \mathbf{0}$, $u_0 = 10$ and $\tilde{\mathbf{u}}$ to a random vector with entries in $[-0.5, 0.5]$, which was afterwards normalized such that $\|\tilde{\mathbf{u}}\|_2$ was a random number in the interval $[0, 10]$.

The vector \mathbf{c} was randomly chosen with entries in $[-0.5, 0.5]$ and \mathbf{A} was chosen as a random $m \times n$ matrix with entries in $[0, 1]$. To guarantee feasibility, we included the midpoint of $[\mathbf{l}, \mathbf{u}]_{\mathcal{L}_n}$ into the feasible region by setting $\mathbf{b} = \frac{1}{2}\mathbf{A}(\mathbf{u} - \mathbf{l})$.

SDP

Based on parameter tuples $\frac{n}{m}$, we generated random instances for $\mathcal{K} = \mathcal{S}_+^n$. We set $\mathbf{L} = \mathbf{0}$ and constructed \mathbf{U} in the following way: We first constructed a random $n \times n$ matrix \mathbf{V} with values in $[0, 1]$ and set $\mathbf{U} = \mathbf{V}\mathbf{V}^\top + \frac{1}{10}\mathbf{I}_n$. Afterwards, \mathbf{U} was normalized such that $\text{tr}(\mathbf{U}) = 10$.

The remaining parameters were chosen in the same way as for SOCP: \mathbf{C} was randomly chosen with entries in $[-0.5, 0.5]$ and \mathcal{A} was chosen as a random $m \times n^2$ linear

operator with entries in $[0, 1]$. To guarantee feasibility, we included the midpoint of $[\mathbf{L}, \mathbf{U}]_{S_+^n}$ into the feasible region by setting $\mathbf{B} = \frac{1}{2}\mathcal{A}(\mathbf{U} - \mathbf{L})$.

Plots

In the following plots, each data point $\frac{n}{m}$ corresponds to the average of 25 instances randomly generated according to the procedure outlined before with parameters $\frac{n}{m}$. The error tolerance for CLPN was set to 10^{-6} and the error tolerance ε given in the plots apply to the subroutine MNP.

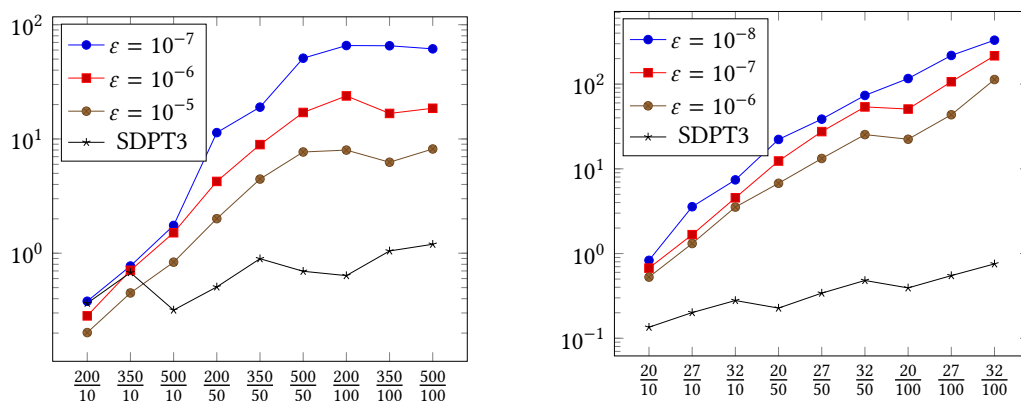


Figure 3.2: Runtime (sec) for parameters $\frac{n}{m}$. Left: SOCP. Right: SDP.

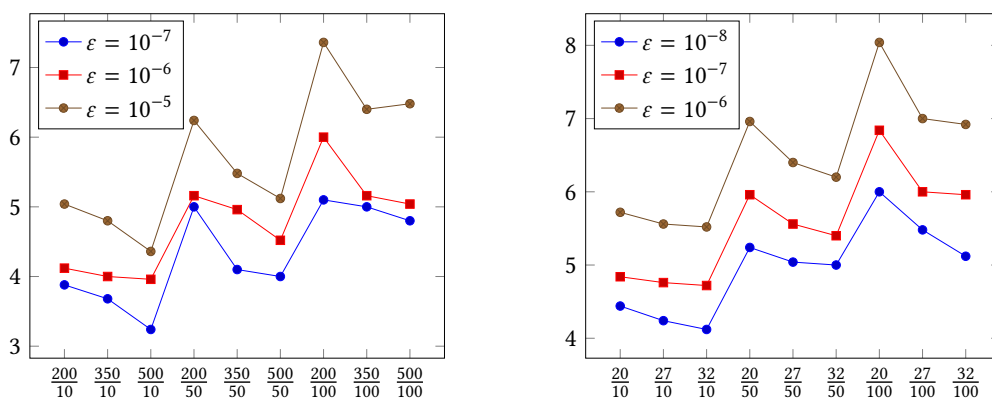


Figure 3.3: Newton-steps for parameters $\frac{n}{m}$. Left: SOCP. Right: SDP.

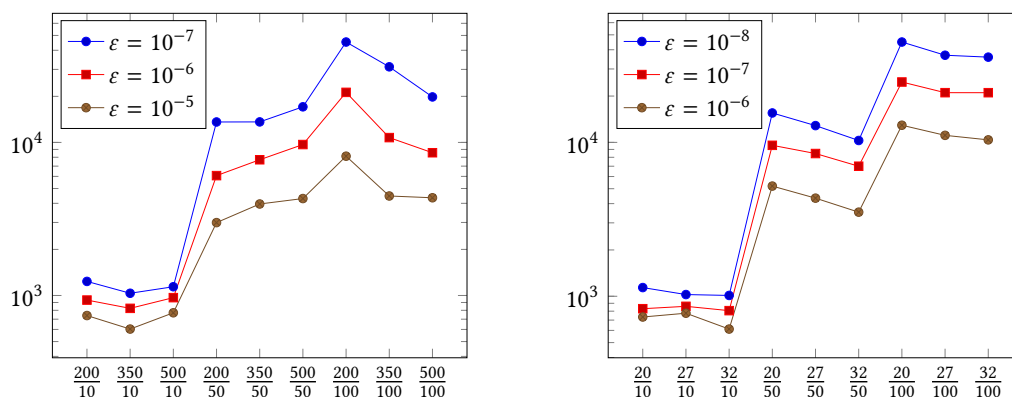


Figure 3.4: MNP computations for parameters $\frac{n}{m}$. Left: SOCP. Right: SDP.

The plots in Figure 3.2 show that the choice of accuracy for MNP has a great impact on the overall running time of the algorithm. While reducing the accuracy can speed up the algorithm significantly, going below the accuracy given in the plots often resulted in major problems in converging to the solution, so care has to be taken by choosing this parameter.

Overall, the data in Figure 3.3 resembles the results of [FHYZ08] for the case of \mathbb{R}_+^n , in the sense that only a few Newton-steps are necessary to get a close approximate solution. Figure 3.4 also shows that, like in the original paper, the number m of constraints seems to have a much larger impact on the performance than the number of the variables n , since much more subproblems have to be solved.

Conclusion

Unfortunately, the experiments show that our implementation of the CLP-Newton method is much slower than the reference algorithm. However, the runtime of the algorithm could be drastically reduced by a faster algorithm for the minimum-norm-point subroutine, since the number of Newton-steps remain small independent of the underlying cones. In particular, the algorithm might be improved by a more rigorous treatment of the necessary accuracy for the subproblems, since the experiments indicate much slower progress with increasing accuracy.

Chapter 4

Partitions and Assignment Matrices

In this chapter, we formally introduce partition problems in Section 4.1 and discuss the (dis)advantages of representing them as assignment matrices in Section 4.2, where we discuss problems related to symmetries and how to solve them with the theory of orbitopes.

4.1 Overview

Partitions

For a fixed number $n \in \mathbb{N}$, a *partition* of $[n]$ is a subset $\{T_j \mid j \in [k]\} \subseteq 2^{[n]}$ such that

- (i) $T_j \neq \emptyset \quad \forall j \in [k]$,
- (ii) $\bigcup_{j \in [k]} T_j = [n]$,
- (iii) $T_i \cap T_j = \emptyset \quad \forall i, j \in [k]$.

The sets that make up a partition are called *parts*. Throughout the thesis, we will use the shorthand notation $\mathcal{T} = \{T_j \mid j \in [k]\}$ and assume that $L \subseteq 2^{[n]}$ is an independence system.

Definition 4.1 (Set of Partitions):

For $k \in [n]$, let \mathcal{P}_k^n denote the set of all partitions of $[n]$ that consist of exactly k parts. Then the set

$$\mathcal{P}_k^n(L) := \{\mathcal{T} \in \mathcal{P}_k^n \mid \mathcal{T} \subseteq L\}$$

denotes the set of all partitions of $[n]$ whose k parts belong to L . More general, we let

$$\mathcal{P}^n(L) = \bigcup_{k \in [n]} \mathcal{P}_k^n(L)$$

be the set of all partitions of $[n]$ with regards to L and we shortly write $\mathcal{P}^n := \mathcal{P}^n(2^{[n]})$ as well.

Lemma 4.2 (Cardinality, [Sta11, p. 73ff.]):

The cardinality of \mathcal{P}_k^n is given by $\left\{ \begin{smallmatrix} n \\ k \end{smallmatrix} \right\}$, the Stirling numbers of the second kind. They have the explicit description

$$\left\{ \begin{smallmatrix} n \\ k \end{smallmatrix} \right\} = \frac{1}{k!} \sum_{j \in [k]} (-1)^{k-j} \binom{k}{j} j^n,$$

and satisfy the recurrence

$$\left\{ \begin{smallmatrix} n \\ k \end{smallmatrix} \right\} = k \cdot \left\{ \begin{smallmatrix} n-1 \\ k \end{smallmatrix} \right\} + \left\{ \begin{smallmatrix} n-1 \\ k-1 \end{smallmatrix} \right\}$$

where $\left\{ \begin{smallmatrix} k \\ k \end{smallmatrix} \right\} = \left\{ \begin{smallmatrix} n \\ 1 \end{smallmatrix} \right\} = 1$ for all $n, k \in \mathbb{N}$.

The cardinality of \mathcal{P}^n is given by the Bell number $B(n)$ and satisfies the formulas

$$B(n) = \sum_{k \in [n]} S(n, k) = \sum_{i=0}^n \binom{n}{i} B(i),$$

where $B(0) = 1$.

k -partition Problems

Let $k \in [n]$. A function $f : \mathcal{P}_k^n(L) \rightarrow \mathbb{R}$ is called k -partition function and its associated optimization problem

$$\min \{ f(\mathcal{T}) \mid \mathcal{T} \in \mathcal{P}_k^n(L) \}$$

is called a k -partition problem. For reference, we will denote this problem as $(f, \mathcal{P}_k^n(L))$.

A k -partition function $f_k : \mathcal{P}_k^n(L) \rightarrow \mathbb{R}$ is separable (with regards to $f : L \rightarrow \mathbb{R}$), if for each $\mathcal{T} \in \mathcal{P}_k^n(L)$ we have

$$f_k(\mathcal{T}) = \sum_{j \in [k]} f(T_j),$$

and the corresponding k -partition problem is called separable as well.

Example 4.3 (Max-Cut):

Consider a graph $G \in \mathcal{G}_n$ with a weighted adjacency matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ to define a function $\text{cut}_{\mathbf{C}} : 2^{[n]} \rightarrow \mathbb{R}$ through

$$\text{cut}_{\mathbf{C}}(T) = \sum_{i,j \in T} c_{ij} \quad \forall T \subseteq [n].$$

Then the famous *Max- k -Cut* problem can be stated as the separable k -partition problem $(\text{cut}_k, \mathcal{P}_2^{[n]})$ where

$$\text{cut}_k(\mathcal{T}) = \sum_{j \in [k]} \text{cut}_{\mathbf{C}}(T_j)$$

is separable with regards to $\text{cut}_{\mathbf{C}}$ as defined above.

Remark 4.4:

The problem Max-2-Cut is often simply denoted as Max-Cut in literature.

As a corollary of the preceding example, we see that partition problems are NP-hard in general.

Theorem 4.5 ([Sch03, Vol. C, Theorem 75.1]):

The problem Max-2-Cut is NP-hard.

Partition Problems

A function $f : \mathcal{P}^n(L) \rightarrow \mathbb{R}$ is simply called *partition function* and the associated problem

$$\min \{f(\mathcal{T}) \mid \mathcal{T} \in \mathcal{P}^n(L)\}$$

is called a *partition problem*, which will be denoted as $(f, \mathcal{P}^n(L))$.

Likewise, a partition function $f_{[n]}$ is *separable* (with regards to $f : L \rightarrow \mathbb{R}$), if the restriction of $f_{[n]}$ to k -partitions is separable with regards to f for every $k \in [n]$, that is

$$f_{[n]}(\mathcal{T}) = \sum_{j \in [k]} f(T_j) \quad \forall \mathcal{T} \in \mathcal{P}_k^n(L)$$

holds for all $k \in [n]$.

A partition function f is called a *counting function* if it only depends on the number of parts in a partition and is strictly increasing. In particular, if f is a counting function, then there is a strictly increasing function $g : \mathbb{N} \rightarrow \mathbb{R}$ such that

$$f(\mathcal{T}) = g(k)$$

holds for all k -partitions $\mathcal{T} \in \mathcal{P}_k^n(L)$. We may abuse notation by denoting the counting function by $g(k)$ directly instead.

A partition problem whose partition function is a counting function is called a *minimum-cover problem* and has the following combinatorial interpretation.

Lemma 4.6:

Let $(g, \mathcal{P}^n(L))$ be a minimum-cover problem. Then solving $(f, \mathcal{P}^n(L))$ is equivalent to finding

$$\min \{k \in \mathbb{N} \mid \mathcal{P}_k^n(L) \neq \emptyset\},$$

and we can assume without loss of generality that $g = \text{Id}$.

Example 4.7 (Graph Colouring):

Consider a graph $G \in \mathcal{G}_n$ and denote by $S_G \subseteq 2^{[n]}$ the *stable sets* of G where

$$S_G = \{S \subseteq [n] \mid E(G[S]) = \emptyset\},$$

that is, $S \in S_G$ if and only if the induced subgraph $G[S]$ does not contain any edges. Note that this way, S_G satisfies the axioms (2.2) of an independence system.

Then the minimum-cover problem $(\text{Id}, \mathcal{P}_k^n(S_G))$ is the problem of finding the smallest partition of $[n]$ into stable sets of G , which is known as the *graph colouring problem*. In particular, the *chromatic number* $\chi(G)$ can be expressed as

$$\chi(G) = \min \{k \in \mathbb{N} \mid \mathcal{P}_k^n(S_G) \neq \emptyset\}.$$

A more detailed treatment of the colouring problem is the topic of Section 6.3.1.

Lastly, we also consider *regularized partition functions* $f_g : \mathcal{P}^n(L) \rightarrow \mathbb{R}$ that can be written as

$$f_g = f_{[n]} + g$$

where $f_{[n]}$ is a separable partition function and g is a counting function that serves as a *regularizer*.

Example 4.8 (Regularized Euclidean clustering):

Consider a set of points $\{\mathbf{b}_i\}_{i \in [n]} \subseteq \mathbb{R}^d$ in Euclidean space. Then we can define an optimal-value function $f : 2^{[n]} \rightarrow \mathbb{R}$ by setting

$$f(T) = \min_{\mathbf{x} \in \mathbb{R}^d} \sum_{i \in T} \|\mathbf{x} - \mathbf{b}_i\|_2^2$$

for all $T \subseteq [n]$ to get a partition function $f_{[n]}$ separable with regards to f . Now using $g(k) = k^2$ and some penalty parameter $\mu > 0$, the regularized partition function $f_{\mu, g}$ corresponds to quadratically regularized *Euclidean clustering*. The treatment of Euclidean clustering problems is the topic of Section 6.3.2.

4.2 Assignment Matrices

In the following, we will describe how to fit partitions into an optimization framework by giving a standard representation of partitions as *assignment matrices*. Given a set $I \subseteq [n]$, its *characteristic vector* $\mathbf{e}_I \in \{0, 1\}^n$ is coordinate-wise defined as

$$(\mathbf{e}_I)_i = \begin{cases} 1 & \text{if } i \in I, \\ 0 & \text{else.} \end{cases}$$

Geometrically, the map $I \mapsto \mathbf{e}_I$ bijectively maps the power set 2^n to the vertices of the n -dimensional unit cube C^n . In particular, we have

$$I \cap J = \emptyset \Leftrightarrow \langle \mathbf{e}_I, \mathbf{e}_J \rangle = 0 \quad (4.1)$$

and

$$\mathbf{e}_{(I \cup J)} = \mathbf{e}_I + \mathbf{e}_J - \mathbf{e}_{(I \cap J)}. \quad (4.2)$$

Our goal is now to express k -partitions through $n \times k$ binary matrices by treating their columns as characteristic vectors. To this end, let

$$V(L) = \{\mathbf{e}_I \in \{0, 1\}^n \mid I \in L_*\} = \{\mathbf{x} \in \{0, 1\}^n \mid \text{supp}(\mathbf{x}) \in L_*\}$$

be the vertex-set associated with L and define the set of L -constrained assignment matrices

$$\mathcal{U}_{n,k}(L) = \{\mathbf{U} \in \{0, 1\}^{n \times k} \mid \mathbf{U}\mathbf{e}_{[k]} = \mathbf{e}_{[n]}, \text{col}(\mathbf{U}) \subseteq V(L)\}. \quad (4.3)$$

We will also use the notation

$$\mathcal{U}_{n,k} := \mathcal{U}_{n,k}(2^{[n]}) = \{\mathbf{U} \in \{0, 1\}^{n \times k} \mid \mathbf{U}\mathbf{e}_{[k]} = \mathbf{e}_{[n]}, \mathbf{0} \notin \text{col}(\mathbf{U})\},$$

which immediately shows the alternative expression

$$\mathcal{U}_{n,k}(L) = \mathcal{U}_{n,k} \cap \mathcal{V}_{\mathbb{R}}(L_*^k)$$

for (4.3).

By examining (4.1) and (4.2), we can verify that each matrix in $\mathcal{U}_{n,k}(L)$ can be turned into a partition by treating the columns as characteristic vectors. Unfortunately, we have multiple matrices in $\mathcal{U}_{n,k}(L)$ that yield the same partition.

Example 4.9:

Both matrices

$$\begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}$$

represent the partition $\{\{1, 2\}, \{3\}\}$ by treating their columns as characteristic vectors.

This lack of unique representation arises since the elements of $\mathcal{U}_{n,k}(L)$ are *ordered* columns, while the elements of $\mathcal{P}_k^n(L)$ are *unordered* sets. More formally, the group \mathfrak{S}_k induces a group action on $\mathcal{U}_{n,k}(L)$ by permutation of its columns. Setting

$$\pi * (\mathbf{u}_1, \dots, \mathbf{u}_k) = (\mathbf{u}_{\pi(1)}, \dots, \mathbf{u}_{\pi(k)}), \quad \forall \pi \in \mathfrak{S}_k, \forall (\mathbf{u}_1, \dots, \mathbf{u}_k) \in \mathcal{U}_{n,k}(L),$$

recall that the *orbit* of a matrix $\mathbf{U} \in \mathcal{U}_{n,k}(L)$ under this *group action* is the set

$$[\mathbf{U}] := \{\pi * \mathbf{U} \mid \pi \in \mathfrak{S}_k\},$$

and we expect each orbit to correspond to a partition in $\mathcal{P}_k^n(L)$. To properly see this, let us introduce an order on $\mathcal{U}_{n,k}(L)$ to define a unique representative for each orbit by choosing the maximal member with regards to that order.

Definition 4.10 (Maximal Orbit Representatives, [KP08]):

Define an order $<_{lex}$ on $\mathbb{R}^{n \times k}$ by setting $\mathbf{A} <_{lex} \mathbf{B}$ if and only if $a_{ij} < b_{ij}$ for the smallest (i, j) w.r.t. $<_{lex}$ such that $a_{ij} \neq b_{ij}$. Then $\mathcal{U}_{n,k}^{lex}(L)$ is defined as the set of matrices $\mathbf{U} \in \mathcal{U}_{n,k}(L)$ which are maximal with regards to $<_{lex}$ in their orbit $[\mathbf{U}]$.

Example 4.11:

The set $\mathcal{U}_{3,2}$ listed in ascending order of $<_{lex}$ is given by

$$\begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{pmatrix} <_{lex} \begin{pmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} <_{lex} \begin{pmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 0 \end{pmatrix} <_{lex} \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix} <_{lex} \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{pmatrix} <_{lex} \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

With this definition at hand, we can state the core representation result.

Theorem 4.12:

We have

$$\mathcal{P}_k^n(L) \leftrightarrow \mathcal{U}_{n,k}^{lex}(L),$$

or rather

$$\mathcal{P}_k^n(L) \times \mathfrak{S}_k \leftrightarrow \mathcal{U}_{n,k}(L). \quad (4.4)$$

Proof. We define the map $\varphi: \mathcal{P}_k^n(L) \rightarrow \mathcal{U}_{n,k}^{lex}(L)$, which assigns $\mathcal{T} \in \mathcal{P}_k^n(L)$ to

$$\varphi(\mathcal{T}) =: (\pi_{\mathcal{T}} * (\mathbf{e}_{T_1}, \dots, \mathbf{e}_{T_k})) \in \mathcal{U}_{n,k}^{lex}(L),$$

where $\pi_{\mathcal{T}} \in \mathfrak{S}_k$ is chosen such that $\varphi(\mathcal{T})$ is maximal in its orbit.

This map is bijective: Its inverse takes the columns \mathbf{u}_i of $\mathbf{U} \in \mathcal{U}_{n,k}^{lex}(L)$, extracts their support $T_i = \text{supp}(\mathbf{u}_i)$ and outputs $\{T_1, \dots, T_k\}$. This inverse is well-defined, since by definition of $\mathbf{u}_i \in V(L)$, we have $T_i \in L_*$ nonempty, and since $\mathbf{U}\mathbf{e}_{[k]} = \mathbf{e}_{[n]}$, they are disjoint and satisfy $\cup_{j \in [k]} T_j = [n]$.

For the second part, it now suffices to argue that

$$\mathcal{U}_{n,k}^{lex}(L) \times \mathfrak{S}_k \leftrightarrow \mathcal{U}_{n,k}(L).$$

To see this, consider the map $\psi: \mathcal{U}_{n,k}^{lex}(L) \times \mathfrak{S}_k \rightarrow \mathcal{U}_{n,k}(L)$ given as

$$\psi(\mathbf{U}, \pi) = \pi * \mathbf{U}.$$

Since any $\mathbf{U} \in \mathcal{U}_{n,k}^{lex}(L)$ has k distinct columns, ψ is injective, and since it is surjective by definition of $\mathcal{U}_{n,k}^{lex}(L)$, the statement follows. \square

Fortunately, for $k = 2$, we can describe $\mathcal{U}_{n,2}^{lex}(L)$ explicitly.

Theorem 4.13:

Let $n \in \mathbb{N}$, then

$$\mathcal{U}_{n,2}^{lex}(L) = \{\mathbf{U} \in \mathcal{U}_{n,2}(L) \mid u_{11} = 1\}.$$

Proof. Let $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2) \in \mathcal{U}_{n,2}(L)$ and $\mathbf{U}' = (\mathbf{u}_2, \mathbf{u}_1)$. Then the orbit of \mathbf{U} is given by $[\mathbf{U}] = \{\mathbf{U}, \mathbf{U}'\}$ and by definition, $(\mathbf{u}_1)_1 \neq (\mathbf{u}_2)_1$. Thus

$$\mathbf{U}' <_{lex} \mathbf{U} \Leftrightarrow (\mathbf{u}_2)_1 < (\mathbf{u}_1)_1 \Leftrightarrow (\mathbf{u}_1)_1 = 1, (\mathbf{u}_2)_1 = 0,$$

which is what we wanted to show. \square

Of course, Theorem 4.13 is not very surprising. Since we already know that the first element will necessarily be part of one of the sets, the theorem merely states that we can list this set as the first one.

4.2.1 Symmetry induced Problems

Theorem 4.12 implies that in general, it is a very bad idea to use $\mathcal{U}_{n,k}(L)$ as a modelling basis. Since we would like to optimize functions defined on $\mathcal{P}_k^n(L)$, we can state (4.4) as

$$\mathcal{P}_k^n(L) \leftrightarrow \mathcal{U}_{n,k}(L)/\mathfrak{S}_k,$$

which means that the representation of any k -partition function f in this setting must be constant on the orbits $\mathcal{U}_{n,k}(L)/\mathfrak{S}_k$. This will necessarily introduce problems in any approach to this problem involving the whole set $\mathcal{U}_{n,k}(L)$, since both the feasible and the optimal set become highly symmetrical without applying any techniques for symmetry breaking.

Since we ultimately want to tackle k -partition problems with a convex formulation, it is paramount to highlight these problems in the convex setting. To this end, we will ignore the constraints given by L for now and assume $L = 2^{[n]}$ to focus on the easiest case $\mathcal{U}_{n,k}$.

The smallest convex body containing $\mathcal{U}_{n,k}$ is naturally $\text{conv}(\mathcal{U}_{n,k})$, but we will use

$$\mathcal{U}_{n,k} = \bigcup_{U \in \mathcal{U}_{n,k}^{\text{lex}}} [U]$$

to consider the subset $\text{conv}([U]) \subseteq \text{conv}(\mathcal{U}_{n,k})$ for $U \in \text{conv}(\mathcal{U}_{n,k})$ as well. To this end, we get the following characterization.

Theorem 4.14:

For all $U \in \mathcal{U}_{n,k}^{\text{lex}}$, we have

$$\text{bar}(\text{conv}(\mathcal{U}_{n,k})) = \text{bar}(\text{conv}([U])) = \frac{1}{k} \mathbf{J}_{n,k}$$

where bar denotes the barycenter. In particular, all convexified orbits intersect in $\frac{1}{k} \mathbf{J}_{n,k}$, and if $f: \mathbb{R}^{n \times k} \rightarrow \mathbb{R}$ is a convex function invariant under the group action of \mathfrak{S}_k , then

$$\min \{f(\mathbf{X}) \mid \mathbf{X} \in \text{conv}(\mathcal{U}_{n,k})\} = \min \{f(\mathbf{X}) \mid \mathbf{X} \in \text{conv}([U])\} = f\left(\frac{1}{k} \mathbf{J}_{n,k}\right).$$

Proof. The barycenter of a set is necessarily invariant under its symmetries, so we consider the Reynolds operator

$$\varphi: U \mapsto \frac{1}{k!} \sum_{\pi \in \mathfrak{S}_k} (\pi * U)$$

that maps any matrix $\mathbf{U} \in \mathcal{U}_{n,k}$ to the average over its orbit. In particular, this means that $\varphi(\mathbf{U}) \in \text{conv}([\mathbf{U}])$. Since

$$\varphi(\mathbf{U}) = \frac{1}{k!} \left(\sum_{\pi \in \mathfrak{S}_k} \mathbf{u}_{\pi(1)}, \dots, \sum_{\pi \in \mathfrak{S}_k} \mathbf{u}_{\pi(k)} \right) = \frac{1}{k} \left(\sum_{j \in [k]} \mathbf{u}_j, \dots, \sum_{j \in [k]} \mathbf{u}_j \right) = \frac{1}{k} \mathbf{J}_{n,k}$$

is constant on $\mathcal{U}_{n,k}$, the point $\frac{1}{k} \mathbf{J}_{n,k}$ is necessarily the barycenter of each convexified orbit $\text{conv}([\mathbf{U}])$ and $\text{conv}(\mathcal{U}_{n,k})$. Furthermore, φ is also constant on the whole of $\text{conv}(\mathcal{U}_{n,k})$ due to linearity.

Now, for any $\mathbf{U} \in \text{conv}(\mathcal{U}_{n,k})$, the second statement follows from

$$f(\mathbf{U}) = \frac{1}{k!} \sum_{\pi \in \mathfrak{S}_k} f(\pi * \mathbf{U}) \geq f(\varphi(\mathbf{U})) = f\left(\frac{1}{k} \mathbf{J}_{n,k}\right),$$

where the first equation follows from the invariance of f . \square

We visualize Theorem 4.14 before we look at its implications.

Example 4.15:

Consider the set

$$\mathcal{U}_{3,2} = \left\{ \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \right\},$$

and recall that it is listed ascending in $<_{lex}$. Then the corresponding set

$$\mathcal{U}_{3,2}^{lex} = \left\{ \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \right\}$$

contains exactly one element of each orbit. For visualization purposes, we will use the linear projection $\rho : \mathbb{R}^{3 \times 2} \rightarrow \mathbb{R}^2$ given by

$$\rho(\mathbf{U}) = \frac{1}{\sqrt{6}} \cdot \text{diag} \left(\begin{pmatrix} \sqrt{3} & 0 & -\sqrt{3} \\ -1 & 2 & -1 \end{pmatrix} \cdot \mathbf{U} \right),$$

and note that the projection is chosen in such a way that the line spanned by $\mathbf{J}_{3,2}$ is mapped to the origin. Then Figure 4.1 shows how the convex hull of each orbit intersects $\frac{1}{2} \mathbf{J}_{3,2}$, and that the convex hull of $\mathcal{U}_{3,2}^{lex}$ is much smaller due to the absence of symmetry.

Theorem 4.14 shows that if we want to work with $\text{conv}(\mathcal{U}_{n,k})$ and a convex objective function f , mild assumptions on f will lead to the trivial result of $\frac{1}{k} \mathbf{J}_{n,k}$, which, interpreted as a probabilistic statement, merely states that each partition is as likely as any other to be the optimal solution to the problem at hand. This is more likely when the objective function f is nonlinear, but may also happen in the linear case, when the underlying algorithm does not guarantee an extreme point as solution.

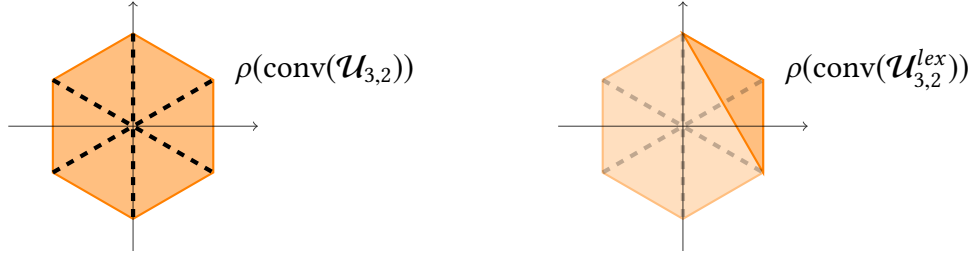


Figure 4.1: Visualization of Theorem 4.14. The shaded area shows the projection of convex hulls of $\mathcal{U}_{3,2}$ and $\mathcal{U}_{3,2}^{\text{lex}}$ respectively. The dashed lines are the projections of the convex hulls of the orbits, highlighting how they meet in the barycenter.

4.2.2 Orbitopes

Even if we work with the simplex-algorithm, which guarantees vertex solutions, Figure 4.1 still shows that it would be beneficial to rather work with the reduced polytope $\text{conv}(\mathcal{U}_{n,k}^{\text{lex}})$ instead. This is particular true for branch and bound algorithms, which may build up an extensive tree due to enumerating members of the same orbit in different branches.

To this end, Kaibel et al. used the term *partition-orbitope* for the set $\text{conv}(\mathcal{U}_{n,k}^{\text{lex}})$ and showed both an explicit description of exponential size [KP08], as well as a compactly lifted representation of linear size [FK09].

Remark 4.16:

The term *orbitope* has also been used for a related, but different construction in literature and should not be confused. In [Stu11], orbitopes are defined as convex hulls of the orbits of a compact algebraic group acting linearly on a real vector space.

This notion does not coincide with our definition, since we consider the convex hull of selected representatives from *multiple* orbits. In particular, $\mathcal{U}_{n,k}^{\text{lex}}(L)$ has a non-empty intersection with each orbitope arising from $\mathcal{U}_{n,k}(L)$ according to the second definition. This can be seen in Figure 4.1, where the dashed lines are projections of the orbitopes of the second definition.

We summarize the main result of the latter paper without getting into the details.

Theorem 4.17 ([FK09]):

There is an extended formulation $P_{n,k} \in \mathbb{R}^{O(nk)}$ for $\text{conv}(\mathcal{U}_{n,k}^{\text{lex}})$ such that $P_{n,k}$ is integral and can be described by $O(nk)$ inequalities. Furthermore, linear optimization over $P_{n,k}$, and, as a consequence, over $\text{conv}(\mathcal{U}_{n,k}^{\text{lex}})$ can be done in time $O(nk)$.

The idea here is that one can identify the matrices in $\mathcal{U}_{n,k}^{\text{lex}}$ with certain flows through a directed, acyclic network arising from a slightly altered $n \times k$ grid graph. In this setting, linear optimization can be reduced to a longest $s - t$ path problem and solved in linear time in the size of the grid.

As a special case of Theorem 4.17, we will explicitly show a description of the simple orbitopes $\text{conv}(\mathcal{U}_{n,2}^{\text{lex}})$.

Theorem 4.18:

Let $n \in \mathbb{N}$, then

$$\text{conv}(\mathcal{U}_{n,2}^{\text{lex}}) = \left\{ \begin{pmatrix} 1 & 0 \\ \mathbf{e} - \mathbf{u} & \mathbf{u} \end{pmatrix} \in \mathbb{R}^{n \times 2} \mid \mathbf{u} \in [0, 1]^{n-1}, \langle \mathbf{u}, \mathbf{e} \rangle \geq 1 \right\}.$$

Proof. Let C denote the set on the right handed side. Then by Theorem 4.13, $\mathcal{U}_{n,2}^{\text{lex}} \subseteq C$, and since C is convex, this holds true when taking the convex hull on both sides. To see the inverse, note that this is equivalent to showing that

$$\text{conv}(\{\mathbf{u} \in \{0, 1\}^{n-1} \mid \mathbf{u} \neq \mathbf{0}\}) = \{\mathbf{u} \in [0, 1]^{n-1} \mid \langle \mathbf{u}, \mathbf{e} \rangle \geq 1\},$$

by the characterization of $\mathcal{U}_{n,2}^{\text{lex}}$ in Theorem 4.13. This can be verified by showing that the vertices of both sets coincide, which follows by inspection of the facet defined by $\langle \mathbf{u}, \mathbf{e} \rangle \geq 1$. \square

While the theory of orbitopes takes care of all k -partition functions that can be expressed as linear functions over $\mathcal{U}_{n,k}$, they don't help much with non-linear functions. In particular, we argue that asymptotically, as n grows for fixed k , optimizing over $\text{conv}(\mathcal{U}_{n,k}^{\text{lex}})$ can only help with the first few elements. To see this, let the aggregation of $\text{conv}(\mathcal{U}_{n,k}^{\text{lex}})$ be defined as

$$\mathbf{A}(n, k) := \sum_{\mathbf{U} \in \mathcal{U}_{n,k}^{\text{lex}}} \mathbf{U}.$$

Lemma 4.19:

We have the recurrence

$$\mathbf{A}(n, k) = \begin{pmatrix} k \cdot \mathbf{A}(n-1, k) \\ \begin{Bmatrix} n-1 \\ k \end{Bmatrix} \cdot \mathbf{e}_{[k]}^\top \end{pmatrix} + \begin{pmatrix} \mathbf{A}(n-1, k-1) & \mathbf{0} \\ \mathbf{0} & \begin{Bmatrix} n-1 \\ k-1 \end{Bmatrix} \end{pmatrix}$$

where

$$\mathbf{A}(n, 1) = \mathbf{e}_{[n]} \text{ and } \mathbf{A}(k, k) = \mathbf{I}_k$$

for all $k, n \in \mathbb{N}$.

Proof. Denote the right hand side by $\mathbf{B}(n, k)$ and decompose it as the sum

$$\mathbf{B}(n, k) = \sum_{j \in [k]} \sum_{\mathbf{U} \in \mathcal{U}_{n-1,k}^{\text{lex}}} \begin{pmatrix} \mathbf{U} \\ \mathbf{e}_j^\top \end{pmatrix} + \sum_{\mathbf{U}' \in \mathcal{U}_{n-1,k-1}^{\text{lex}}} \begin{pmatrix} \mathbf{U}' & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix}.$$

We argue that each matrix in this sum belongs $\mathcal{U}_{n,k}^{\text{lex}}$ and occurs only once. To see the latter, note that both individual summations clearly involve different matrices, and we

can never have $\mathbf{U} = (\mathbf{U}' \ \mathbf{0})$ for any pair $\mathbf{U} \in \mathcal{U}_{n-1,k}^{lex}, \mathbf{U}' \in \mathcal{U}_{n-1,k-1}^{lex}$, since the last column of \mathbf{U} needs to have at least one non-zero entry to represent a k -partition.

To see membership in $\mathcal{U}_{n,k}^{lex}$, note that $\begin{pmatrix} \mathbf{U} \\ \mathbf{e}_j \end{pmatrix}$ is maximal among its orbit w.r.t. $<_{lex}$ since \mathbf{U} already is. For $\begin{pmatrix} \mathbf{U}' & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{pmatrix}$, the same is true since \mathbf{U}' is maximal among its \mathfrak{S}_{k-1} -orbit w.r.t. $<_{lex}$ and since the last column $\begin{pmatrix} \mathbf{0} \\ \mathbf{1} \end{pmatrix}$ is the smallest vector in this order.

We finish the proof by claiming that $\mathbf{B}(n, k)$ sums over as many matrices as $\mathbf{A}(n, k)$. Since $\mathcal{U}_{n,k}^{lex} \cong \mathcal{P}_k^n$, we necessarily have $\mathbf{A}(n, k)\mathbf{e}_{[k]} = \left\{ \begin{matrix} n \\ k \end{matrix} \right\} \cdot \mathbf{e}_{[n]}$. At the same time, we have

$$\mathbf{B}(n, k)\mathbf{e}_{[k]} = (k \cdot \left\{ \begin{matrix} n-1 \\ k \end{matrix} \right\} + \left\{ \begin{matrix} n-1 \\ k-1 \end{matrix} \right\}) \cdot \mathbf{e}_{[n]} = \left\{ \begin{matrix} n \\ k \end{matrix} \right\} \cdot \mathbf{e}_{[n]}$$

by Lemma 4.2, which shows the claim. \square

Since the barycenter of $\text{conv}(\mathcal{U}_{n,k})$ was given as the average over its vertices by Theorem 4.14, we want to investigate the average of vertices of $\text{conv}(\mathcal{U}_{n,k}^{lex})$, which is given by

$$\hat{\mathbf{A}}(n, k) = \mathbf{A}(n, k) / \left\{ \begin{matrix} n \\ k \end{matrix} \right\}.$$

Denote with $\hat{\mathbf{a}}(n, k, i)$ the i -th row of $\hat{\mathbf{A}}(n, k)$. It follows from Lemma 4.19 that for $1 \leq i < n$, we have the recurrence

$$\hat{\mathbf{a}}(n, k, i) = \frac{k \cdot \left\{ \begin{matrix} n-1 \\ k \end{matrix} \right\}}{\left\{ \begin{matrix} n \\ k \end{matrix} \right\}} \cdot \hat{\mathbf{a}}(n-1, k, i) + \frac{\left\{ \begin{matrix} n-1 \\ k-1 \end{matrix} \right\}}{\left\{ \begin{matrix} n \\ k \end{matrix} \right\}} \cdot (\hat{\mathbf{a}}(n-1, k-1, i) \ 0). \quad (4.5)$$

While giving explicit formulas is very hard in the general case, we have the following simple theorem for the case of bi-partitions.

Theorem 4.20:

For $k = 2$ and $1 < i \leq n$, we have

$$\hat{\mathbf{a}}(n, 2, i) = \begin{pmatrix} \frac{2^{n-2}-1}{2^{n-1}-1} & \frac{2^{n-2}}{2^{n-1}-1} \end{pmatrix}.$$

In particular,

$$\lim_{n \rightarrow \infty} \hat{\mathbf{a}}(n, 2, i) = \frac{1}{2} \cdot \mathbf{e}_{[2]}.$$

Proof. We show this result by induction on n . For $n = 2$, it follows by definition that $\mathbf{A}(2, 2) = \mathbf{I}_2$ and so

$$\hat{\mathbf{a}}(2, 2, 2) = (0 \ 1) = \begin{pmatrix} \frac{1-1}{2-1} & \frac{1}{2-1} \end{pmatrix}.$$

For $n > 2$, we use the fact that $\left\{ \begin{matrix} n \\ 2 \end{matrix} \right\} = 2^{n-1} - 1$. To see this, note that we can enumerate all 2-partitions of $[n]$ twice by enumerating the pairs $(T, [n] \setminus T)$, where $T \in 2^{[n]} \setminus \{\emptyset, [n]\}$.

Then for $1 < i < n$, (4.5) turns into

$$\begin{aligned}\hat{\mathbf{a}}(n, 2, i) &= \frac{2 \cdot \binom{n-1}{2}}{\binom{n}{2}} \cdot \hat{\mathbf{a}}(n-1, 2, i) + \frac{\binom{n-1}{1}}{\binom{n}{2}} \cdot (\hat{\mathbf{a}}(n-1, 1, i) \ 0) \\ &= \frac{2 \cdot (2^{n-2} - 1)}{2^{n-1} - 1} \cdot \begin{pmatrix} 2^{n-3} - 1 & 2^{n-3} \\ 2^{n-2} - 1 & 2^{n-2} - 1 \end{pmatrix} + \frac{1}{2^{n-1} - 1} \cdot (1 \ 0) \\ &= \begin{pmatrix} 2^{n-2} - 1 & 2^{n-2} \\ 2^{n-1} - 1 & 2^{n-1} - 1 \end{pmatrix}.\end{aligned}$$

For $i = n$, we also have

$$\hat{\mathbf{a}}(n, 2, i) = \frac{\binom{n-1}{2}}{\binom{n}{2}} \mathbf{e}_{[2]} + \frac{\binom{n-1}{1}}{\binom{n}{2}} \mathbf{e}_2 = \begin{pmatrix} 2^{n-2} - 1 & 2^{n-2} \\ 2^n - 1 & 2^n - 1 \end{pmatrix}.$$

□

In general, we can observe a similar process for $k > 2$ in the sense that

$$\left\| \hat{\mathbf{a}}(n, k, i) - \frac{1}{k} \mathbf{e}_{[k]} \right\|_2$$

is quickly decreasing both in i and in n . In particular, Figure 4.2 suggests the following conjecture.

Conjecture 4.21:

The limits

$$\lim_{n \rightarrow \infty} \text{dist} \left(\text{conv}(\mathcal{U}_{n,k}^{\text{lex}}), \text{bar}(\text{conv}(\mathcal{U}_{n,k})) \right) \leq \lim_{n \rightarrow \infty} \left\| \hat{\mathbf{A}}(n, k) - \frac{1}{k} \mathbf{J}_{n,k} \right\|_2$$

exist, are finite and grow sublinearly with k .

Unfortunately, this conjecture would imply that for fixed k and growing values of n , there is a good chance for the optimum of non-linear convex functions over $\text{conv}(\mathcal{U}_{n,k}^{\text{lex}})$ to stay relatively close to $\text{bar}(\text{conv}(\mathcal{U}_{n,k}))$, since the distance is independent of the dimension of the underlying space. In particular, using $\text{conv}(\mathcal{U}_{n,k}^{\text{lex}})$ over $\text{conv}(\mathcal{U}_{n,k})$ might only lead to marginal improvements, which is why we will consider reformulation techniques to break the underlying symmetry in the following chapters.

One main obstacle in proving the conjecture is the unavailability of sharp bounds for the fractions $\binom{n-1}{k} / \binom{n}{k}$ that appear in the convex combination used to update the rows in (4.5). These fractions play a major role in statistics to compute an unbiased estimation with minimum variation of certain distributions, as explained in [Ber75]. For fixed fractions $\frac{k}{n}$, we have the asymptotic behaviour

$$\frac{\binom{n-1}{k}}{\binom{n}{k}} \sim \frac{\alpha(\frac{k}{n})}{n},$$

where $\alpha(x)$ is the solution of $\frac{1-\exp(-\alpha)}{\alpha} = x$ according to [Har68], which implies that

$$\lim_{n \rightarrow \infty} \frac{k \binom{n-1}{k}}{\binom{n}{k}} \rightarrow 1.$$

This means that new rows $\hat{\mathbf{a}}(n, k, n)$ are already initialized increasingly close to $\frac{1}{k} \cdot \mathbf{e}_{[k]}$ and do not change much through the update (4.5), supporting the conjecture.

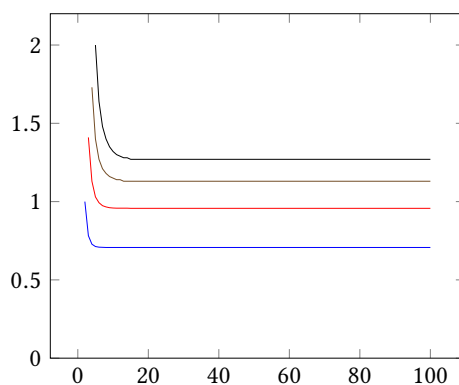


Figure 4.2: Plots of $\|\hat{\mathbf{A}}(n, k) - \frac{1}{k} \mathbf{J}_{n,k}\|_2$ against $n = k, k + 1, \dots, 100$ for increasing values of k . The plots for $k = 2, 3, 4, 5$ are shown from bottom to top respectively.

Chapter 5

Partition Matrices

This chapter treats partition matrices, a class of matrices that have been long used in combinatorial optimization for clustering problems like Max- k -cut or graph colouring. We introduce them in Section 5.1 and show how they relate to combinatorial moment matrices in Section 5.2. Finally, in Section 5.3 we show how these matrices are connected to assignment matrices through MM.

5.1 Overview

Another approach to remove symmetry from partition problems is to map the set $\mathcal{U}_{n,k}$ of assignment matrices onto the set of k -partition matrices, which are defined as follows.

Definition 5.1:

A binary matrix $\mathbf{W} \in \{0, 1\}^{n \times n}$ is called (k -)partition matrix, if there are k integers n_1, \dots, n_k and a permutation $\pi \in \mathfrak{S}_n$ such that $n = \sum_{j \in [k]} n_j$ and

$$\pi(\mathbf{W}) = \begin{pmatrix} \mathbf{J}_{n_1} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{J}_{n_2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} & \mathbf{J}_{n_k} \end{pmatrix}.$$

The set of all k -partition matrices will be denoted as PM_k^n .

The idea here is that partition matrices contain the information whether two elements belong to the same part of a given partition. To this end, note that for a given partition matrix \mathbf{W} , we have

$$w_{ij} = \begin{cases} 1 & \text{if } i \text{ and } j \text{ belong to the same set,} \\ 0 & \text{else.} \end{cases}$$

As a consequence, we have $\text{diag}(\mathbf{W}) = \mathbf{e}$ and each block corresponds to a part of the partition.

We can explicitly describe the set of k -partition matrices with the following lemma.

Lemma 5.2 ([Ren10, Lemma 18.3]):

The set PM_k^n can be explicitly described as

$$\text{PM}_k^n = \{\mathbf{W} \in \{0, 1\}^{n \times n} \mid \text{diag}(\mathbf{W}) = \mathbf{e}, \text{rank}(\mathbf{W}) = k, \mathbf{W} \geq \mathbf{0}\}.$$

With this representation in hand, we can now extend this concept to L -constrained k -partitions.

Definition 5.3:

The L -constrained k -partition matrices are given as

$$\text{PM}_k^n(L) := \{\mathbf{W} \in \text{PM}_k^n \mid \text{col}(\mathbf{W}) \subseteq V(L)\} = \text{PM}_k^n \cap \mathcal{V}_{\mathbb{R}}(L_*).$$

By definition, we only allow columns that represent nonempty sets belonging to L . We are now able to show that there is a one to one correspondence between partitions and partition matrices.

Lemma 5.4:

We have $\mathcal{U}_{n,k}^{\text{lex}}(L) \leftrightarrow \text{PM}_k^n(L)$ via the surjection $\varphi : \mathcal{U}_{n,k}(L) \rightarrow \text{PM}_k^n(L)$ given by

$$\mathbf{U} \mapsto \mathbf{U}\mathbf{U}^\top.$$

Furthermore, φ is invariant under orthogonal transformations and constant on the orbits contained in $\mathcal{U}_{n,k}(L)/\mathfrak{S}_k$.

Proof. Let $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_k)$ and $\varphi(\mathbf{U}) = \mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_n)$. Since the rows of \mathbf{U} are binary unit vectors, we immediately get $\mathbf{W} \in \{0, 1\}^{n \times n}$ and $\text{diag}(\mathbf{W}) = \mathbf{e}$. Thus $\mathbf{W} \in \text{PM}_k^n$, since $\text{rank}(\mathbf{W}) = k$ and $\mathbf{W} \geq \mathbf{0}$ follow by definition from $\mathbf{W} = \mathbf{U}\mathbf{U}^\top$ and the pairwise orthogonality of the columns of \mathbf{U} . To see that $\mathbf{W} \in \text{PM}_k^n(L)$, note that $\mathbf{e}_i^\top \mathbf{U} = \mathbf{e}_j^\top$ for some $j \in [k]$ and thus

$$\mathbf{W}\mathbf{e}_i = \mathbf{U}\mathbf{U}^\top \mathbf{e}_i = \mathbf{U}\mathbf{e}_j \in V(L) \quad (5.1)$$

by assumption.

To see the invariance of φ under orthogonal transformations \mathbf{Q} , note that

$$\varphi(\mathbf{U}\mathbf{Q}) = \mathbf{U}\mathbf{Q}\mathbf{Q}^\top \mathbf{U}^\top = \mathbf{U}\mathbf{U}^\top,$$

which also shows that φ is constant on the orbits induced by column-permutations.

Injectivity of φ on $\mathcal{U}_{n,k}^{\text{lex}}(L)$ follows by noting that the sets

$$\{\mathbf{u}_1, \dots, \mathbf{u}_k\} = \{\mathbf{w}_1, \dots, \mathbf{w}_n\}$$

coincide due to (5.1) and $\mathbf{0} \notin \text{col}(\mathbf{U})$.

To see surjectivity of φ , let $\mathbf{W} \in \text{PM}_k^n(L)$ and consider the Gram representation $\mathbf{W} = \mathbf{U}_0 \mathbf{U}_0^\top$ with $\mathbf{U}_0 \in \mathbb{R}^{n \times k}$ and $\text{rank}(\mathbf{U}_0) = k$. Then there is an invertible $k \times k$ submatrix \mathbf{M} and we assume without loss of generality that

$$\mathbf{U}_0 = \begin{pmatrix} \mathbf{M} \\ \mathbf{U}'_0 \end{pmatrix}.$$

Then $\text{rank}(\mathbf{M}\mathbf{M}^\top) = \text{rank}(\mathbf{M}) = k$, and $\mathbf{M}\mathbf{M}^\top$ is binary with $\text{diag}(\mathbf{M}\mathbf{M}^\top) = \mathbf{e}$ as a submatrix of \mathbf{W} . Then $\mathbf{M}\mathbf{M}^\top = \mathbf{I}_k = \mathbf{M}^\top \mathbf{M}$ is orthogonal and

$$\mathbf{U}_1 := \begin{pmatrix} \mathbf{I}_k \\ \mathbf{U}'_1 \end{pmatrix} = \begin{pmatrix} \mathbf{M}\mathbf{M}^\top \\ \mathbf{U}'_0 \mathbf{M}^\top \end{pmatrix} = \mathbf{U}_0 \mathbf{M}^\top$$

shows that $\mathbf{U}_1 \mathbf{U}_1^\top = \mathbf{U}_0 \mathbf{M}\mathbf{M}^\top \mathbf{U}_0^\top = \mathbf{W}$ where \mathbf{U}_1 contains \mathbf{I}_k as submatrix where \mathbf{M} used to be. This shows \mathbf{U}_1 is binary, since

$$\begin{pmatrix} \mathbf{I}_k & \mathbf{U}'_1{}^\top \\ \mathbf{U}'_1 & \mathbf{U}'_1 \mathbf{U}'_1{}^\top \end{pmatrix} = \begin{pmatrix} \mathbf{I}_k \\ \mathbf{U}'_1 \end{pmatrix} \begin{pmatrix} \mathbf{I}_k \\ \mathbf{U}'_1 \end{pmatrix}^\top = \mathbf{U}_1 \mathbf{U}_1^\top = \mathbf{W} \in \{0, 1\}^{n \times n}.$$

Now (5.1) shows both $\text{col}(\mathbf{U}_1) \subseteq V(L)$ as well as the last equation in

$$[n] = \text{supp}(\text{diag}(\mathbf{W})) \subseteq \text{supp}(\mathbf{W}\mathbf{e}) = \text{supp}(\mathbf{U}_1 \mathbf{e}) \subseteq [n],$$

which implies $\mathbf{U}_1 \mathbf{e} = \mathbf{e}$. Then $\mathbf{U}_1 \in \mathcal{U}_{n,k}(L)$ and we are done. \square

As a consequence, we get a nice characterization of the set of partitions $\mathcal{P}^n(L)$ as

$$\text{PM}^n(L) = \bigcup_{k \in [n]} \text{PM}_k^n(L) = \{ \mathbf{W} \in \{0, 1\}^{n \times n} \mid \text{diag}(\mathbf{W}) = \mathbf{e}, \mathbf{W} \geq \mathbf{0}, \text{col}(\mathbf{W}) \subseteq V(L) \}.$$

5.2 Connection to Combinatorial Moment Matrices

At this point, the description of PM_k^n is difficult to work with due to both the binary constraint as well as the rank constraint, which both lead to NP-hardness in general. Since Theorem 4.5 showed that Max-Cut is NP-hard, so is optimizing over PM_k^n , since a Max-Cut instance with objective \mathbf{C} is equivalent to solving

$$\max \{ \langle \mathbf{C}, \mathbf{W} \rangle \mid \mathbf{W} \in \text{PM}_2^n \}.$$

As such, we can not expect to remove both of these constraints, although it is possible to remove the rank constraint as shown in the following lemma.

Lemma 5.5 ([Ren10, Lemma 18.4]):

The set of k -partition matrices has the explicit description

$$\begin{aligned} \text{PM}_k^n &= \{\mathbf{W} \in \{0, 1\}^{n \times n} \mid \text{diag}(\mathbf{W}) = \mathbf{e}, (t\mathbf{W} \geq J \Leftrightarrow t \geq k)\} \\ &= \left\{ \mathbf{W} \in \{0, 1\}^{n \times n} \mid \text{diag}(\mathbf{W}) = \mathbf{e}, \left[\begin{pmatrix} t & \mathbf{e}^\top \\ \mathbf{e} & \mathbf{W} \end{pmatrix} \geq \mathbf{0} \Leftrightarrow t \geq k \right] \right\}. \end{aligned} \quad (5.2)$$

This lemma is important, since it highlights a link to the theory of combinatorial moment matrices from Subsection 3.1. Lemma 5.4 shows that starting from an assignment matrix \mathbf{U} , we can get a partition matrix by setting

$$\mathbf{W} = \mathbf{U}\mathbf{U}^\top = \sum_{j \in [k]} \mathbf{u}_j \mathbf{u}_j^\top,$$

where each \mathbf{u}_j is the characteristic vector of a part of the partition. By adding a new entry 1, we can consider $\begin{pmatrix} 1 & \mathbf{u}_j^\top \end{pmatrix}$ as a truncated moment sequence up to degree 1, and the collection of these sequences yields an extended assignment matrix

$$\begin{pmatrix} \mathbf{e}^\top \\ \mathbf{U} \end{pmatrix} \in \{0, 1\}^{(n+1) \times k}$$

which factors into an extended partition matrix as

$$\begin{pmatrix} k & \mathbf{e}^\top \\ \mathbf{e} & \mathbf{W} \end{pmatrix} = \begin{pmatrix} k & \mathbf{e}^\top \mathbf{U}^\top \\ \mathbf{U} \mathbf{e} & \mathbf{U} \mathbf{U}^\top \end{pmatrix} = \begin{pmatrix} \mathbf{e}^\top \\ \mathbf{U} \end{pmatrix} (\mathbf{e} \quad \mathbf{U}^\top) = \sum_{j \in [k]} \begin{pmatrix} 1 & \mathbf{u}_j^\top \\ \mathbf{u}_j & \mathbf{u}_j \mathbf{u}_j^\top \end{pmatrix}.$$

The important observation here is that an extended k -partition matrix is equal to the sum of the k first-order combinatorial moment matrices corresponding to the individual sets in its k -partition.

Example 5.6:

Consider again the partition $\{\{1, 2\}, \{3\}\} \in \mathcal{P}_2^3$ from Example 4.9. Lexicographic order of the corresponding characteristic vectors leads to the assignment and partition matrices

$$\mathbf{U} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{W} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

The extended partition matrix can be decomposed as sum of first-order combinatorial moment matrices as

$$\left(\begin{array}{c|ccc} 2 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{array} \right) = \left(\begin{array}{c|ccc} 1 & 1 & 1 & 0 \\ \hline 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right) + \left(\begin{array}{c|ccc} 1 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \end{array} \right).$$

5.3 Convexification

In the end, we are interested in solving optimization problems over $\text{PM}_k^n(L)$, so we are naturally interested in its convex hull $\text{conv}(\text{PM}_k^n(L))$. However, as we have already seen, Max- k -Cut can be described as

$$\max \{ \langle \mathbf{C}, \mathbf{X} \rangle \mid \mathbf{X} \in \text{PM}_k^n \},$$

which is equivalent to

$$\max \{ \langle \mathbf{C}, \mathbf{X} \rangle \mid \mathbf{X} \in \text{conv}(\text{PM}_k^n) \},$$

due to the linearity of the objective. Consequently, we should not expect to find a compact description of $\text{conv}(\text{PM}_k^n(L))$, as it is likely out of reach.

In order to tackle this problem anyway, it is easiest to work with convex relaxations of $\text{conv}(\text{PM}_k^n)$. The standard convex relaxation of the set PM_k^n in literature arises from formulation (5.2) by relaxing the integrality constraint $\mathbf{W} \in \{0, 1\}^{n \times n}$ away. This leads us to the set

$$\text{RPM}_k^n(L) = \left\{ \mathbf{W} \in \mathbb{R}^{n \times n} \mid \text{diag}(\mathbf{W}) = \mathbf{e}, \begin{pmatrix} k & \mathbf{e}^\top \\ \mathbf{e} & \mathbf{W} \end{pmatrix} \geq \mathbf{0}, \langle \mathbf{W}, \Omega_L \rangle = 0 \right\}, \quad (5.3)$$

which is a convex relaxation of $\text{PM}_k^n(L)$. The constraint $\langle \mathbf{W}, \Omega_L \rangle = 0$ arises as a truncation of $\text{col}(\mathbf{W}) \subseteq V(L)$ and only excludes the sets of size 2 that are forbidden by L .

5.3.1 Applying the Method of Moments to Partition Matrices

A way to interpret Lemma 5.4 is to consider the set $\text{PM}_k^n(L)$ as the image of the second order moments of $\mathcal{U}_{n,k}(L)$ under a linear map. It follows then, since $\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))$ is a convex relaxation of the second order moments of $\mathcal{U}_{n,k}(L)$, that applying the same linear map to $\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))$ will lead us to a convex relaxation of $\text{PM}_k^n(L)$ as well. The goal of this section is to show that surprisingly, this approach yields a relaxation that is equivalent to the straightforward relaxation $\text{RPM}_k^n(L)$ from (5.3).

To see this, first recall the set of assignment matrices

$$\mathcal{U}_{n,k}(L) = \{ \mathbf{U} \in \{0, 1\}^{n \times k} \mid \mathbf{U}\mathbf{e}_{[k]} = \mathbf{e}_{[n]}, \text{col}(\mathbf{U}) \subseteq V(L) \}$$

from (4.3). In order to state the corresponding relaxation of $\text{PM}_k^n(L)$, we will first investigate the structure of $\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))$, where any assignment matrix $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_k)$ is represented as

$$\begin{pmatrix} 1 \\ \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_k \end{pmatrix} \begin{pmatrix} 1 \\ \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_k \end{pmatrix}^\top = \begin{pmatrix} 1 & \mathbf{u}_1^\top & \mathbf{u}_2^\top & \dots & \mathbf{u}_k^\top \\ \mathbf{u}_1 & \mathbf{u}_1\mathbf{u}_1^\top & \mathbf{u}_1\mathbf{u}_2^\top & \dots & \mathbf{u}_1\mathbf{u}_k^\top \\ \mathbf{u}_2 & \mathbf{u}_2\mathbf{u}_1^\top & \mathbf{u}_2\mathbf{u}_2^\top & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \mathbf{u}_{k-1}\mathbf{u}_k^\top \\ \mathbf{u}_k & \mathbf{u}_k\mathbf{u}_1^\top & \dots & \mathbf{u}_k\mathbf{u}_{k-1}^\top & \mathbf{u}_k\mathbf{u}_k^\top \end{pmatrix}. \quad (5.4)$$

For the rest of this section, let us denote matrices with this block structure as

$$\mathbf{U}(n, k) := \begin{pmatrix} 1 & \mathbf{u}_1^\top & \mathbf{u}_2^\top & \cdots & \mathbf{u}_k^\top \\ \mathbf{u}_1 & \mathbf{U}_{11} & \mathbf{U}_{12} & \cdots & \mathbf{U}_{1k} \\ \mathbf{u}_2 & \mathbf{U}_{21} & \mathbf{U}_{22} & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \mathbf{U}_{(k-1)k} \\ \mathbf{u}_k & \mathbf{U}_{k1} & \cdots & \mathbf{U}_{k(k-1)} & \mathbf{U}_{kk} \end{pmatrix} \in \mathbb{R}^{(1+nk) \times (1+nk)}. \quad (5.5)$$

In order to describe the set $\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))$, we use the representation

$$\mathcal{U}_{n,k}(L) = \mathcal{U}_{n,k} \cap \mathcal{V}_{\mathbb{R}}(L^k)$$

to work with

$$\mathcal{N}_1^*(\mathcal{U}_{n,k}(L)) = \mathcal{N}_1^*(\mathcal{U}_{n,k}) \cap \mathcal{N}_1^*(\mathcal{I}(L^k)).$$

Then $\mathcal{N}_1^*(\mathcal{U}_{n,k})$ can be explicitly described as

$$\mathcal{N}_1^*(\mathcal{U}_{n,k}) = \left\{ \mathbf{U}(n, k) \geq \mathbf{0} \mid \begin{array}{l} \sum_{j \in [k]} \mathbf{u}_j = \mathbf{e}, \\ \sum_{j \in [k]} \mathbf{U}_{ij} = \mathbf{u}_i \mathbf{e}^\top, \end{array} \quad \begin{array}{l} \langle \mathbf{u}_i, \mathbf{e} \rangle \geq 1 \\ \text{diag}(\mathbf{U}_{ii}) = \mathbf{u}_i \end{array} \forall i \in [k] \right\}. \quad (5.6)$$

Furthermore, it follows from Example 2.13 that

$$\mathcal{N}_1^*(\mathcal{I}(L^k)) = \{ \mathbf{U}(n, k) \geq \mathbf{0} \mid \langle \mathbf{U}_{ii}, \boldsymbol{\Omega}_L \rangle = 0 \forall i \in [k] \}, \quad (5.7)$$

whose block structure repeats with k .

Now if $\mathbf{U}(n, k)$ arises from a matrix $\mathbf{U} \in \mathcal{U}_{n,k}(L)$, then the linear map φ_2 given by

$$\varphi_2 : \mathbf{U}(n, k) \mapsto \sum_{j \in [k]} \mathbf{U}_{jj}$$

yields the same result as φ from Lemma 5.4, since under these circumstances,

$$\varphi_2(\mathbf{U}(n, k)) = \sum_{j \in [k]} \mathbf{U}_{jj} = \sum_{j \in [k]} \mathbf{u}_j \mathbf{u}_j^\top = \mathbf{U} \mathbf{U}^\top = \varphi(\mathbf{U}).$$

As a consequence, we get a convex relaxation of $\text{PM}_k^n(L)$ by noting that

$$\text{PM}_k^n(L) \subseteq \varphi_2(\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))).$$

The rest of this section is dedicated to show Theorem 5.8, which will state that already

$$\varphi_2(\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))) = \text{RPM}_k^n(L).$$

Due to the symmetry, we can construct multiple valid representations of (5.4) by permuting the columns of the assignment matrix \mathbf{U} before the construction. To see this, let the group \mathfrak{S}_k act on $\mathbb{R}^{(1+nk) \times (1+nk)}$ by setting

$$\pi * \mathbf{U}(n, k) = \begin{pmatrix} 1 & \mathbf{u}_{\pi(1)}^\top & \cdots & \mathbf{u}_{\pi(k)}^\top \\ \mathbf{u}_{\pi(1)} & \mathbf{U}_{\pi(1)\pi(1)} & \cdots & \mathbf{U}_{\pi(1)\pi(k)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{u}_{\pi(k)} & \mathbf{U}_{\pi(k)\pi(1)} & \cdots & \mathbf{U}_{\pi(k)\pi(k)} \end{pmatrix}$$

for each $\pi \in \mathfrak{S}_k$ and each $\mathbf{U}(n, k) \in \mathbb{R}^{(1+nk) \times (1+nk)}$. It then follows from the symmetrical nature of $\mathcal{U}_{n,k}(L)$ that whenever $\mathbf{U}(n, k)$ arises from $\mathbf{U} \in \mathcal{U}_{n,k}(L)$ as in (5.4), the matrix $\pi * \mathbf{U}(n, k)$ arises from $\pi(\mathbf{U})$. This naturally leads to $k!$ distinct representations $\pi * \mathbf{U}(n, k)$ of each partition, which makes the set $\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))$ difficult to work with.

However, we can reduce it to its symmetry invariant subspace. To this end, let us introduce the linear *Reynolds operator* $\psi: \mathbb{R}^{(1+nk) \times (1+nk)} \rightarrow \mathbb{R}^{(1+nk) \times (1+nk)}$ given as

$$\psi(\mathbf{U}(n, k)) := \frac{1}{k!} \sum_{\pi \in \mathfrak{S}_k} (\pi * \mathbf{U}(n, k)) = \begin{pmatrix} 1 & \mathbf{a}^\top & \cdots & \cdots & \mathbf{a}^\top \\ \mathbf{a} & \mathbf{A} & \mathbf{B} & \cdots & \mathbf{B} \\ \vdots & \mathbf{B} & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \mathbf{B} \\ \mathbf{a} & \mathbf{B} & \cdots & \mathbf{B} & \mathbf{A} \end{pmatrix} =: \mathbf{U}(\mathbf{A}, \mathbf{B}, \mathbf{a}),$$

where

$$\begin{aligned} \mathbf{a} &= \frac{1}{k!} \sum_{\pi \in \mathfrak{S}_k} \mathbf{u}_{\pi(i)} &= \frac{1}{k} \sum_{j \in [k]} \mathbf{u}_j & \quad \forall i \in [k], \\ \mathbf{A} &= \frac{1}{k!} \sum_{\pi \in \mathfrak{S}_k} \mathbf{U}_{\pi(j)\pi(j)} &= \frac{1}{k} \sum_{j \in [k]} \mathbf{U}_{jj} & \quad \forall i \in [k], \\ \mathbf{B} &= \frac{1}{k!} \sum_{\pi \in \mathfrak{S}_k} \mathbf{U}_{\pi(i)\pi(j)} &= \frac{1}{k(k-1)} \sum_{i, j \in [k], i \neq j} \mathbf{U}_{\pi(i)\pi(j)} & \quad \forall i, j \in [k]. \end{aligned}$$

For later, note that we get the identity

$$\varphi_2(\mathbf{U}(n, k)) = \sum_{j \in [k]} \mathbf{U}_{jj} = k\mathbf{A} = \varphi_2(\mathbf{U}(\mathbf{A}, \mathbf{B}, \mathbf{a})), \quad (5.8)$$

which shows that $\varphi_2 \circ \psi = \varphi_2$ as a side result.

Denoting the linear subspace given by the image of ψ by $H_\psi := \varphi(\mathbb{R}^{(1+nk) \times (1+nk)})$, it has the explicit description

$$H_\psi = \{ \mathbf{U}(\mathbf{A}, \mathbf{B}, \mathbf{a}) \mid \mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}, \mathbf{a} \in \mathbb{R}^n \},$$

as

$$\pi * \mathbf{U}(\mathbf{A}, \mathbf{B}, \mathbf{a}) = \mathbf{U}(\mathbf{A}, \mathbf{B}, \mathbf{a})$$

holds for all $\pi \in \mathfrak{S}_k$ and all $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, $\mathbf{a} \in \mathbb{R}^n$. The intersection $H_\psi \cap \mathcal{S}_+^{(1+nk)}$ is particularly nice, as is shown by the following lemma.

Lemma 5.7 ([GL08, Lemma 2.8]):

Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ and

$$\mathbf{X} = \mathbf{I}_k \otimes \mathbf{A} + \mathbf{J}_k \otimes \mathbf{B} = \begin{pmatrix} \mathbf{A} + \mathbf{B} & \mathbf{B} & \dots & \mathbf{B} \\ \mathbf{B} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{B} \\ \mathbf{B} & \dots & \mathbf{B} & \mathbf{A} + \mathbf{B} \end{pmatrix} \in \mathbb{R}^{kn \times kn}.$$

Then $\mathbf{X} \geq \mathbf{0}$ if and only if $\mathbf{A} \geq \mathbf{0}$ and $\mathbf{A} + k\mathbf{B} \geq \mathbf{0}$.

We are now ready to show the central result of this section.

Theorem 5.8:

It holds that

$$\varphi_2(\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))) = \text{RPM}_k^n(L).$$

Proof. We first argue that

$$\psi(\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))) = \mathcal{N}_1^*(\mathcal{U}_{n,k}(L)) \cap H_\psi.$$

For " \subseteq ", note that $\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))$ is convex and that by definition, ψ maps matrices $\mathbf{U}(n, k) \in \mathcal{N}_1^*(\mathcal{U}_{n,k}(L))$ into the convex hull of their orbits, so

$$\psi(\mathcal{N}_1^*(\mathcal{U}_{n,k}(L))) \subseteq \mathcal{N}_1^*(\mathcal{U}_{n,k}(L)).$$

For " \supseteq ", note that ψ acts as the identity on H_ψ , and so $\mathbf{U}(n, k) \in \mathcal{N}_1^*(\mathcal{U}_{n,k}(L)) \cap H_\psi$ implies $\mathbf{U}(n, k) = \psi(\mathbf{U}(n, k)) \in \psi(\mathcal{N}_1^*(\mathcal{U}_{n,k}(L)))$.

It follows then from (5.6) that

$$\psi(\mathcal{N}_1^*(\mathcal{U}_{n,k})) = \left\{ \mathbf{U}(\mathbf{A}, \mathbf{B}, \frac{1}{k}\mathbf{e}) \geq \mathbf{0} \mid \mathbf{A} + (k-1)\mathbf{B} = \frac{1}{k}\mathbf{J}_n, \text{diag}(\mathbf{A}) = \frac{1}{k}\mathbf{e} \right\}.$$

Using the Schur complement Lemma 2.5 on $\mathbf{U}(\mathbf{A}, \mathbf{B}, \frac{1}{k}\mathbf{e})$ shows

$$\mathbf{U}(\mathbf{A}, \mathbf{B}, \frac{1}{k}\mathbf{e}) \geq \mathbf{0} \Leftrightarrow \mathbf{I}_k \otimes (\mathbf{A} - \mathbf{B}) + \mathbf{J}_k \otimes (\mathbf{B} - \frac{1}{k^2}\mathbf{J}_n) \geq \mathbf{0},$$

which is equivalent to

$$\mathbf{A} - \mathbf{B} \geq \mathbf{0}, (\mathbf{A} - \mathbf{B}) + k(\mathbf{B} - \frac{1}{k^2}\mathbf{J}_n) \geq \mathbf{0}$$

by Lemma 5.7. Since $\mathbf{A} + (k-1)\mathbf{B} = \frac{1}{k}\mathbf{J}_n$ holds by assumption, the second psd. condition is satisfied and solving for \mathbf{B} yields

$$\mathbf{B}(\mathbf{A}) := \frac{1}{k-1} \left(\frac{1}{k}\mathbf{J}_n - \mathbf{A} \right),$$

which can be substituted into $\mathbf{A} - \mathbf{B} \geq \mathbf{0}$ to yield $\mathbf{A} \geq \frac{1}{k^2}\mathbf{J}_n$ after rescaling.

We can thus simplify

$$\begin{aligned} \psi(\mathcal{N}_1^*(\mathcal{U}_{n,k})) &= \left\{ \mathbf{U}(\mathbf{A}, \mathbf{B}(\mathbf{A}), \frac{1}{k}\mathbf{e}) \mid \begin{pmatrix} 1 & \frac{1}{k}\mathbf{e}^\top \\ \frac{1}{k}\mathbf{e} & \mathbf{A} \end{pmatrix} \geq \mathbf{0}, \text{diag}(\mathbf{A}) = \frac{1}{k}\mathbf{e} \right\} \\ &= \left\{ \mathbf{U}(\mathbf{A}, \mathbf{B}(\mathbf{A}), \frac{1}{k}\mathbf{e}) \mid k\mathbf{A} \in \text{RPM}_k^n \right\}, \end{aligned}$$

and (5.8) shows the claim for $L = 2^{\lfloor n \rfloor}$.

To finish the proof, we need to show that this structural result is unchanged by introducing a non-trivial L . But this follows by definition, since adding an independence system L only introduces the constraint $\langle \mathbf{A}, \Omega_L \rangle = 0$ in both sets, as can be seen by comparing (5.3) and (5.7). \square

This shows that one way to understand the approximation quality of $\text{RPM}_k^n(L)$ is the fact that it reduces to computing the first stage of MM for $\mathcal{U}_{n,k}(L)$.

5.3.2 Applying the Method of Moments to Orbitopes

As we have seen in the preceding section in Theorem 5.8, the symmetry of $\mathcal{U}_{n,k}(L)$ leads to a compact formulation of the first stage of MM, and in turn to the compact relaxation $\text{RPM}_k^n(L)$ for $\text{PM}_k^n(L)$. While this is good in terms of computation, the question is whether this also leads to a decrease in quality of the solutions, as is often the case with relaxations defined on the symmetry-invariant subspace. To this end, the next logical step is to apply MM to the symmetry-free $\mathcal{U}_{n,k}^{\text{lex}}(L)$ instead, and ask if the image under φ_2 leads to tighter relaxations than $\text{RPM}_k^n(L)$.

In general, this would be an ambitious approach, since applying MM to the extended formulation of $\mathcal{U}_{n,k}^{\text{lex}}(L)$ does not lead to sets with intuitive descriptions. Fortunately, we can explicitly analyse this construction for the case of $k = 2$, due to the simplified description of $\mathcal{U}_{n,2}^{\text{lex}}(L)$ in Theorem 4.13.

Remark 5.9:

Although $k = 2$ may seem limiting at first, it can actually be considered as the hardest setting by noting that good approximation ratios for Max- k -Cut are available for growing k through the following, simple heuristic:

Choosing a partition matrix $\mathbf{W} \in \text{PM}_k^n$ uniformly at random, the probability of $w_{ij} = 1$ equals $\frac{1}{k}$ for any pair $i \neq j \in [n]$. This translates to an expected approximation

ratio of $1 - \frac{1}{k}$ for Max- k -Cut given as

$$\max \{ \langle \mathbf{C}, \mathbf{W} \rangle \mid \mathbf{W} \in \text{PM}_k^n \}, \quad (5.9)$$

which is strictly increasing for growing k .

Before we go on, it is important to point out that in the case of $L \neq 2^{[n]}$, we can simplify the problem drastically.

Lemma 5.10:

Let $L \neq 2^{[n]}$, $G_L \in \mathcal{G}_n$ be the graph whose adjacency matrix is given by Ω_L and let $m = |\text{CC}(G_L)|$ be the number of its connected components. Then

$$\mathcal{U}_{n,2}(L) \leftrightarrow \mathcal{U}_{m,2}^0 := \{ \mathbf{U} \in \{0, 1\}^{m \times 2} \mid \mathbf{U}\mathbf{e}_{[2]} = \mathbf{e}_{[m]} \}.$$

Proof. By the central assumption that $\mathcal{P}_2^n(L) \neq \emptyset$, G_L is bipartite and can be decomposed into the union of its bipartite connected components. Then for $j \in [m]$, the j -th connected component has a unique partition $\{T_1^j, T_2^j\}$ with $T_1^j \neq \emptyset$, which can be found by the following well-known algorithm:

Starting with any node i of the component, assume $i \in T_1^j$ to conclude for its neighbourhood $N(i)$ that $N(i) \subseteq T_2^j$. But then $N(N(i)) \subseteq T_1^j$, and iterating this process at most n steps, we cover all nodes of the component and arrive at a full list for $\{T_1^j, T_2^j\}$.

Now any element of $\{T_1, T_2\} \in \mathcal{P}_2^n(L)$ is completely described by an assignment

$$\{T_1^j \mid j \in [m]\} \rightarrow \{T_1, T_2\},$$

since for $\{a, b\} = [2]$, we have $T_1^j \subseteq T_a$ if and only if $T_2^j \subseteq T_b$. Thus $\mathcal{U}_{n,2}(L)$ can be represented by all $m \times 2$ assignment matrices, where we allow for $\mathbf{0}$ columns. \square

Recall that according to (5.6) and (5.7), the first stage of MM for $\mathcal{U}_{n,2}(L)$ is

$$\mathcal{N}_1^*(\mathcal{U}_{n,2}) = \left\{ \mathbf{U}(n, 2) \left| \begin{array}{ll} \mathbf{u}_1 + \mathbf{u}_2 = \mathbf{e}, & \mathbf{U}_{i1} + \mathbf{U}_{i2} = \mathbf{u}_i \mathbf{e}^\top, \\ \text{diag}(\mathbf{U}_{11}) = \mathbf{u}_1, & \text{diag}(\mathbf{U}_{22}) = \mathbf{u}_2, \\ \langle \mathbf{u}_1, \mathbf{e} \rangle \geq 1, & \langle \mathbf{u}_2, \mathbf{e} \rangle \geq 1, \\ \mathbf{U}(n, 2) \geq \mathbf{0} \end{array} \right. \right\},$$

where the inequalities $\langle \mathbf{u}_i, \mathbf{e} \rangle \geq 1$ are the only addition to $\mathcal{N}_1^*(\mathcal{U}_{n,2})$ over $\mathcal{N}_1^*(\mathcal{U}_{n,2}^0)$.

Theorem 5.11:

For all $n \in \mathbb{N}$, it holds that

$$\mathcal{N}_1^*(\mathcal{U}_{n,2}^{0,lex}) \cong \mathcal{N}_1^*(\mathcal{U}_{n-1,2}^0),$$

where

$$\mathcal{U}_{n,2}^{0,lex} = \{ \mathbf{U} \in \mathcal{U}_{n,2}^0 \mid u_{11} = 1 \}.$$

Proof. Let $\mathbf{U} \in \mathcal{N}_1^*(\mathcal{U}_{n,2}^{0,lex})$. Then $\mathbf{U} \in \mathcal{N}_1^*(\mathcal{U}_{n,2}^0)$ and by definition of MM, we have the additional constraints $(\mathbf{u}_1)_1 = 1$ and $\mathbf{U}_{i1}\mathbf{e}_1 = \mathbf{u}_i$ for $i = 1, 2$. In particular,

$$\mathbf{U} = \left(\begin{array}{c|c|c} 1 & \mathbf{u}_1^\top & \mathbf{u}_2^\top \\ \hline \mathbf{u}_1 & \mathbf{U}_{11} & \mathbf{U}_{12} \\ \hline \mathbf{u}_2 & \mathbf{U}_{21} & \mathbf{U}_{22} \end{array} \right) = \left(\begin{array}{c|c|c|c|c} 1 & 1 & \mathbf{u}'_1{}^\top & 0 & \mathbf{u}'_2{}^\top \\ \hline 1 & 1 & \mathbf{u}'_1{}^\top & 0 & \mathbf{u}'_2{}^\top \\ \hline \mathbf{u}'_1 & \mathbf{u}'_1 & \mathbf{U}'_{11} & \mathbf{0} & \mathbf{U}'_{12} \\ \hline 0 & 0 & \mathbf{0} & 0 & \mathbf{0} \\ \hline \mathbf{u}'_2 & \mathbf{u}'_2 & \mathbf{U}'_{21} & \mathbf{0} & \mathbf{U}'_{22} \end{array} \right), \quad (5.10)$$

where the blocks of both matrices correspond to each other. Then we can reduce $\mathbf{U} \geq \mathbf{0}$ to $\mathbf{U}' \geq \mathbf{0}$, where

$$\mathbf{U}' = \begin{pmatrix} 1 & \mathbf{u}'_1{}^\top & \mathbf{u}'_2{}^\top \\ \mathbf{u}'_1 & \mathbf{U}'_{11} & \mathbf{U}'_{12} \\ \mathbf{u}'_2 & \mathbf{U}'_{21} & \mathbf{U}'_{22} \end{pmatrix},$$

due to Corollary 3.3.

But then $\mathbf{U}' \in \mathcal{N}_1^*(\mathcal{U}_{n-1,2}^0)$, which can be verified by applying the equations from $\mathcal{N}_1^*(\mathcal{U}_{n,2}^0)$ on the blocks given by the right hand side of (5.10). \square

While it is true that from a structural point of view, we could now use Theorem 5.8 on the reformulation given by Lemma 5.10 to end up with RPM_2^{m-1} , the reduction breaks down when we consider the objective function as well. When formulating (5.9) in the setting of $\mathcal{N}_1^*(\mathcal{U}_{n,2})$, we get the objective matrix

$$\left(\begin{array}{c|c|c} 0 & \mathbf{0} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{W} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} & \mathbf{W} \end{array} \right) = \left(\begin{array}{c|c|c|c|c} 0 & 0 & \mathbf{0} & 0 & \mathbf{0} \\ \hline 0 & w_{11} & \mathbf{w}'_1{}^\top & 0 & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{w}'_1 & \mathbf{W}'_{11} & \mathbf{0} & \mathbf{0} \\ \hline 0 & 0 & \mathbf{0} & w_{11} & \mathbf{w}'_1{}^\top \\ \hline \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{w}'_1{}^\top & \mathbf{W}'_{11} \end{array} \right) \mapsto \mathbf{W}' = \begin{pmatrix} w_{11} & \mathbf{w}'_1{}^\top & \mathbf{0} \\ \mathbf{w}'_1 & \mathbf{W}'_{11} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{W}'_{11} \end{pmatrix},$$

for which we cannot find a new objective \mathbf{W}'' such that

$$\langle \mathbf{U}', \mathbf{W}' \rangle = \langle \varphi_2(\mathbf{U}'), \mathbf{W}'' \rangle.$$

This shows that applying MM to orbitopes does indeed break the symmetry and prevents a compact formulation similar to RPM_2^{m-1} , though at the cost of an SDP of twice the size. As a corollary, the same is true for the case when $L = 2^{\lfloor n \rfloor}$, with the minor addition of the inequality $\langle \mathbf{u}_2, \mathbf{e} \rangle \geq 1$, since $\langle \mathbf{u}_1, \mathbf{e} \rangle \geq 1$ is satisfied by assumption.

Corollary 5.12:

For all $n \in \mathbb{N}$, it holds that

$$\mathcal{N}_1^*(\mathcal{U}_{n,2}^{lex}) \cong \left\{ \mathbf{U}(n-1, 2) \left| \begin{array}{l} \mathbf{u}_1 + \mathbf{u}_2 = \mathbf{e}, \quad \mathbf{U}_{i1} + \mathbf{U}_{i2} = \mathbf{u}_i \mathbf{e}^\top, \\ \text{diag}(\mathbf{U}_{11}) = \mathbf{u}_1, \quad \text{diag}(\mathbf{U}_{22}) = \mathbf{u}_2, \\ \langle \mathbf{u}_2, \mathbf{e} \rangle \geq 1, \quad \mathbf{U}(n-1, 2) \geq \mathbf{0} \end{array} \right. \right\}.$$

Chapter 6

Projection Matrices

So far, we have seen that the symmetrical nature of assignment matrices $\mathcal{U}_{n,k}(L)$ has both led to meaningless solutions as well as more compact formulations of otherwise large problems. The reason for this is the 'global' representation of the partition:

Given the i -th row of an assignment matrix, we only know the label of the part that contains i , but not how this set actually looks like. Since we need the complete assignment matrix in order to reconstruct the part, in a sense, this information is locally unknown.

The idea of this chapter is to make this information locally available - if we assign an actual description of the part containing i to i instead of a label, symmetry will cease to be an issue, although at the price of a redundant description.

To this end, we will express the partitions as projection matrices in Section 6.1, which will have nice properties concerning the hierarchy constructed by MM in Section 6.2. In Section 6.3, we will apply this representation to both the graph-colouring and Euclidean k -clustering problems. As a main result of the thesis, we give a new formulation of the graph colouring number $\chi(G)$ and relate its relaxations to well-studied ones from literature.

6.1 Overview

Throughout this chapter, let $\mathbf{R} = (\mathbf{r}_1, \dots, \mathbf{r}_n) = \{r_{ij}\}_{i,j \in [n]} \in \mathbb{R}^{n \times n}$. We first recall some basic properties about projection matrices arising from orthogonal projections onto subspaces.

Definition 6.1:

Let $\mathbf{X} \in \mathbb{R}^{n \times k}$ be a matrix with full column rank to define its corresponding *projection matrix* as

$$\rho(\mathbf{X}) := \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \in \mathbb{R}^{n \times n}.$$

The set of symmetric projection matrices of size n is defined as

$$\text{SProM}^n := \{\mathbf{R} \in \mathbb{R}^{n \times n} \mid \mathbf{R}^2 = \mathbf{R}, \mathbf{R} = \mathbf{R}^\top\} \subseteq \mathcal{S}_+^n.$$

Note that by assumption, $\mathbf{X}^\top \mathbf{X}$ is invertible and $\rho(\mathbf{X}) \in \text{SProM}^n$. We proceed by recalling some fundamental properties of projection matrices.

Lemma 6.2:

Let $\mathbf{R} \in \text{SProM}^n$ and $\mathbf{X} \in \mathbb{R}^{n \times k}$ with full column rank. Then the following holds:

- (i) Binary eigenvalues: $\text{spct}(\mathbf{R}) \subseteq \{0, 1\}$,
- (ii) Rank equation: $\text{tr}(\mathbf{R}) = \text{rank}(\mathbf{R})$,
- (iii) Projection property: $\rho(\mathbf{X})\mathbf{X} = \mathbf{X}$, $\rho(\mathbf{X}) \cdot \mathbb{R}^n = \mathbf{X} \cdot \mathbb{R}^k$,
- (iv) Isometry invariance: $\rho(\mathbf{X}) = \rho(\mathbf{X}\mathbf{Q})$ for all orthogonal $\mathbf{Q} \in \mathbb{R}^{n \times n}$.

Proof. (i) follows from the Cayley-Hamilton theorem applied to the matrix polynomial $\mathbf{R}^2 = \mathbf{R}$. For (ii), note that $\text{tr}(\mathbf{R})$ is equivalent to the sum of its eigenvalues, which count the rank due to (i). (iii) and (iv) follow directly from the definition. \square

Remark 6.3:

The set of projection matrices can be bijectively mapped to the set of all linear subspaces of \mathbb{R}^n by assigning a projection matrix to its image. In particular, the matrix $\rho(\mathbf{X})$ maps to the column space spanned by \mathbf{X} as can be seen by Lemma 6.2.

More specifically, we can recover the Grassmannian $\text{Gr}(k, \mathbb{R}^n)$ as a real variety by intersecting SProM^n with the affine hyperplane

$$\{\mathbf{R} \in \mathbb{R}^{n \times n} \mid \text{tr}(\mathbf{R}) = k\}$$

and evoking the rank equation in Lemma 6.2.

Example 6.4:

Recall Example 5.6, where we looked at the assignment matrix \mathbf{U} and corresponding partition matrix $\mathbf{U}\mathbf{U}^\top$ of the partition $\{\{1, 2\}, \{3\}\}$ given by

$$\mathbf{U} = \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{pmatrix} \text{ and } \mathbf{U}\mathbf{U}^\top = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Applying ρ to \mathbf{U} leads us to the corresponding projection matrix

$$\rho(\mathbf{U}) = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}^{-1} \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Denoting the columns as $\rho(\mathbf{U}) = (\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$, we observe for the following that

- $\mathbf{U}^\top \mathbf{U}$ is a diagonal matrix which contains the size of parts of the partition,
- $\rho(\mathbf{U})$ is a rescaled version of the partition matrix $\mathbf{U}\mathbf{U}^\top$,
- the columns $\{\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3\}$ consist of all eigenvectors of $\rho(\mathbf{U})$ with repetition,
- $\{\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3\} \subseteq \Delta^3$,
- $\|\mathbf{r}_i\|_0 \cdot \|\mathbf{r}_i\|_\infty = 1$ holds for each column $i \in [3]$,
- $\mathbf{r}_i = \mathbf{r}_j$ if and only if $(\mathbf{r}_i)_j = (\mathbf{r}_j)_i > 0$.

The preceding example motivates the following definition of a special class of doubly stochastic matrices.

Definition 6.5:

The set of *combinatorial projection matrices* $\text{CProM}_k^n(L)$ is given as

$$\text{CProM}_k^n(L) = \left\{ \mathbf{R} \in \mathbb{R}^{n \times n} \left| \begin{array}{l} \text{col}(\mathbf{R}) \subseteq \Delta_L^n \\ r_{ij} \cdot (\mathbf{r}_i - \mathbf{r}_j) = \mathbf{0} \\ \mathbf{R} = \mathbf{R}^\top \\ \text{tr}(\mathbf{R}) = k \end{array} \right. \forall i, j \in [n] \right\},$$

where

$$\Delta_L^n := \{\mathbf{r} \in \Delta^n \mid \text{supp}(\mathbf{r}) \in L\}.$$

The quadratic equations

$$r_{ij} \cdot (\mathbf{r}_i - \mathbf{r}_j) = \mathbf{0} \quad \forall i, j \in [n]$$

will be called *block-inducing* in the following.

Unsurprisingly, this set turns out to be exactly the projection matrices corresponding to assignment matrices, as shown by the following two lemmas.

Lemma 6.6:

Combinatorial projection matrices are projection matrices

$$\text{CProM}_k^n(L) \subseteq \text{SProM}^n$$

and have strictly positive diagonal $\text{diag}(\mathbf{R}) > 0$. In particular, each column \mathbf{r}_i is uniquely determined by its support via

$$r_{ji} = \delta_{j \in \text{supp}(\mathbf{r}_i)} \frac{1}{\|\mathbf{r}_i\|_0}$$

and satisfies the equation

$$\|\mathbf{r}_i\|_0 \cdot \|\mathbf{r}_i\|_\infty = 1.$$

Proof. Let $\mathbf{R} \in \text{CProM}_k^n(L)$ and choose any $i, j, l \in [n]$. By assumption, the block-inducing equations $r_{il} \cdot (\mathbf{r}_i - \mathbf{r}_l) = \mathbf{0}$ and $r_{ij} \cdot (\mathbf{r}_i - \mathbf{r}_j) = \mathbf{0}$ show the identity

$$r_{il}r_{jl} = r_{il}r_{ij} = r_{ij}r_{il} = r_{ij}r_{jl},$$

using the symmetry of \mathbf{R} . Then $\mathbf{R}^2 = \mathbf{R}$ follows from

$$(\mathbf{R}^2)_{ij} = \sum_{l \in [n]} r_{il}r_{jl} = \sum_{l \in [n]} r_{ij}r_{jl} = r_{ij} \left(\sum_{l \in [n]} r_{jl} \right) = r_{ij} = (\mathbf{R})_{ij} \quad \forall i, j \in [n].$$

Furthermore, the block-inducing equations $r_{ij} \cdot (\mathbf{r}_i - \mathbf{r}_j) = \mathbf{0}$ show $r_{ii}(r_{ij} - r_{ii}) = 0$ for all $j \in [n]$, which implies $r_{ij} \in \{0, r_{ii}\}$. Since $\mathbf{r}_i \in \Delta^n$, we necessarily have $r_{ii} \cdot \|\mathbf{r}_i\|_0 = 1$, which finishes the proof. \square

Theorem 6.7:

The combinatorial projection matrices are precisely the projection matrices corresponding to assignment matrices:

$$\rho(\mathcal{U}_{n,k}(L)) = \text{CProM}_k^n(L) \subseteq \text{SProM}^n.$$

In particular, ρ is a bijection that shows

$$\mathcal{U}_{n,k}^{\text{lex}}(L) \leftrightarrow \text{CProM}_k^n(L).$$

Proof. We first show $\rho(\mathcal{U}_{n,k}(L)) \subseteq \text{CProM}_k^n(L)$, so let $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_k) \in \mathcal{U}_{n,k}(L)$. We already know $\rho(\mathbf{U}) \in \text{SProM}^n$ and the corresponding properties from Lemma 6.2, so we only need to check the first two properties outlined in Definition 6.5.

For the first, assume $(\mathbf{u}_j)_i = 1$ to see

$$\rho(\mathbf{U})\mathbf{e}_i = \mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{e}_i = \mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{e}_j = \mathbf{U} \frac{\mathbf{e}_j}{\langle \mathbf{u}_j, \mathbf{u}_j \rangle} = \frac{\mathbf{u}_j}{\langle \mathbf{e}, \mathbf{u}_j \rangle} \in \Delta_L^n,$$

where we used that \mathbf{u}_j is binary in the second to last equation. For later, note the implication that

$$\{\text{supp}(\rho(\mathbf{U})\mathbf{e}_i) \mid i \in [n]\} = \{\text{supp}(\mathbf{U}\mathbf{e}_j) \mid j \in [k]\}. \quad (6.1)$$

For the second, assume $\mathbf{e}_i^\top \rho(\mathbf{U})\mathbf{e}_j > 0$, since $\rho(\mathbf{U}) \geq 0$ and there is nothing to show otherwise. Assuming $(\mathbf{u}_{j'})_j = 1$ and $(\mathbf{u}_{i'})_i = 1$, we have

$$0 < \mathbf{e}_i^\top \rho(\mathbf{U})\mathbf{e}_j = \mathbf{e}_i^\top \mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{e}_j = \mathbf{e}_{i'}^\top (\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{e}_{j'} = \delta_{i',j'} \cdot (\mathbf{U}^\top \mathbf{U})_{i',i'}^{-1},$$

since $(\mathbf{U}^\top \mathbf{U})^{-1}$ is diagonal. It follows that $i' = j'$, which in turn means

$$\mathbf{U}^\top \mathbf{e}_i = \mathbf{e}_{i'} = \mathbf{e}_{j'} = \mathbf{U}^\top \mathbf{e}_j.$$

In particular,

$$\rho(\mathbf{U})(\mathbf{e}_i - \mathbf{e}_j) = \mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1} (\mathbf{U}^\top (\mathbf{e}_i - \mathbf{e}_j)) = \mathbf{0},$$

which we wanted to show.

The reverse inclusion $\rho(\mathcal{U}_{n,k}(L)) \supseteq \text{CProM}_k^n(L)$ follows if we can show that ρ is a bijection from $\mathcal{U}_{n,k}^{\text{lex}}(L)$ to $\text{CProM}_k^n(L)$, since $\rho(\mathcal{U}_{n,k}(L)) = \rho(\mathcal{U}_{n,k}^{\text{lex}}(L))$ by the isometry invariance of Lemma 6.2.

To see that ρ is injective, note that the invariant (6.1) immediately shows that whenever $\mathbf{U}, \mathbf{U}' \in \mathcal{U}_{n,k}^{\text{lex}}(L)$ are assignment matrices of distinct partitions $\mathcal{T} \neq \mathcal{T}'$, then also $\rho(\mathbf{U}) \neq \rho(\mathbf{U}')$.

To show surjectivity, we argue that for $\mathbf{R} \in \text{CProM}_k^n(L)$, the map

$$\mathbf{R} \rightarrow \mathcal{T}(\mathbf{R}) =: \{\text{supp}(\mathbf{r}_i) \mid i \in [n]\} \in \mathcal{P}_k^n(L)$$

is well-defined and injective, since then the corresponding assignment-matrices define a proper inverse $\rho^{-1}(\mathbf{R})$.

Since $\mathbf{r}_i \in \Delta_L^n$, we have $\mathcal{T}(\mathbf{R}) \subseteq L_*$. Disjointness of the sets in $\mathcal{T}(\mathbf{R})$ follows from the implications of the block-inducing equations, and coverage of $[n]$ follows due to \mathbf{R} being doubly stochastic and so $\mathcal{T}(\mathbf{R}) \in \mathcal{P}_k^n(L)$ is well-defined.

Finally, the map $\mathbf{R} \rightarrow \mathcal{T}(\mathbf{R})$ is injective since \mathbf{R} can be uniquely reconstructed from $\mathcal{T}(\mathbf{R})$ by starting with $i \in \text{supp}(\mathbf{r}_i)$ and applying Lemma 6.6. \square

6.2 Convexification

As in the case of partition matrices, we are interested in the complexity of optimizing a linear function over $\text{CProM}_k^n(L)$. Unfortunately, Section 6.3 shows examples where this is NP-hard, so we can not expect a algorithmically exploitable description. For this reason, we again revert to studying relaxations of the convex hull.

The easiest convex relaxation for $\text{CProM}_k^n(L)$ is given by replacing the non-linear constraints

$$r_{ij} \cdot (\mathbf{r}_i - \mathbf{r}_j) = \mathbf{0} \quad \forall i, j \in [n]$$

with a psd. constraint $\mathbf{R} \geq \mathbf{0}$ emerging from the fact that since \mathbf{R} is a projection matrix, we have $\mathbf{R}^2 = \mathbf{R}$. We thus get the set

$$\text{RProM}_k^n(L)^0 = \{\mathbf{R} \in \mathbb{R}_+^{n \times n} \mid \text{tr}(\mathbf{R}) = k, \mathbf{R}\mathbf{e} = \mathbf{e}, \langle \mathbf{R}, \Omega_L \rangle = 0, \mathbf{R} \geq \mathbf{0}\}.$$

Just like in the preceding chapters, the idea now is to take the variety $\text{CProM}_k^n(L)$ and apply MM to it. For this reason, we will need to investigate the monomial structure of the underlying ideal of $\text{CProM}_k^n(L)$, which we will call \mathcal{I}_ρ and is defined as

$$\mathcal{I}_\rho = \left\langle \{r_{ij} - r_{ji}, \text{tr}(\mathbf{R}) - k, \langle \mathbf{r}_i, \mathbf{e} \rangle - 1, r_{ij}(r_{il} - r_{jl}) \mid \forall i, j, l \in [n]\} \right\rangle \subseteq \mathbb{R}[\mathbf{R}].$$

Unfortunately, the ideal \mathcal{I}_ρ is an example where stating a reduced Gröbner basis and a corresponding basis for $\mathbb{R}[\mathbf{R}]/\mathcal{I}_\rho$ explicitly becomes difficult and would distract from the underlying problem structure. For this reason, we will only state a generating system instead of a basis.

Lemma 6.8:

A monomial generating system of $\mathbb{R}_t[\mathbf{R}]/\mathcal{I}_\rho$ can be indexed by the set

$$\text{ProMo}_{n,t} = \left\{ (G, \psi) \in \mathcal{G}_n \times \mathbb{N}^{\text{CC}(G)} \mid \frac{1}{2}|E(G)| + \sum_{T \in \text{CC}(G)} \psi(T) = t \right\}.$$

Proof. Let $\mathbf{R}^\Psi \in \mathbb{R}_t[\mathbf{R}]$, so in particular $\Psi \in \mathbb{N}_t^{n \times n}$. We will show that each equivalence class $[\mathbf{R}^\Psi] \in \mathbb{R}_t[\mathbf{R}]/\mathcal{I}_\rho$ can be identified with a pair (G, ψ) as described in the set above. To simplify notation, we will note that the congruence $\mathbf{R}^\Psi \equiv \mathbf{R}^{\Psi'} \pmod{\mathcal{I}_\rho}$ can equivalently be expressed as a congruence on the exponent matrices $\Psi \cong \Psi'$, whose integer entries describe the powers of the corresponding variables.

First, since the columns of \mathbf{R} are normalized, we have the inclusion

$$\mathbb{R}[\mathbf{R}]/\mathcal{I}_\rho \subseteq \prod_{i \in [n]} \mathbb{R}[\mathbf{r}_i]/\mathcal{I}(\Delta^n)$$

and we can use homogenization (2.5) on any column \mathbf{r}_i to see that \mathbf{R}^Ψ can be written as a sum of several $\mathbf{R}^{\Psi'}$, where $\Psi' \in \mathbb{N}^{n \times n}$ and $\langle \Psi', \mathbf{J}_n \rangle = t$.

Second, we can find $\Psi' \cong \Psi''$ where Ψ'' is upper triangular, since $r_{ij} = r_{ji}$, and so everything can be expressed in terms of r_{ij} with $i \leq j$.

Third, we can use the identities $r_{ij}(r_{il} - r_{jl}) = 0$ with $l = j$ to see that

$$r_{jj} \cdot r_{ij} = r_{ij}^2 = r_{ii} \cdot r_{ij}, \quad (6.2)$$

which shows that we can decrease the off-diagonal entries of Ψ exceeding 1 to increase the main-diagonal of the corresponding column or row. We thus can find upper-triangular $\Psi''' \cong \Psi''$ such that

$$\Psi''' \in \{0, 1\}^{n \times n} + \text{Diag}(\mathbb{N}^n), \quad \langle \Psi''', \mathbf{J}_n \rangle = t,$$

e.g. all off-diagonal exponents are binary.

In this form, the off-diagonal entries of Ψ''' define the adjacency matrix of a directed graph, which we can identify with its underlying undirected graph G . We define a map $\psi \in \mathbb{N}^{\text{CC}(G)}$ by setting

$$\psi(T) = \sum_{i \in T} \Psi_{ii} \quad \forall T \in \text{CC}(G).$$

Then (6.2) shows that the construction of ψ is independent on how we choose to construct Ψ''' from Ψ'' , since we can decrease any Ψ'''_{ii} by an integer to increase Ψ'''_{jj} by the

same amount whenever the edge (i, j) is contained in G without leaving its equivalence class. In particular, this extends to any connecting path between i and j , such that the sum of the diagonal entries over a connected component are invariant. We can thus identify Ψ' with the couple (G, ψ) as claimed, allowing us to write any \mathbf{R}^Ψ as a linear combination of such monomials. \square

One can verify that $\text{ProMo}_{n,1} \cong \{(i, j) \in [n]^2 \mid i \leq j\}$, as graphs containing a single edge (i, j) can be identified with that edge, and the connected components of edgeless graphs can be identified with their vertices i to yield (i, i) . We thus recover the monomials $\{r_{ij} \mid i \leq j\}$ we started with.

Remark 6.9:

The monomial generating system indexed by $\text{ProMo}_{n,t}$ in Lemma 6.8 is not isomorphic to a basis for $\mathbb{R}_t[\mathbf{R}]/\mathcal{I}_\rho$, since further reductions via the unused equation $\text{tr}(\mathbf{R}) - k$ are possible. We neglect any such simplifications for the sake of clarity.

While considering $\mathbb{R}_t[\mathbf{R}]$ with symmetry equations would lead to $z_{\binom{n}{2}}(t)$ monomials, the reduced system indexed by $\text{ProMo}_{n,t}$ grows much slower, as can be seen in Table 6.1.

	1	2	3	4	5		1	2	3	4	5
5	15	100	390	1000	1863	5	16	136	816	3876	15504
6	21	201	1156	4451	12230	6	22	253	2024	12650	65780
7	28	364	2905	15876	62972	7	29	435	4495	35960	237336

Table 6.1: Left: $|\text{ProMo}_{n,t}|$ indexed by pairs (n, t) . Right: Values $\dim(\mathbb{R}_t[\mathbf{R}])$ indexed by pairs (n, t) .

Table 6.1 shows that the number of monomials for approximating $\text{CProM}_k^n(L)$ grows much too fast. Unfortunately, even the very first stage of MM needs second order moments, which still grow too fast even for small inputs. To remedy this, we can relax the first stage of MM by discarding some of the off-diagonal entries. One way to do this is by starting with all monomials in $\mathbb{R}[\mathbf{R}]$, discarding off-diagonal entries to end up with a block-diagonal structure and then reducing the remaining matrices by using a smaller monomial base suited for $\mathbb{R}[\mathbf{R}]/\mathcal{I}_\rho$.

The main matrix in $\mathcal{N}_1^*(\text{CProM}_k^n(L))$ is given by

$$\mathbf{M}_1(\mathbf{R}) = \begin{pmatrix} 1 & \mathbf{r}_1^\top & \dots & \mathbf{r}_n^\top \\ \mathbf{r}_1 & \mathbf{R}_{11} & \dots & \mathbf{R}_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{r}_n & \mathbf{R}_{n1} & \dots & \mathbf{R}_{nn} \end{pmatrix} \geq \mathbf{0}.$$

The easiest way to relax this condition is to discard all matrices \mathbf{R}_{ij} with $i \neq j$, which relaxes $\mathbf{M}_1(\mathbf{R}) \geq \mathbf{0}$ to the system

$$\mathbf{M}_1(\mathbf{r}_i) = \begin{pmatrix} 1 & \mathbf{r}_i^\top \\ \mathbf{r}_i & \mathbf{R}_{ii} \end{pmatrix} \geq \mathbf{0} \quad \forall i \in [n].$$

Invoking the results of Lemma 6.8, we see that we can discard the first row and column from this condition. Applying the remaining constraints for each \mathbf{R}_{ii} , we end up with feasible sets

$$R_i = \left\{ \mathbf{R}_{ii} \in \mathbb{R}_+^{n \times n} \mid \langle \mathbf{R}_{ii}, \mathbf{J}_n \rangle = 1, \mathbf{R}_{ii} \mathbf{e}_i = \text{diag}(\mathbf{R}_{ii}), \mathbf{R}_{ii} \geq \mathbf{0} \right\},$$

where the condition $\mathbf{R}_{ii} \mathbf{e}_i = \text{diag}(\mathbf{R}_{ii})$ breaks the symmetry between different $i \in [n]$. By invoking the block inducing equations, one can show that additionally, the matrices \mathbf{R}_{ii} are linked through the constraints

$$\sum_{i \in [n]} \mathbf{R}_{ii} = \mathbf{R}$$

and

$$(\mathbf{R}_{ii})_{jl} = (\mathbf{R}_{jj})_{il} = (\mathbf{R}_{ll})_{ij},$$

which implies that we can consider each matrix \mathbf{R}_{ii} as a slice of a symmetrical third-order tensor, as shown in Figure 6.1.

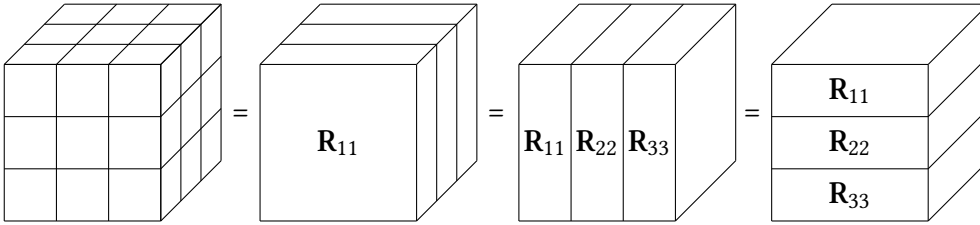


Figure 6.1: Different ways to align the variables of the various \mathbf{R}_{ii} into the same symmetrical third order tensor, each cube representing one variable.

Denoting this structure with the set

$$R = \left\{ \mathbf{R} \in \mathbb{R}^{n \times n} \mid \mathbf{R} = \sum_{i \in [n]} \mathbf{R}_{ii}, (\mathbf{R}_{ii})_{jl} = (\mathbf{R}_{jj})_{il}, \mathbf{R}_{ii} \in R_i \quad \forall i, j, l \in [n] \right\},$$

a relaxation of the first stage of MM is given by

$$\text{RProM}_k^n(L) = \{ \mathbf{R} \in R \mid \text{tr}(\mathbf{R}) = k, \mathbf{R} \mathbf{e} = \mathbf{e}, \langle \mathbf{R}, \mathbf{\Omega}_L \rangle = 0 \},$$

where $\text{RProM}_k^n(L)^0$ relaxes this set by using $R \subseteq \mathbb{R}_+^{n \times n} \cap \mathcal{S}_+^n$.

Remark 6.10:

The approach of constructing a symmetrical third-order tensor to relax a moment matrix of order 1 was already proposed in [GL08] to relax a combinatorial moment matrix of order 1 for the stable set problem. In particular, each entry of our tensor can be identified with a set $\{i, j, l\}$ via its indices, which can in turn be understood as a rescaling of combinatorial moments up to the third order.

While this process can be generalized to turn higher order combinatorial moment matrices of rank 1 into symmetrical tensors of appropriate order, it does not work in our setting. The main reason is the assumption $\text{rank}(\mathbf{R}) = k > 1$, which destroys the symmetry; in general, the equation $(\mathbf{R}_{ij})_{lm} = (\mathbf{R}_{il})_{jm}$ will not hold, so we can not derive unique entries for the entries associated with $\{i, j, l, m\}$.

6.3 Applications

6.3.1 Graph Colouring

We will shortly recall the basics about graph colouring before we show how to use projection matrices for a new formulation of this problem. For an extensive survey on convex relaxations for graph colouring, consider [GL08].

Stable Sets

Definition 6.11:

A *stable set* of a graph $G \in \mathcal{G}_n$ is a subset $S \subseteq [n]$ such that the subgraph induced by S does not contain any edges. The set of stable sets in G is denoted by $S_G \subseteq 2^{[n]}$. The *stable set number* $\alpha(G)$ is the biggest size $|S|$ of any stable set in G and defines a function $\alpha : \mathcal{G}_n \rightarrow \mathbb{N}$. The *clique number* $\omega(G)$ of G is the biggest size $|S|$ of a stable set in the complement graph G , so $\omega(G) := \alpha(\overline{G})$.

Formally, we can express $\alpha(G)$ as the solution to an integer problem by working with characteristic vectors $\mathbf{x} \in \{0, 1\}^n$ via

$$\begin{aligned} \alpha(G) &= \max \{ \langle \mathbf{x}, \mathbf{e} \rangle \mid x_i x_j = 0 \quad \forall (i, j) \in E, \mathbf{x} \in \{0, 1\}^n \} \\ &= \max \{ \langle \mathbf{x}, \mathbf{e} \rangle \mid \mathbf{x}^\top \mathbf{A}_G \mathbf{x} = 0, \mathbf{x} \in \{0, 1\}^n \}, \end{aligned}$$

where \mathbf{A}_G denotes the adjacency matrix of G . For later, we will note the following crucial property of S_G .

Lemma 6.12:

The stable sets S_G form an independence system.

Proof. Follows directly by checking the properties in the definition (2.2). □

Graph Colourings

Definition 6.13:

A k -colouring of a graph $G = ([n], E)$ is a k -partition $\mathcal{T} = (T_1, \dots, T_k)$ of the node set $[n]$, such that each T_i is a stable set in G . The *chromatic number* $\chi(G)$ is the smallest value k such that G has a k -colouring.

As shown back in Example 4.7, computing $\chi(G)$ is a minimum cover problem, but it can not be modelled with assignment matrices alone. While there is a model involving assignment matrices, we will consider another standard formulation in terms of extended partition matrices instead.

$$\begin{aligned} \chi(G) &= \min \{k \in \mathbb{N} \mid \mathbf{X} \in \text{PM}_k^n(S_G)\} \\ &= \min \left\{ k \in \mathbb{R} \mid \begin{pmatrix} k & \mathbf{e}^\top \\ \mathbf{e} & \mathbf{X} \end{pmatrix} \geq \mathbf{0}, \langle \mathbf{X}, \mathbf{A}_G \rangle = 0, \text{diag}(\mathbf{X}) = \mathbf{e}, \mathbf{X} \in \{0, 1\}^{n \times n} \right\} \end{aligned} \quad (6.3)$$

Unfortunately, we have the following hardness result.

Theorem 6.14 ([Sch03, BS94]):

Computing $\alpha(G)$ or $\chi(G)$ is NP-hard. Furthermore, it is NP-hard to approximate $\chi(G)$ within $n^{1/14-\epsilon}$ for any $\epsilon > 0$.

In spite of being NP-hard to compute in general, $\chi(G)$ behaves much nicer restricted to special classes of graphs with regards to the following construction.

Definition 6.15 (Lovász theta number [Lov79]):

The *Lovász theta number* $\vartheta(G)$ is the optimal value of the following primal/dual SDP problems

$$\begin{aligned} \max \{ \langle \mathbf{X}, \mathbf{J}_n \rangle \mid \langle \mathbf{X}, \mathbf{I}_n \rangle = 1, \langle \mathbf{X}, \mathbf{A}_G \rangle = 0, \mathbf{X} \geq \mathbf{0} \}, & \quad (\vartheta\text{-P}) \\ \min \left\{ k \mid \text{diag}(\mathbf{Y}) = \mathbf{e}, \langle \mathbf{Y}, \mathbf{A}_G \rangle = 0, \begin{pmatrix} k & \mathbf{e}^\top \\ \mathbf{e} & \mathbf{Y} \end{pmatrix} \geq \mathbf{0} \right\}. & \quad (\vartheta\text{-D}) \end{aligned}$$

Theorem 6.16 (Lovász sandwich Theorem [Lov87]):

For all $G \in \mathcal{G}_n$, we have

$$\omega(G) \leq \vartheta(\overline{G}) \leq \chi(G).$$

This result shines in the context of *perfect graphs*, where $\omega(H) = \chi(H)$ for all induced subgraphs $H \subseteq G$, and $\chi(G)$ can be computed in polynomial time. For more on the theory of perfect graphs, consider the recent survey [Tro15].

In spite of the hardness results in Theorem 6.14, a lot of work has been done to improve $\vartheta(G)$ to get better bounds for imperfect graphs. We exemplarily show the following variants taken from [DR07] and [GL08].

Nonnegativity

Adding $X \geq \mathbf{0}$ to (ϑ -P) and $Y \geq \mathbf{0}$ to (ϑ -D) results in the Schrijver number $\vartheta^-(G)$ [Sch79] and the Szegedy number $\vartheta^+(G)$ [Sze94] respectively.

Applying MM

Following Section 5.3, we can interpret $\text{PM}_k^n(S_G)$ in (6.3) as variety in terms of $\mathcal{U}_{n,k}(S_G)$, and MM can be applied as shown before. Denoting the corresponding optimal value of the t -th stage as $\psi^{(t)}(G)$, we get a converging hierarchy of lower bounds for $\chi(G)$.

Note that we actually have $\psi^{(1)}(G) = \vartheta(\overline{G})$, since the feasible set of formulation (ϑ -D) is precisely the relaxation $\text{RPM}_k^n(S_{\overline{G}})$ in this case. In particular, the Lovász theta number arises from the underlying MM hierarchy, and so MM can be seen as a natural generalization.

Summarizing the results so far, [GL08] shows

$$\psi^{(1)}(G) = \vartheta(\overline{G}) \leq \vartheta^+(\overline{G}) \leq \psi^{(2)}(G) \leq \chi(G).$$

There are several other relaxations for $\chi(G)$ available in [GL08], where among others the computationally tractable $\psi(G)$ is introduced.

The projection ϑ -number

We now propose a new formulation in terms of projection matrices and relate it to $\vartheta(G)$.

Theorem 6.17:

We have the following characterization of the chromatic number:

$$\begin{aligned} \chi(G) &= \min \{k \mid \text{CProM}_k^n(S_G) \neq \emptyset\} \\ &= \min \left\{ \text{tr}(\mathbf{R}) \left| \begin{array}{l} \text{col}(\mathbf{R}) \subseteq \Delta^n \\ r_{ij} = 0 \quad \forall (i, j) \in E(G) \\ r_{ij} \cdot (\mathbf{r}_i - \mathbf{r}_j) = \mathbf{0} \quad \forall i, j \in [n] \\ \mathbf{R} = \mathbf{R}^\top \end{array} \right. \right\}. \end{aligned}$$

Proof. The first equation follows by definition, so we will show the second. Let $\chi'(G)$ denote the optimal value of the second optimization problem and let \mathcal{T} be a minimal colouring with $\chi(G)$ colours. Then the projection matrix $\mathbf{R}(\mathcal{T})$ corresponding to \mathcal{T} is feasible for the second optimization problem, since $\text{supp}(\mathbf{R}(\mathcal{T})\mathbf{e}_j) \in S_G$ and $\text{diag}(\mathbf{R}(\mathcal{T})) > \mathbf{0}$ imply $\mathbf{e}_i^\top \mathbf{R}(\mathcal{T})\mathbf{e}_j = 0$ whenever $(i, j) \in E(G)$, and so $\chi'(G) \leq \chi(G)$.

To see the other inequality, first note that $\chi'(G) \geq 0$ is integral, since whenever \mathbf{R} is feasible, $\mathbf{R} \in \text{SProm}^n$ as shown in the proof of Lemma 6.6, and so the eigenvalues of

\mathbf{R} are binary. It suffices to show that $\mathbf{R} \in \text{CProM}_{\text{tr}(\mathbf{R})}^n(S_G)$ then, and the only thing that is left to prove is $\text{supp}(\mathbf{r}_i) \in S_G$.

Suppose that $\text{supp}(\mathbf{r}_i) \notin S_G$, so there is $(j, l) \in E(G)$ such that $(\mathbf{r}_i)_j, (\mathbf{r}_i)_l > 0$. It follows from $r_{ij} \cdot (\mathbf{r}_i - \mathbf{r}_j) = \mathbf{0}$ that $(\mathbf{r}_j)_l = (\mathbf{r}_i)_l > 0$, which contradicts $(\mathbf{r}_j)_l = 0$. \square

The importance of the preceding theorem lies in the fact that we made the constraint $\text{supp}(\mathbf{r}_i) \in S_G$ tractable, so that we can use the hierarchy given by MM to approximate $\chi(G)$. Following Section 6.2, we can immediately give the 0-th stage of our relaxations.

Definition 6.18:

The projection ϑ -number is given as

$$\begin{aligned} \hat{\vartheta}(G) &:= \min \{k \in \mathbb{R}_+^n \mid \text{RProM}_k^n(S_G)^0 \neq \emptyset\} \\ &= \min \{\text{tr}(\mathbf{R}) \mid \mathbf{R} \geq \mathbf{0}, \mathbf{R}\mathbf{e} = \mathbf{e}, \langle \mathbf{R}, \mathbf{A}_{\overline{G}} \rangle = 0, \mathbf{R} \geq \mathbf{0}\}. \end{aligned} \quad (6.4)$$

We should note that by non-negativity of \mathbf{R} , we have

$$\langle \mathbf{R}, \mathbf{A}_{\overline{G}} \rangle = 0 \quad \Leftrightarrow \quad r_{ij} = 0 \quad \forall (i, j) \in E(G),$$

and we will prefer this notation. For the following, it will be important to explicitly write down the Szegedy number mentioned before, which is given as

$$\vartheta^+(G) := \min \left\{ x_0 \mid \begin{pmatrix} x_0 & \mathbf{e}^\top \\ \mathbf{e} & \mathbf{X} \end{pmatrix} \geq \mathbf{0}, \text{diag}(\mathbf{X}) = \mathbf{e}, \langle \mathbf{X}, \mathbf{A}_{\overline{G}} \rangle = 0, \mathbf{X} \geq \mathbf{0} \right\}. \quad (6.5)$$

Since both bounds have similar computational complexity, it makes sense to compare their quality. Unfortunately, it turns out that the classical approach is always at least as good as using projection matrices, as shown in the following theorem.

Theorem 6.19:

Let $G \in \mathcal{G}_n$, then $\vartheta^+(G) \geq \hat{\vartheta}(G)$.

Proof. Let $x_0 = \vartheta^+(G)$ and \mathbf{X} be an optimal solution to (6.5). Since \mathbf{X} is symmetrical and non-negative, it follows from [BPS66] that there exists a diagonal matrix $\mathbf{D} = \text{Diag}(\mathbf{d})$ with $\mathbf{d} > \mathbf{0}$ such that

$$\begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{pmatrix} \begin{pmatrix} \vartheta^+(G) & \mathbf{e}^\top \\ \mathbf{e} & \mathbf{X} \end{pmatrix} \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{pmatrix} = \begin{pmatrix} \vartheta^+(G) & \mathbf{d}^\top \\ \mathbf{d} & \mathbf{DXD} \end{pmatrix} \geq \mathbf{0},$$

where $\mathbf{R} := \mathbf{DXD}$ is doubly stochastic, so $\mathbf{R}\mathbf{e} = \mathbf{e}$. In particular, \mathbf{R} is feasible for (6.4), and we have $\text{tr}(\mathbf{R}) = \|\mathbf{d}\|_2^2$. Using positive semidefiniteness, for any $\mu \in \mathbb{R}$, we have

$$0 \leq \begin{pmatrix} -1 \\ \mu \mathbf{d} \end{pmatrix}^\top \begin{pmatrix} \vartheta^+(G) & \mathbf{d}^\top \\ \mathbf{d} & \mathbf{R} \end{pmatrix} \begin{pmatrix} -1 \\ \mu \mathbf{d} \end{pmatrix} = \vartheta^+(G) - 2\mu\|\mathbf{d}\|_2^2 + \mu^2(\mathbf{d}^\top \mathbf{R} \mathbf{d})$$

which becomes

$$\begin{aligned}\vartheta^+(G) &\geq 2\mu\|\mathbf{d}\|_2^2 - \mu^2(\mathbf{d}^\top \mathbf{R} \mathbf{d}) = \|\mathbf{d}\|_2^2 \left(2\mu - \mu^2 \left(\frac{\mathbf{d}^\top \mathbf{R} \mathbf{d}}{\|\mathbf{d}\|_2} \frac{\mathbf{d}}{\|\mathbf{d}\|_2} \right) \right) \\ &\geq \|\mathbf{d}\|_2^2 (2\mu - \mu^2) = \text{tr}(\mathbf{R}) (2\mu - \mu^2).\end{aligned}$$

Since $\max \{2\mu - \mu^2 \mid \mu \in \mathbb{R}\} = 1$, it follows that \mathbf{R} is a feasible solution to (6.4) that satisfies $\vartheta^+(G) \geq \text{tr}(\mathbf{R})$ and the theorem follows. \square

A crucial question in understanding the projection-model is the magnitude of the gap $\vartheta^+(G) - \hat{\vartheta}(G)$. It turns out that we can state an explicit asymptotic lowerbound.

Theorem 6.20:

The worst case gap $\vartheta^+(G) - \hat{\vartheta}(G)$ has asymptotic behaviour

$$\max \left\{ \vartheta^+(G) - \hat{\vartheta}(G) \mid G \in \mathcal{G}_n \right\} = \Omega(n).$$

In particular,

$$\limsup_{n \rightarrow \infty} \max \left\{ \frac{\vartheta^+(G) - \hat{\vartheta}(G)}{n} \mid G \in \mathcal{G}_n \right\} \geq \left(\frac{3}{\sqrt{2}} - 2 \right) \approx 0,1213.$$

Proof. We will explicitly construct a graph family for which the bound on the gap holds true asymptotically. To this end, consider for any two integers $n_1, n_2 \in \mathbb{N}$, the graph

$$G(n_1, n_2) := K_{n_1} \cup K_{n_2},$$

which we define as the union of two complete graphs with n_1 and n_2 nodes respectively. In particular, the corresponding adjacency matrix is given as

$$\mathbf{A}_{G(n_1, n_2)} = \begin{pmatrix} \mathbf{J}_{n_1} - \mathbf{I}_{n_1} & \mathbf{0} \\ \mathbf{0} & \mathbf{J}_{n_2} - \mathbf{I}_{n_2} \end{pmatrix},$$

and since these graphs are perfect, we explicitly have

$$\chi(G(n_1, n_2)) = \vartheta^+(\overline{G(n_1, n_2)}) = \omega(G(n_1, n_2)) = \max(n_1, n_2).$$

For $\hat{\vartheta}(\overline{G(n_1, n_2)})$, we will construct the optimal solution \mathbf{R} in closed form. Since this graph has multiple symmetries, we can assume \mathbf{R} to be symmetry invariant and parametrize the feasible set with only three parameters α, β, γ by setting

$$\mathbf{R} = \begin{pmatrix} \alpha \mathbf{I}_{n_1} & \beta \mathbf{J}_{n_1, n_2} \\ \beta \mathbf{J}_{n_2, n_1} & \gamma \mathbf{I}_{n_2} \end{pmatrix}.$$

Now that the constraint $\langle \mathbf{R}, \mathbf{A}_{G(n_1, n_2)} \rangle = 0$ is satisfied, the remaining affine constraints turn into

$$\begin{aligned} \mathbf{R} \geq \mathbf{0} &\Leftrightarrow \alpha, \beta, \gamma \geq 0, \\ \mathbf{R}\mathbf{e} = \mathbf{e} &\Leftrightarrow \alpha + n_2\beta = 1, \quad \gamma + n_1\beta = 1. \end{aligned}$$

Assuming $\alpha, \gamma \neq 0$ to make the condition $\mathbf{R} \geq \mathbf{0}$ non-trivial, we can use the Schur complement Lemma 2.5 to rewrite

$$\begin{aligned} \mathbf{R} \geq \mathbf{0} &\Leftrightarrow \gamma \mathbf{I}_{n_2} - \frac{\beta^2}{\alpha} \mathbf{J}_{n_2, n_1} \cdot \mathbf{J}_{n_1, n_2} \geq \mathbf{0} \Leftrightarrow \gamma \mathbf{I}_{n_2} - \frac{\beta^2}{\alpha} n_1 \mathbf{J}_{n_2} \geq \mathbf{0} \\ &\Leftrightarrow \gamma \geq \frac{\beta^2}{\alpha} n_1 n_2, \end{aligned}$$

where we used the fact that \mathbf{J}_n only has one non-zero eigenvalue given by n . Lastly, we can explicitly express the objective function as

$$\text{tr}(\mathbf{R}) = \alpha n_1 + \gamma n_2 = n_1 + n_2 - 2n_1 n_2 \beta,$$

using the affine constraints. Ignoring the constants, the resulting problem of computing the number $\hat{\vartheta}(\overline{G(n_1, n_2)})$ is equivalent to

$$\max \left\{ \beta \mid \alpha + n_2\beta = 1, \quad \gamma + n_1\beta = 1, \quad \alpha\gamma \geq \beta^2 n_1 n_2, \quad \alpha, \beta, \gamma \geq 0 \right\}.$$

Using the equations, one can show that the unique solution is

$$\begin{pmatrix} \alpha \\ \beta \\ \gamma \end{pmatrix} = \frac{1}{n_1 + n_2} \begin{pmatrix} n_1 \\ 1 \\ n_2 \end{pmatrix}$$

and

$$\hat{\vartheta}(\overline{G(n_1, n_2)}) = \frac{n_1^2 + n_2^2}{n_1 + n_2}.$$

In particular, we now have the gap

$$\Delta(n_1, n_2) := \vartheta^+(\overline{G(n_1, n_2)}) - \hat{\vartheta}(\overline{G(n_1, n_2)}) = \max(n_1, n_2) - \frac{n_1^2 + n_2^2}{n_1 + n_2}. \quad (6.6)$$

W.l.o.g., let $n_1 = m, n_2 = \mu m \in \mathbb{N}$ for some $\mu \in [0, 1]$. Then (6.6) reads

$$\Delta(m, \mu m) = m - \frac{(1 + \mu^2)m^2}{(1 + \mu)m} = \mu \left(\frac{1 - \mu}{1 + \mu} \right) m.$$

Finally, optimizing the choice $\mu \in [0, 1]$ yields the biggest theoretical gap for $\mu = \sqrt{2} - 1$ and

$$\max_{\mu \in [0, 1]} \Delta(m, \mu m) = (3 - 2\sqrt{2})m.$$

For growing m , we can approximate this gap arbitrarily well by choosing $n_1 = m$ and $n_2 = \lceil (\sqrt{2} - 1)m \rceil$ to get a graph of size $\lceil \sqrt{2}m \rceil$ with asymptotic relative gap

$$\limsup_{m \rightarrow \infty} \frac{\Delta(m, \lceil (\sqrt{2} - 1)m \rceil)}{\lceil \sqrt{2}m \rceil} = \frac{3 - 2\sqrt{2}}{\sqrt{2}} \approx 0,1213.$$

□

As a side effect of the preceding theorem, we get the following corollary, which might explain the discrepancy by relating $\hat{\vartheta}(G)$ to Theorem 6.16.

Corollary 6.21:

The inequality $\omega(G) \leq \hat{\vartheta}(\bar{G})$ does not hold in general. In particular, $\hat{\vartheta}(G)$ is not necessarily exact for perfect graphs G .

Even though the preceding results show severe disadvantages over the classical formulation, we can recover some useful properties for the class of vertex-transitive graphs. We first cite the property in question.

Theorem 6.22 ([Sze94]):

For all $G \in \mathcal{G}_n$, the inequality

$$\vartheta^+(G) \cdot \vartheta^-(\bar{G}) \geq n$$

holds, with equality if G is vertex-transitive.

Surprisingly, we can use $\hat{\vartheta}(G)$ to sharpen this inequality.

Lemma 6.23:

For all $G \in \mathcal{G}_n$, the inequality $\hat{\vartheta}(G) \cdot \vartheta^-(\bar{G}) \geq n$ holds.

Proof. Let \mathbf{R} be an optimal solution to (6.4). Recalling the definition

$$\vartheta^-(\bar{G}) = \max \{ \langle \mathbf{X}, \mathbf{J} \rangle \mid \text{tr}(\mathbf{X}) = 1, \langle \mathbf{X}, \mathbf{A}_{\bar{G}} \rangle = 0, \mathbf{X} \geq \mathbf{0}, \mathbf{X} \geq \mathbf{0} \}$$

as (ϑ -P) with non-negativity constraints, we see that $\mathbf{X} := \frac{1}{\text{tr}(\mathbf{R})}\mathbf{R}$ is a feasible solution and as desired,

$$\vartheta^-(\bar{G}) \geq \langle \mathbf{X}, \mathbf{J} \rangle = \frac{\langle \mathbf{R}, \mathbf{J} \rangle}{\text{tr}(\mathbf{R})} = \frac{n}{\hat{\vartheta}(G)}.$$

□

Theorem 6.24:

For vertex-transitive $G \in \mathcal{G}_n$, we have

$$\vartheta^+(G) = \hat{\vartheta}(G).$$

In particular, $\omega(G) = \hat{\vartheta}(\bar{G}) = \chi(G)$ for vertex-transitive perfect $G \in \mathcal{G}_n$.

Proof. Using Theorem 6.19, Lemma 6.23 and Theorem 6.22, we see that

$$\vartheta^+(G) \geq \hat{\vartheta}(G) \geq \frac{n}{\vartheta^-(G)} = \vartheta^+(G)$$

whenever $G \in \mathcal{G}_n$ is vertex-transitive. \square

Since vertex-transitive graphs are known to be examples for which the relaxations $\vartheta(G)$ and $\vartheta^+(G)$ perform badly, it is surprising to see that we can recover an analogue of Theorem 6.16 for $\hat{\vartheta}(G)$ in this case. The implication is that the advantage of $\vartheta(G)$ over $\hat{\vartheta}(G)$ is related to exploiting the lacking symmetries of a given graph.

Moving up in the hierarchy

We close this section by noting that even when we refine $\hat{\vartheta}(G)$ by using the computationally expensive $\text{RProM}_k^n(S_G)$ instead of $\text{RProM}_k^n(S_G)^0$, we still cannot guarantee a lower-bound of $\vartheta(G)$. To this end, let

$$\hat{\vartheta}'(G) := \min \{k \in \mathbb{R}_+^n \mid \text{RProM}_k^n(S_G) \neq \emptyset\}$$

be the corresponding strengthening of $\hat{\vartheta}(G)$.

Empirical evidence points to the fact that $\hat{\vartheta}'(G)$ yields the correct answer for the graph class $G(n_1, n_2)$ used as counter example in Theorem 6.20. However, we can easily generalize this counter example to find another class of perfect graphs for which we get results where $\vartheta(G) > \hat{\vartheta}'(G)$. Let $n_1, n_2, n_3 \in \mathbb{N}$ such that

$$G(n_1, n_2, n_3) = K_{n_1} \cup K_{n_2} \cup K_{n_3}$$

is the union of three complete graphs. Again, this class of graphs is perfect, and as such we get

$$\vartheta(\overline{G(n_1, n_2, n_3)}) = \omega(G(n_1, n_2, n_3)) = \max\{n_1, n_2, n_3\}.$$

Continuing in this fashion, we also define the class of perfect graphs given by the union of a complete graph on n_1 nodes together with m isolated nodes, or more formally

$$G(n_1, \mathbf{e}_m) = K_{n_1} \cup \left(\bigcup_{i \in [m]} K_1 \right),$$

where

$$\vartheta(\overline{G(n_1, \mathbf{e}_m)}) = \omega(G(n_1, \mathbf{e}_m)) = n_1.$$

Fixing the number of total nodes to 9, the tables in Figure 6.2 and 6.3 show how the various relaxations behave. While there is a definite increase in quality by going from

n_1	n_2	n_3	$\hat{\vartheta}$	$\hat{\vartheta}'$	$\vartheta = \chi$
3	3	3	3	3	3
4	3	2	3.222	3.968	4
4	4	1	3.666	4	4
5	2	2	3.666	4.972	5
5	3	1	3.888	4.983	5
6	2	1	4.555	5.983	6
7	1	1	5.666	6.985	7

Table 6.2: Relaxations for $G(n_1, n_2, n_3)$.

n_1	m	$\hat{\vartheta}$	$\hat{\vartheta}'$	$\vartheta = \chi$
2	7	1.222	1.772	2
3	6	1.666	2.792	3
4	5	2.333	3.851	4
5	4	3.222	4.905	5
6	3	4.333	5.951	6
7	2	5.666	6.986	7
8	1	7.222	8	8

Table 6.3: Relaxations for $G(n_1, e_m)$.

$\hat{\vartheta}(G)$ to $\hat{\vartheta}'(G)$, the difference between $\hat{\vartheta}'(G)$ and $\vartheta(G)$ grows roughly linear with the number of connected components.

The main reason for this bad performance of relaxations based on projection matrices is the fact that the relaxations are *not monotone* in terms of subgraphs. In particular, if H is a subgraph of G , then we have the implication

$$H \leq G \quad \Rightarrow \quad \vartheta(H) \leq \vartheta(G),$$

which is not true for either $\hat{\vartheta}(G)$ or $\hat{\vartheta}'(G)$. Unfortunately, it is not straightforward to add this property to the functions in question, leaving them at a inherent disadvantage compared to the original ϑ -function. In particular, this begs the question if in general, approximating matrices with binary eigenvalues is harder than approximating matrices with binary entries.

6.3.2 Euclidean k -Clustering

Surprisingly, projection matrices also occur naturally in *Euclidean clustering*, where the goal is to cluster points in euclidean space into sets such that their deviation from a central point, or centroid, is minimal. More formally, we have the following definition.

Definition 6.25:

Let $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n) \in \mathbb{R}^{d \times n}$ and define the corresponding *squared error set-function* $\text{SE}_{\mathbf{B}} : 2^{[n]} \rightarrow \mathbb{R}$ as

$$\text{SE}_{\mathbf{B}}(T) = \min \left\{ \sum_{i \in T} \|\mathbf{b}_i - \mathbf{x}\|_2^2 \mid \mathbf{x} \in \mathbb{R}^d \right\}.$$

Then Euclidean k -clustering is the separable partition problem given by

$$\min \left\{ \sum_{T \in \mathcal{T}} \text{SE}_{\mathbf{B}}(T) \mid \mathcal{T} \in \mathcal{P}_k^n \right\}. \quad (6.7)$$

Unfortunately, Euclidean k -clustering is known to be NP-hard [ADHP09]. We cite the following, stronger result regarding its approximation.

Theorem 6.26 ([ACKS15]):

There exists a constant $\varepsilon > 0$ such that it is NP-hard to approximate the Euclidean k -clustering problem (6.7) to a factor better than $(1 + \varepsilon)$.

In alternative to (6.7), we can restate the problem by using assignment matrices via

$$\begin{aligned} \min \left\{ \sum_{T \in \mathcal{T}} \text{SE}_B(T) \mid \mathcal{T} \in \mathcal{P}_k^n \right\} &= \min \left\{ \sum_{T \in \mathcal{T}} \min \left\{ \sum_{i \in T} \|\mathbf{b}_i - \mathbf{x}\|_2^2 \mid \mathbf{x} \in \mathbb{R}^d \right\} \mid \mathcal{T} \in \mathcal{P}_k^n \right\} \\ &= \min \left\{ \sum_{T \in \mathcal{T}} \sum_{i \in T} \|\mathbf{b}_i - \mathbf{x}_T\|_2^2 \mid \mathcal{T} \in \mathcal{P}_k^n, \{\mathbf{x}_T\}_{T \in \mathcal{T}} \subseteq \mathbb{R}^d \right\} \\ &= \min \left\{ \sum_{j \in [k]} \sum_{i \in [n]} u_{ij} \|\mathbf{b}_i - \mathbf{x}_j\|_2^2 \mid \mathbf{U} \in \mathcal{U}_{n,k}, \mathbf{X} \in \mathbb{R}^{d \times k} \right\}. \end{aligned}$$

It is important to note that due to symmetry, the minimizers of this formulation are not unique. While we could remedy this by working with $\mathcal{U}_{n,k}^{\text{lex}}$, it would be difficult to apply MM due to the newly introduced centroid variables x_j . Instead, we will reformulate the problem based on the following lemma.

Lemma 6.27:

The function SE_B has the following explicit form:

$$\text{SE}_B(T) = \sum_{i \in T} \|\mathbf{b}_i\|_2^2 - \frac{1}{|T|} \left\| \sum_{i \in T} \mathbf{b}_i \right\|_2^2$$

In particular, the unique minimizer of $\text{SE}_B(T)$ is given by

$$\mathbf{x}_T^* = \frac{\sum_{i \in T} \mathbf{b}_i}{|T|}.$$

Proof. For a fixed $T \in 2^{[n]}$, let $g(\mathbf{x}) = \sum_{i \in T} \|\mathbf{b}_i - \mathbf{x}\|_2^2$. Since g is convex, it suffices to find a local minimizer by finding a point where the gradient vanishes. This leads to

$$\nabla g(\mathbf{x}) = 2 \sum_{i \in T} (\mathbf{x} - \mathbf{b}_i),$$

which vanishes if and only if $|T| \cdot \mathbf{x} = \sum_{i \in T} \mathbf{b}_i$. Consequently,

$$\begin{aligned} \text{SE}_{\mathbf{B}}(T) = g(\mathbf{x}_T^*) &= \sum_{i \in T} \left\| \mathbf{b}_i - \frac{\sum_{i \in T} \mathbf{b}_i}{|T|} \right\|_2^2 \\ &= \sum_{i \in T} \left(\|\mathbf{b}_i\|_2^2 - \frac{2}{|T|} \langle \mathbf{b}_i, \sum_{j \in T} \mathbf{b}_j \rangle + \left\| \frac{\sum_{i \in T} \mathbf{b}_i}{|T|} \right\|_2^2 \right) \\ &= \sum_{i \in T} \|\mathbf{b}_i\|_2^2 - \frac{1}{|T|} \left\| \sum_{i \in T} \mathbf{b}_i \right\|_2^2. \end{aligned}$$

□

Corollary 6.28:

For each $T \in 2^{[n]}$, the unique minimizer of $\text{SE}_{\mathbf{B}}(T)$ is contained in $\text{conv}(\{\mathbf{b}_i \mid i \in [n]\})$, the convex hull of the data points. Since each vertex \mathbf{b}_i is a minimizer of $\text{SE}_{\mathbf{B}}(\{i\})$, this set coincides with the convex hull of all minimizers of $\text{SE}_{\mathbf{B}}$, which shows that it has a compact description as a polytope on n vertices.

Lemma 6.27 allows us to formulate Euclidean k -clustering as a linear optimization problem over the set of combinatorial projection matrices as follows.

Theorem 6.29 ([PX05]):

The Euclidean k -clustering problem (6.7) is equivalent to

$$\min \{ \langle \mathbf{B}^T \mathbf{B}, \mathbf{I}_n - \mathbf{R} \rangle \mid \mathbf{R} \in \text{CProM}_k^n \}.$$

Proof. Let $\mathcal{T} \in \mathcal{P}_k^n$ and let $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_k) \in \mathcal{U}_{n,k}$ be the corresponding assignment matrix. Applying Lemma 6.27 to (6.7) shows

$$\begin{aligned} \sum_{T \in \mathcal{T}} \text{SE}_{\mathbf{B}}(T) &= \sum_{T \in \mathcal{T}} \sum_{i \in T} \|\mathbf{b}_i\|_2^2 - \sum_{T \in \mathcal{T}} \frac{1}{|T|} \left\| \sum_{i \in T} \mathbf{b}_i \right\|_2^2 = \sum_{i \in [n]} \|\mathbf{b}_i\|_2^2 - \sum_{j \in [k]} \frac{1}{\|\mathbf{u}_j\|_2^2} \|\mathbf{B} \mathbf{u}_j\|_2^2 \\ &= \langle \mathbf{B}^T \mathbf{B}, \mathbf{I}_n \rangle - \sum_{j \in [k]} \frac{1}{\|\mathbf{u}_j\|_2^2} \langle \mathbf{B}^T \mathbf{B}, \mathbf{u}_j \mathbf{u}_j^T \rangle = \left\langle \mathbf{B}^T \mathbf{B}, \mathbf{I}_n - \sum_{j \in [k]} \frac{\mathbf{u}_j \mathbf{u}_j^T}{\|\mathbf{u}_j\|_2^2} \right\rangle \\ &= \langle \mathbf{B}^T \mathbf{B}, \mathbf{I}_n - \rho(\mathbf{U}) \rangle. \end{aligned}$$

Since $\rho(\mathcal{U}_{n,k}) = \text{CProM}_k^n$ by Theorem 6.7, the result follows. □

This formulation was already proposed in [PX05], where it was shown that the same is true for a slightly relaxed formulation of CProM_k^n , where the block inducing equations are replaced with the more general projection constraint $\mathbf{R}^2 = \mathbf{R}$.

Convexification

In [PW07], it was shown that $\text{RProM}_k^{n,0}$ is a suitable relaxation for CProM_k^n and that using techniques similar to spectral clustering, one can use the solution of

$$\min \{ \langle \mathbf{B}^\top \mathbf{B}, \mathbf{I}_n - \mathbf{R} \rangle \mid \mathbf{R} \in \text{RProM}_k^{n,0} \}. \quad (6.8)$$

to arrive at a Euclidean k -clustering problem in dimension \mathbb{R}^{k-1} . Solving the new problem is, for the common situation where $d \gg k$, significantly easier, and it is shown that such a solution can be used to reconstruct a 2-approximation of the original problem. In particular, for the special case of $k = 2$, the resulting clustering problem can be solved in time $\mathcal{O}(n \log n)$.

While it is possible to improve the results of (6.8) by going from $\text{RProM}_k^{n,0}$ to the sharper relaxation RProM_k^n , the runtimes of other approximation algorithms give little incentives to do so. For a proper overview over viable alternatives, we direct the reader to the paper [PW07].

Chapter 7

Affine Euclidean Clustering

This chapter covers the treatment of what we call the affine Euclidean k -clustering problem by using a simplicial cover of the feasible set. This is the largest chapter of this thesis, since this problem is very general and can express many different kinds of partition problems.

In Section 7.1, we introduce the problem and show how it is an overall much harder generalization of Euclidean k -clustering. In order to tackle this problem anyway, we introduce the concept of simplicial cover in Section 7.2 and show how to use this to reformulate the problem into a much more structured form, ready to be exploited in Section 7.3 to formulate convex relaxations. At this point, we are able to relate our relaxation to existing work in Section 7.4 before we go on by showing how to extract feasible solutions in Section 7.5. After that, we show how our approach can be generalized to various variants of the problem in Section 7.6 before we finish the chapter with a discussion of applications in Section 7.7.

7.1 Overview

This section outlines the problem of affine Euclidean k -clustering and the associated difficulties in solving it.

7.1.1 Problem Formulation

In Section 6.3.2, we considered set functions SE_B , which mapped a subset of given points to their least-square centroid as partition criterion. In this chapter, we will examine a generalization of these functions with increased expressiveness and start by introducing them.

Definition 7.1:

A set function $SE_B^{\mathcal{A}} : 2^{[n]} \rightarrow \mathbb{R}$ is called an *affine squared error function*, if there are

matrices $\mathcal{A} := \{\mathbf{A}_i\}_{i \in [n]} \subseteq \mathbb{R}^{l \times d}$ and $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n) \in \mathbb{R}^{l \times n}$ such that

$$\text{SE}_{\mathbf{B}}^{\mathcal{A}}(T) = \min \left\{ \sum_{i \in T} \|\mathbf{A}_i \mathbf{x} - \mathbf{b}_i\|_2^2 \mid \mathbf{x} \in \mathbb{R}^d \right\}. \quad (7.1)$$

Setting $d = l$ and using $\mathbf{A}_i = \mathbf{I}_d$ for all $i \in [n]$, we recover $\text{SE}_{\mathbf{B}}$. However, compared to $\text{SE}_{\mathbf{B}}$, the map $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$ is more complex, as it can in general not be geometrically understood as computing a centroid of the point set $\{\mathbf{b}_i\}_{i \in T}$. Following Section 6.3.2, we define the *affine Euclidean k -clustering problem* as the separable partition problem given by

$$\min \left\{ \sum_{T \in \mathcal{T}} \text{SE}_{\mathbf{B}}^{\mathcal{A}}(T) \mid \mathcal{T} \in \mathcal{P}_k^n \right\}. \quad (7.2)$$

Alternatively, we can restate the problem by using assignment matrices via

$$\begin{aligned} \min \left\{ \sum_{T \in \mathcal{T}} \text{SE}_{\mathbf{B}}^{\mathcal{A}}(T) \mid \mathcal{T} \in \mathcal{P}_k^n \right\} &= \min \left\{ \sum_{T \in \mathcal{T}} \min \left\{ \sum_{i \in T} \|\mathbf{A}_i \mathbf{x} - \mathbf{b}_i\|_2^2 \mid \mathbf{x} \in \mathbb{R}^d \right\} \mid \mathcal{T} \in \mathcal{P}_k^n \right\} \\ &= \min \left\{ \sum_{T \in \mathcal{T}} \sum_{i \in T} \|\mathbf{A}_i \mathbf{x}_T - \mathbf{b}_i\|_2^2 \mid \mathcal{T} \in \mathcal{P}_k^n, \{\mathbf{x}_T\}_{T \in \mathcal{T}} \subseteq \mathbb{R}^d \right\} \\ &= \min \left\{ \sum_{j \in [k]} \sum_{i \in [n]} u_{ij} \|\mathbf{A}_i \mathbf{x}_j - \mathbf{b}_i\|_2^2 \mid \mathbf{U} \in \mathcal{U}_{n,k}, \mathbf{X} \in \mathbb{R}^{d \times k} \right\}. \end{aligned}$$

For later reference, the formulation

$$\min \left\{ \sum_{j \in [k]} \sum_{i \in [n]} u_{ij} \|\mathbf{A}_i \mathbf{x}_j - \mathbf{b}_i\|_2^2 \mid \mathbf{U} \in \mathcal{U}_{n,k}, \mathbf{X} \in \mathbb{R}^{d \times k} \right\}. \quad (7.3)$$

will be referred to as *affine Euclidean k -clustering* from now on and will be the focus of this chapter. We start with a basic observation.

Lemma 7.2:

Problem (7.3) is a polynomial optimization problem.

Proof. Since

$$\|\mathbf{A}_i \mathbf{x}_j - \mathbf{b}_i\|_2^2 = \mathbf{x}_j^\top (\mathbf{A}_i^\top \mathbf{A}_i) \mathbf{x}_j - 2(\mathbf{b}_i^\top \mathbf{A}_i) \mathbf{x}_j + \|\mathbf{b}_i\|_2^2$$

and

$$u_{ij} \in \{0, 1\} \quad \Leftrightarrow \quad u_{ij}(1 - u_{ij}) = 0$$

for all $i \in [n]$ and $j \in [k]$, we need to optimize a polynomial over a real variety. \square

We point out that again due to symmetry, the minimizers of (7.3) are not unique, since any permutation acting on the columns of both \mathbf{U} and \mathbf{X} simultaneously will leave the objective value invariant.

While we could hope on using $\mathcal{U}_{n,k}^{lex}$ to remedy this, it would be difficult to apply MM due to the newly introduced centroid variables \mathbf{X} . In particular, the total degree of the objective function is 3 due to \mathbf{X} , so that we would need the prohibitive big second stage of MM, which is not available for practical computations. Furthermore, it is hard to extract feasible solutions from lower levels of MM, and the symmetrical structure of the problem makes this even harder.

Instead, a better strategy would be to first reduce the problem to a more compact reformulation. To this end, it makes sense to follow Section 6.3.2 again and to try finding a closed form for $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$. This leads to the following lemma.

Lemma 7.3:

The minimizers of $\text{SE}_{\mathbf{B}}^{\mathcal{A}}(T)$ are given as the solution space of the linear system

$$\sum_{i \in T} \mathbf{A}_i^{\top} \mathbf{b}_i = \left(\sum_{i \in T} \mathbf{A}_i^{\top} \mathbf{A}_i \right) \mathbf{x}.$$

Proof. For a fixed set $T \in 2^{[n]}$, let $g(\mathbf{x}) = \sum_{i \in T} \|\mathbf{A}_i \mathbf{x} - \mathbf{b}_i\|_2^2$. Since g is continuous and convex, it suffices to find a local minimizer by finding a point where the gradient vanishes. This leads to

$$\nabla g(\mathbf{x}) = 2 \sum_{i \in T} \mathbf{A}_i^{\top} (\mathbf{A}_i \mathbf{x} - \mathbf{b}_i),$$

which vanishes if and only if \mathbf{x} solves the linear system $\sum_{i \in T} \mathbf{A}_i^{\top} \mathbf{b}_i = \left(\sum_{i \in T} \mathbf{A}_i^{\top} \mathbf{A}_i \right) \mathbf{x}$. \square

Unfortunately, the introduction of the linear maps \mathbf{A}_i makes it very difficult to give an explicit formula for the minimizer of $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$, so that we can not follow the approach from Section 6.3.2. Before pursuing a different line of thought in the following section, we finish with a theorem that highlights the potential complexity of $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$.

Theorem 7.4:

We can find matrices $\{\mathbf{A}_i\}_{i \in [n]}$ and vectors $\{\mathbf{b}_i\}_{i \in [n]}$ such that the convex hull of the minimizers of $\text{SE}_{\mathbf{B}}^{\mathcal{A}}(T)$ is a polytope with an exponential number of vertices and facets in the input size.

Proof. We will explicitly construct such an instance for a given number n of pairs (\mathbf{A}, \mathbf{b}) . Before we start, first note that if we partition

$$\mathbf{A}_i = \begin{pmatrix} \mathbf{A}_{i,1} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{i,2} \end{pmatrix}, \quad \mathbf{b}_i = \begin{pmatrix} \mathbf{b}_{i,1} \\ \mathbf{b}_{i,2} \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}$$

into blocks of consistent dimension, we can separate

$$\text{SE}_{\mathbf{B}}^{\mathcal{A}}(T) = \text{SE}_{\mathbf{B}_1}^{\mathcal{A}_1}(T) + \text{SE}_{\mathbf{B}_2}^{\mathcal{A}_2}(T)$$

where $\mathcal{A}_j = \{\mathbf{A}_{i,j}\}_{i \in [n]} \subseteq \mathbb{R}^{l_j \times d_j}$ and $\mathbf{B}_j = (\mathbf{b}_{1,j}, \dots, \mathbf{b}_{n,j}) \in \mathbb{R}^{l_j \times n}$ for $j = 1, 2$ to compute \mathbf{x}_1 and \mathbf{x}_2 separately. Furthermore, if we can partition $[n] = T_1 \cup T_2$ such that

$$\begin{aligned} i \in T_1 &\Rightarrow (\mathbf{A}_{i,2}, \mathbf{b}_{i,2}) = (\mathbf{0}, \mathbf{0}), \\ i \in T_2 &\Rightarrow (\mathbf{A}_{i,1}, \mathbf{b}_{i,1}) = (\mathbf{0}, \mathbf{0}), \end{aligned}$$

then

$$\arg \min \left(\text{SE}_{\mathbf{B}}^{\mathcal{A}}(T) \right) = \arg \min \left(\text{SE}_{\mathbf{B}_1}^{\mathcal{A}_1}(T \cap T_1) \right) \times \arg \min \left(\text{SE}_{\mathbf{B}_2}^{\mathcal{A}_2}(T \cap T_2) \right).$$

It thus suffices to construct one instance that results in an exponential number of vertices and one that results in an exponential number of facets to show the result, since vertices and facets are preserved under taking the convex hull of their Cartesian product.

For the first instance, we create the n -dimensional unit-cube C^n by setting $d = n$, $l = 1$ and

$$(\mathbf{A}_i, \mathbf{b}_i) = (\mathbf{e}_i^\top, 1) \in \mathbb{R}^{1 \times n} \times \mathbb{R} \quad \forall i \in [n].$$

Then for any set $T \in 2^{[n]}$, we can set $\mathbf{x}_T = \mathbf{e}_T$ to see that since

$$\sum_{i \in T} \mathbf{A}_i^\top \mathbf{b}_i = \sum_{i \in T} \mathbf{e}_i \cdot 1 = \sum_{i \in T} \mathbf{e}_i \langle \mathbf{e}_i, \mathbf{e}_T \rangle = \left(\sum_{i \in T} \mathbf{e}_i \mathbf{e}_i^\top \right) \mathbf{e}_T = \left(\sum_{i \in T} \mathbf{A}_i^\top \mathbf{A}_i \right) \mathbf{x}_T,$$

the solution \mathbf{x}_T is optimal for $\text{SE}_{\mathbf{B}}^{\mathcal{A}}(T)$ according to Lemma 7.3. This takes care of the exponential number of vertices in n , since we get a unique solution for each $T \in 2^{[n]}$, resulting in 2^n vertices.

For the second instance, we construct the n -dimensional cross polytope given by

$$(C^n)^* = \text{conv}(\{\pm \mathbf{e}_i \mid i \in [n]\})$$

by setting $d = l = n$ and

$$(\mathbf{A}_{i,\pm}, \mathbf{b}_{i,\pm}) = (\mathbf{I}_n, \pm \mathbf{e}_i) \in \mathbb{R}^{n \times n} \times \mathbb{R}^n \quad \forall i \in [n].$$

This way, we model the special case of $\text{SE}_{\mathbf{B}}$, where

$$\mathbf{B} = (+\mathbf{e}_1, -\mathbf{e}_1, \dots, +\mathbf{e}_n, -\mathbf{e}_n) \in \mathbb{R}^{n \times 2n}.$$

By Corollary 6.28, it follows immediately that the convex hull of the minimizers of $\text{SE}_{\mathbf{B}}$ is equal to $(C^n)^*$. This takes care of the exponential number of facets in n , since $(C^n)^*$ is the dual-polytope of C^n , and thus has one facet for each of the 2^n vertices of C^n .

Combining both instances through their Cartesian product results in an instance with $3n$ pairs (\mathbf{A}, \mathbf{b}) and at least 2^n vertices and facets, thus completing the proof. \square

7.2 Simplicial Covers

In its general form, it seems very hard to tackle problem (7.3). Instead, we will introduce some assumptions that make the problem easier to handle, starting with a bound on the minimizers of $\text{SE}_B^{\mathcal{A}}$. We will need a definition first.

Definition 7.5 (Simplicial cover):

A finite set of d -dimensional simplices $\{P_s\}_{s \in [q]}$ with disjoint interior is called a *simplicial cover* of $\text{SE}_B^{\mathcal{A}}$ if the inclusion

$$\{\mathbf{x}_T\}_{T \in 2^{[n]}} \subseteq \mathcal{P} = \bigcup_{s \in [q]} P_s$$

holds for the minimizers \mathbf{x}_T of $\text{SE}_B^{\mathcal{A}}(T)$.

In the following, we will abuse notation by using \mathcal{P} to denote both the set of simplices $\{P_s \mid s \in [q]\}$ as well as their union.

Remark 7.6:

The concept of simplicial cover is motivated by what is known in literature as a *simplicial complex*. A simplicial complex is a collection of simplices $\mathcal{P} = \bigcup_{s \in [q]} P_s$ with the property that whenever F_s is a face of P_s and $F_{s'}$ is a face of $P_{s'}$, their intersection $F_s \cap F_{s'}$ is a face of both P_s and $P_{s'}$, independent of the choice of $s, s' \in [q]$. In particular, such a simplicial complex defines an independence system on the set of vertices of a simplicial cover \mathcal{P} , where the independent sets are precisely the sets of vertices that lie on a common face of any simplex P_s . By our definition, each simplicial complex gives rise to a simplicial cover, but the reverse is not true in general, as shown by Figure 7.1.

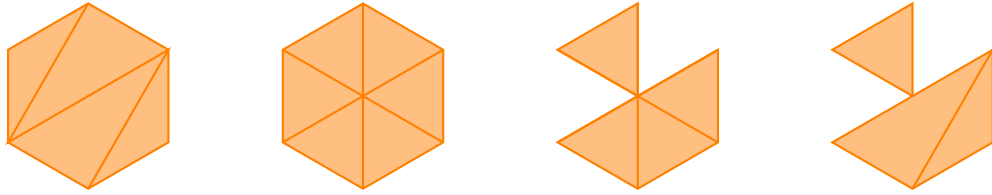


Figure 7.1: Four examples of simplicial covers. All simplicial covers but the last one are also simplicial complexes.

A natural question is to ask about the availability of such a simplicial cover. Since our goal is to use this concept to solve the optimization problem (7.3), the construction of such an simplicial cover should be much cheaper than computing each of the exponentially many minimizers $\{\mathbf{x}_T\}_{T \in 2^{[n]}}$ first. To this end, we will use the following underlying assumption for now.

Assumption 7.7 (Simplicial cover assumption):

Given $\{(A_i, \mathbf{b}_i)\}_{i \in [n]}$, we can construct a simplicial cover of $\text{SE}_B^{\mathcal{A}}$ in polynomial time.

The crucial property of Assumption 7.7 is not the existence but rather the efficiency of the underlying construction. In particular, a simplicial cover always exists, since each point \mathbf{x}_T can be computed by Lemma 7.3, and the convex hull of these points is a polytope that can be turned into a simplicial complex via various methods [BEF00].

Alternatively, if a bound on the norm of all minimizers is available, then a single simplex containing all minimizers can easily be constructed. Such a bound is not unlikely to be available in applications [XWI05].

Example 7.8:

Setting each matrix \mathbf{A}_i to the identity matrix, we recover the original Euclidean clustering problem (6.7). Then Assumption 7.7 is satisfied by Corollary 6.28, since

$$\max \{ \|\mathbf{b}_i\|_2 \mid i \in [n] \}$$

is a bound for the norm of all points in the convex hull of the minimizers.

We will skip the realization of Assumption 7.7 for now and postpone a closer investigation to Section 7.7 where we give examples that satisfy the assumption naturally. For now, we proceed by introducing notation based on a fixed simplicial cover \mathcal{P} for $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$. Let $\mathbf{V}_s \in \mathbb{R}^{d \times (d+1)}$ be the matrix whose columns denote the vertices of P_s . Then $\text{conv}(\mathbf{V}_s) = P_s$ and we gather all vertices in $\mathbf{V} = (\mathbf{V}_1, \dots, \mathbf{V}_q) \in \mathbb{R}^{d \times m}$, where

$$m := \sum_{s \in [q]} |\mathbf{V}_s| = q(d+1).$$

Furthermore, we implicitly define index sets $v(s)$ for each P_s by partitioning

$$[m] = \bigcup_{s \in [q]} v(s).$$

Then membership $\mathbf{x} \in \mathcal{P}$ naturally implies the existence of

$$\boldsymbol{\lambda}^\top := (\boldsymbol{\lambda}_{v(1)}^\top, \dots, \boldsymbol{\lambda}_{v(q)}^\top) \in \Delta^m$$

where

$$\mathbf{x} = \mathbf{V}\boldsymbol{\lambda} = \sum_{s \in [q]} \mathbf{V}_s \boldsymbol{\lambda}_{v(s)}.$$

Remark 7.9:

If the simplices P_s share common vertices, then \mathbf{V} has multiple identical columns across different \mathbf{V}_s . This is intentional, and the potential benefit of removing multiples of these columns will be explored in Section 7.6.

Now \mathcal{P} can be expressed as the linear image of a constrained standard simplex through

$$\mathcal{P} = \left\{ \mathbf{x} \in \mathbb{R}^d \mid \exists \boldsymbol{\lambda} \in \Delta^m : \mathbf{x} = \mathbf{V}\boldsymbol{\lambda}, \quad \boldsymbol{\lambda}_{v(r)} \boldsymbol{\lambda}_{v(s)}^\top = \mathbf{0} \quad \forall r, s \in [q], r \neq s \right\}. \quad (7.4)$$

The nonlinear orthogonality constraint

$$\boldsymbol{\lambda}_{v(r)} \boldsymbol{\lambda}_{v(s)}^\top = \mathbf{0} \quad \forall r, s \in [q], r \neq s \quad (7.5)$$

ensures that exactly one $\boldsymbol{\lambda}_{v(s)}$ is nonzero, which implies $\mathbf{x} = \mathbf{V}\boldsymbol{\lambda} = \mathbf{V}_s \boldsymbol{\lambda}_{v(s)} \in P_s$ for some $s \in [q]$. Since $\boldsymbol{\lambda} \geq \mathbf{0}$, we can see that (7.5) is equivalent to the quadratic constraint

$$\boldsymbol{\lambda}^\top \boldsymbol{\Omega} \boldsymbol{\lambda} = 0,$$

where

$$\boldsymbol{\Omega} := (\mathbf{J}_q - \mathbf{I}_q) \otimes \mathbf{J}_{d+1} \in \{0, 1\}^{m \times m}$$

induces a block structure with zero block-diagonal. This suggests to set

$$\Delta_\Omega^m := \{\boldsymbol{\lambda} \in \Delta^m \mid \boldsymbol{\lambda}^\top \boldsymbol{\Omega} \boldsymbol{\lambda} = 0\}$$

so that in particular, we can write

$$\mathcal{P} = \mathbf{V} \cdot \Delta_\Omega^m$$

as a shorthand for (7.4).

Remark 7.10:

The index sets $\{v(s) \mid s \in [q]\} \subseteq 2^{[m]}$ induce a relatively simple independence set

$$L = \bigcup_{s \in [q]} 2^{v(s)} \subseteq 2^{[m]}$$

due to their disjointness, and our choice of $\boldsymbol{\Omega}$ coincides with the $\boldsymbol{\Omega}_L$ matrix defined for this independence system L , which makes this problem another instance of problems constrained by independence systems.

Remark 7.11:

The parametrization (7.4) of points in \mathcal{P} is unique almost everywhere. This is implicit in Theorem 2.8, which guarantees uniqueness for all interior points $\mathbf{x} \in \bigcup_{s \in [q]} \text{int}(P_s)$. However, since simplicial covers allow intersections of boundaries $\text{bd}(P_r) \cap \text{bd}(P_s)$ to be nonempty, we can not assume unique representations for boundary points once $q > 1$.

Given the new notation, we can proceed by parametrizing (7.1) in terms of (7.4). Using $\mathbf{x} = \mathbf{V}\boldsymbol{\lambda}$, we can homogenize the objective by using $1 = \langle \boldsymbol{\lambda}, \mathbf{e} \rangle$ to end up with

$$\|\mathbf{A}_i \mathbf{x} - \mathbf{b}_i\|_2^2 = \langle \boldsymbol{\lambda}^j, \mathbf{C}_i \boldsymbol{\lambda}^j \rangle,$$

where

$$\mathbf{C}_i := \mathbf{V}^\top \mathbf{A}_i^\top \mathbf{A}_i \mathbf{V} - (\mathbf{e} \mathbf{b}_i^\top \mathbf{A}_i \mathbf{V} + \mathbf{V}^\top \mathbf{A}_i^\top \mathbf{b}_i \mathbf{e}^\top) + \|\mathbf{b}_i\|_2^2 \cdot \mathbf{J}_m \quad \forall i \in [n].$$

This leads to the reformulations

$$\text{SE}_{\mathbf{B}}^{\mathcal{A}}(T) = \min \left\{ \sum_{i \in T} \langle \boldsymbol{\lambda}, \mathbf{C}_i \boldsymbol{\lambda} \rangle \mid \boldsymbol{\lambda} \in \Delta_{\Omega}^m \right\}$$

of (7.1) and

$$\min \left\{ \sum_{i \in [n]} \sum_{j \in [k]} u_{ij} \langle \boldsymbol{\lambda}^j, \mathbf{C}_i \boldsymbol{\lambda}^j \rangle \mid \mathbf{U} \in \mathcal{U}_{n,k}, \boldsymbol{\lambda}^j \in \Delta_{\Omega}^m \quad \forall j \in [k] \right\} \quad (\text{R1})$$

of (7.3).

7.2.1 Separating Simplicial Covers

So far, the goal of every chapter has been to reformulate the partition problems by replacing the assignment matrices with a class of matrices that is better suited to deal with the underlying structure of the problem at hand.

We will continue to follow this paradigm and start by noting that in the context of (R1), the assignment matrix \mathbf{U} can be interpreted as assigning an optimal minimizer of $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$ to each data point. To see this, consider the equation

$$\sum_{j \in [k]} u_{ij} \langle \boldsymbol{\lambda}^j, \mathbf{C}_i \boldsymbol{\lambda}^j \rangle = \left\langle \left(\sum_{j \in [k]} u_{ij} \boldsymbol{\lambda}^j \right), \mathbf{C}_i \left(\sum_{j \in [k]} u_{ij} \boldsymbol{\lambda}^j \right) \right\rangle,$$

which holds since each row of \mathbf{U} is a binary unit vector. Letting $\boldsymbol{\lambda}_i(\mathbf{U}) := \sum_{j \in [k]} u_{ij} \boldsymbol{\lambda}^j$ for each $i \in [n]$ then shows that

$$\mathbf{U} \in \mathcal{U}_{n,k} \Leftrightarrow \{\boldsymbol{\lambda}_i(\mathbf{U}) \mid i \in [n]\} = \{\boldsymbol{\lambda}^j \mid j \in [k]\}.$$

In particular, (R1) can be restated as

$$\min \left\{ \sum_{i \in [n]} \langle \boldsymbol{\lambda}_i, \mathbf{C}_i \boldsymbol{\lambda}_i \rangle \mid \boldsymbol{\lambda}_i \in \{\boldsymbol{\lambda}^j\}_{j \in [k]} \subseteq \Delta_{\Omega}^m \quad \forall i \in [n] \right\}, \quad (7.6)$$

where $\boldsymbol{\lambda}_i(\mathbf{U})$ is replaced by a new vector $\boldsymbol{\lambda}_i$ subject to the membership constraint $\boldsymbol{\lambda}_i \in \{\boldsymbol{\lambda}^j\}_{j \in [k]}$.

At this point, we want to recall again that while it is important to model this membership constraint in a more tractable formulation, using the set $\mathcal{U}_{n,k}$ introduced the inherent problematic symmetries mentioned in Section 4.2.1 by inducing an arbitrary order on the set $\{\boldsymbol{\lambda}^j\}_{j \in [k]}$.

For this reason, our goal is to propose a different, symmetry-free way to formalize this membership, and we will use the following property as a starting point.

Definition 7.12 (Separating simplicial cover):

Let $\mathcal{P} = \bigcup_{s \in [q]} P_s$ be a simplicial cover of $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$. Then a set of points $X \subseteq \mathcal{P}$ is called *separated* by \mathcal{P} if

$$|X \cap P_s| \leq 1 \quad \forall s \in [q].$$

Likewise, a set of parameters $\Lambda \subseteq \Delta_{\Omega}^m$ is *separated* by \mathcal{P} if their image under \mathbf{V} is contained in \mathcal{P} and separated by \mathcal{P} .

In addition, let \mathbf{x}_T be the minimizer of $\text{SE}_{\mathbf{B}}^{\mathcal{A}}(T)$ for all $T \in 2^{[n]}$. Then we say \mathcal{P} *separates* $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$ if there is an optimal solution \mathcal{T} of (7.2) such that $\{\mathbf{x}_T \mid T \in \mathcal{T}\}$ is separated by \mathcal{P} .

The idea behind separating simplicial covers is to guarantee pairwise orthogonality of the vectors $\boldsymbol{\lambda}^j$ in (7.6), as can be seen in the following lemma.

Lemma 7.13:

Let $\mathcal{P} = \bigcup_{s \in [q]} P_s$ be a simplicial cover of $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$. For a set $X \subseteq \mathcal{P}$, let $\Lambda \subseteq \Delta_{\Omega}^m$ be its representation in (7.4). Then X is separated by \mathcal{P} if and only if for each pair $\boldsymbol{\lambda} \neq \boldsymbol{\lambda}' \in \Lambda$,

$$\boldsymbol{\lambda}'_{v(s)} \boldsymbol{\lambda}_{v(s)}^{\top} = \mathbf{0} \quad \forall s \in [q]. \quad (7.7)$$

In particular, each pair in $\boldsymbol{\lambda} \neq \boldsymbol{\lambda}' \in \Lambda$ is coordinatewise orthogonal and has disjoint support.

Proof. Let $\mathbf{x} \neq \mathbf{x}' \in \mathcal{P}$ with $\mathbf{x} = \mathbf{V}\boldsymbol{\lambda}$ and $\mathbf{x}' = \mathbf{V}\boldsymbol{\lambda}'$. Then for some fixed $t, t' \in [q]$,

$$\mathbf{x} \in P_t \Leftrightarrow \langle \boldsymbol{\lambda}_{v(s)}, \mathbf{e} \rangle = \delta_{st} \quad \text{and} \quad \mathbf{x}' \in P_{t'} \Leftrightarrow \langle \boldsymbol{\lambda}'_{v(s)}, \mathbf{e} \rangle = \delta_{st'}$$

for all $s \in [q]$. Now due to non-negativity, each equation of (7.7) can be equivalently expressed as

$$\begin{aligned} \boldsymbol{\lambda}'_{v(s)} \boldsymbol{\lambda}_{v(s)}^{\top} = \mathbf{0} &\Leftrightarrow 0 = \langle \boldsymbol{\lambda}'_{v(s)} \boldsymbol{\lambda}_{v(s)}^{\top}, \mathbf{J}_{d+1} \rangle = \langle \boldsymbol{\lambda}'_{v(s)}, \mathbf{e} \rangle \langle \boldsymbol{\lambda}_{v(s)}, \mathbf{e} \rangle = \delta_{st'} \delta_{st} \\ &\Leftrightarrow (s \neq t) \vee (s \neq t'). \end{aligned}$$

Then letting s run through $[q]$ shows that (7.7) is equivalent to $t \neq t'$, which means that \mathbf{x} and \mathbf{x}' are separated by \mathcal{P} . \square

A direct consequence of Lemma 7.13 is that if \mathcal{P} separates $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$, then (7.7) necessarily holds for an optimal solution of (R1). Furthermore, its simple observation implies that we can express the membership constraint in (7.6) with linear inequalities and quadratic constraints.

Lemma 7.14:

Let \mathcal{P} be a simplicial cover for $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$ and let

$$\boldsymbol{\lambda}_* := \sum_{j \in [k]} \boldsymbol{\lambda}^j \quad (7.8)$$

for a separated set $\{\boldsymbol{\lambda}^j\}_{j \in [k]}$ feasible for (R1). Then the equivalence

$$\boldsymbol{\lambda} \in \{\boldsymbol{\lambda}^j\}_{j \in [k]} \quad \Leftrightarrow \quad \boldsymbol{\lambda} \leq \boldsymbol{\lambda}_* \quad (7.9)$$

holds for $\boldsymbol{\lambda} \in \Delta_{\Omega}^m$.

Proof. Implication " \Rightarrow " is straightforward. Consider the reverse direction for $\boldsymbol{\lambda} \in \Delta_{\Omega}^m$ and let $s_j \in [q]$ for $j \in [k]$ such that $\text{supp}(\boldsymbol{\lambda}^j) \subseteq v(s_j)$. Then it follows from $\boldsymbol{\lambda}^\top \Omega \boldsymbol{\lambda} = 0$, $\mathbf{0} \leq \boldsymbol{\lambda} \leq \boldsymbol{\lambda}_*$ and Lemma 7.13 that

$$\text{supp}(\boldsymbol{\lambda}) \subseteq v(s^*) \subseteq \text{supp}(\boldsymbol{\lambda}_*) = \bigcup_{j \in [k]} \text{supp}(\boldsymbol{\lambda}^j) = \bigcup_{j \in [k]} v(s_j),$$

which implies $s^* = s_{j'}$ for some $j' \in [k]$. Therefore,

$$\mathbf{0} \leq \boldsymbol{\lambda} \leq \boldsymbol{\lambda}^{j'} \quad \text{and} \quad \langle \boldsymbol{\lambda}, \mathbf{e} \rangle = 1 = \langle \boldsymbol{\lambda}^{j'}, \mathbf{e} \rangle,$$

which is only possible if $\boldsymbol{\lambda} = \boldsymbol{\lambda}^{j'}$. \square

Lemma 7.14 reduces the membership in $\{\boldsymbol{\lambda}^j\}_{j \in [k]}$ to a more tractable relation involving its sum $\boldsymbol{\lambda}_*$. The idea is that if we can also give a tractable characterization of those $\boldsymbol{\lambda}_*$ that correspond to a sum of $\{\boldsymbol{\lambda}^j\}_{j \in [k]}$ satisfying the orthogonality constraint (7.7), then we can restate (R1) without the assignment matrices when $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$ is separated by \mathcal{P} . Fortunately, we can give such a tractable characterization, as is shown in the following lemma.

Lemma 7.15:

Let

$$\mathcal{L} := \{ \{\boldsymbol{\lambda}^j\}_{j \in [k]} \subseteq \Delta_{\Omega}^m \mid (7.7) \text{ holds} \}$$

and

$$\mathcal{L}' := \{ \boldsymbol{\lambda} \in k \cdot \Delta^m \mid \langle \boldsymbol{\lambda}_{v(s)}, \mathbf{e} \rangle \boldsymbol{\lambda}_{v(s)} = \boldsymbol{\lambda}_{v(s)} \quad \forall s \in [q] \}. \quad (7.10)$$

Then $\mathcal{L} \leftrightarrow \mathcal{L}'$ through the bijection $\phi : \mathcal{L} \rightarrow \mathcal{L}'$ given by

$$\phi(\{\boldsymbol{\lambda}^j\}_{j \in [k]}) := \sum_{j \in [k]} \boldsymbol{\lambda}^j.$$

Proof. To verify that ϕ is well-defined, we recall from the proof of Lemma 7.13 that

$$\langle \boldsymbol{\lambda}_{v(s)}^j, \mathbf{e} \rangle = \delta_{s, s_j}$$

for some $s_j \in [q]$, and that the s_j are pairwise distinct. In particular, the identity

$$\delta_{s, s_j} \boldsymbol{\lambda}_{v(s)}^j = \boldsymbol{\lambda}_{v(s)}^j$$

holds for all $j \in [k]$ and all $s \in [q]$. Then

$$\left\langle \sum_{j \in [k]} \lambda_{v(s)}^j, \mathbf{e} \right\rangle \left(\sum_{j \in [k]} \lambda_{v(s)}^j \right) = \left(\sum_{j \in [k]} \delta_{s,s_j} \right) \left(\sum_{j \in [k]} \delta_{s,s_j} \lambda_{v(s)}^j \right) = \sum_{j \in [k]} \delta_{s,s_j} \lambda_{v(s)}^j = \sum_{j \in [k]} \lambda_{v(s)}^j$$

since $\delta_{s,s_j} \delta_{s,s_{j'}} = \delta_{j,j'} \delta_{s,s_j}$, and so ϕ is well-defined.

We proceed by constructing an inverse-function $\psi : \mathcal{L}' \rightarrow \mathcal{L}$, so let $\lambda \in \mathcal{L}'$. By taking the scalar product with \mathbf{e} on each of the defining equations in (7.10), we see the membership $\langle \lambda_{v(s)}, \mathbf{e} \rangle \in \{0, 1\}$ for all $s \in [q]$. By definition, there is a set $\{s_j\}_{j \in [k]} \subseteq [q]$ such that $\langle \lambda_{v(s_j)}, \mathbf{e} \rangle = 1$ for all $j \in [k]$. Now define vectors $\{\lambda^j\}_{j \in [k]}$ according to

$$\lambda_{v(s)}^j := \begin{cases} \lambda_{v(s)} & \text{if } s = s_j, \\ \mathbf{0} & \text{else} \end{cases} \quad \forall s \in [q], \quad \forall j \in [k]$$

and set $\psi(\lambda) := \{\lambda^j\}_{j \in [k]}$. It is easy to check that ψ is well-defined. Now for $\lambda \in \mathcal{L}'$ one has

$$\phi(\psi(\lambda))_{v(s)} = \sum_{j \in [k]} \lambda_{v(s)}^j = \begin{cases} \lambda_{v(s)} & \text{if } s \in \{s_j\}_{j \in [k]}, \\ \mathbf{0} = \lambda_{v(s)} & \text{else,} \end{cases}$$

which shows $\phi \circ \psi = \text{Id}_{\mathcal{L}'}$.

For $\{\lambda^j\}_{j \in [k]} \in \mathcal{L}$, let $\lambda = \phi(\{\lambda^j\}_{j \in [k]})$. Then we can choose $\{s_j\}_{j \in [k]}$ such that

$$1 = \langle \lambda_{v(s_j)}, \mathbf{e} \rangle = \sum_{j' \in [k]} \langle \lambda_{v(s_j)}^{j'}, \mathbf{e} \rangle = \langle \lambda_{v(s_j)}^j, \mathbf{e} \rangle$$

by Lemma 7.13, which shows $\psi \circ \phi = \text{Id}_{\mathcal{L}}$. □

With all the tools available from the preceding results, we can restate the variant (7.6) of (R1) as the following, symmetry-free polynomial optimization problem

$$\min \left\{ \sum_{i \in [n]} \langle \lambda_i, C_i \lambda_i \rangle \left| \begin{array}{l} \lambda_* \in k \cdot \Delta^m, \\ \langle (\lambda_*)_{v(s)}, \mathbf{e} \rangle (\lambda_*)_{v(s)} = (\lambda_*)_{v(s)} \quad \forall s \in [q], \\ \lambda_i \in \Delta_{\Omega}^m \quad \forall i \in [n], \\ \lambda_i \leq \lambda_* \quad \forall i \in [n] \end{array} \right. \right\}. \quad (\text{R2})$$

Formally, we have the following connection between (R1) and (R2).

Theorem 7.16:

Problem (R2) computes the optimal separated solution of (R1). In particular, both problems are equivalent if and only if \mathcal{P} separates $\text{SE}_{\mathbf{B}}^{\mathcal{A}}$.

Proof. Starting from (7.6), the optimal separated solution of (R1) is given as

$$\min \left\{ \sum_{i \in [n]} \langle \lambda_i, C_i \lambda_i \rangle \mid \lambda_i \in \{\lambda^j\}_{j \in [k]} \quad \forall i \in [n], \{\lambda^j\}_{j \in [k]} \in \mathcal{L} \right\}$$

by Lemma 7.13. Applying Lemma 7.14, this turns into

$$\min \left\{ \sum_{i \in [n]} \langle \lambda_i, C_i \lambda_i \rangle \mid \lambda_i \leq \sum_{j \in [k]} \lambda^j \quad \forall i \in [n], \{\lambda^j\}_{j \in [k]} \in \mathcal{L} \right\},$$

which in turn is equal to (R2) given as

$$\min \left\{ \sum_{i \in [n]} \langle \lambda_i, C_i \lambda_i \rangle \mid \lambda_i \leq \lambda_* \quad \forall i \in [n], \lambda_* \in \mathcal{L}' \right\}$$

by Lemma 7.15. By definition, there is an optimal separated solution of (R1) if and only if \mathcal{P} separates $\text{SE}_B^{\mathcal{A}}$, so the rest of the theorem follows. \square

Remark 7.17:

While the separation assumption in Theorem 7.16 makes the choice of \mathcal{P} harder when we want to find the global optimum of (7.2), it can also be used as a powerful modelling paradigm. In particular, we can construct \mathcal{P} to search for the optimal solution in predetermined regions, which may be desirable in practice. For example, various facility location problems may be modelled this way.

7.3 Convexification

The goal of this section is to construct convex relaxations of (R2).

Description of (R2) as semi-algebraic set

Before we attempt to relax anything, it is necessary to fix an explicit list of polynomial inequalities to describe (R2). Given the preceding discussion, we will use the following system, which is taken directly from (R2).

$$\min_{\lambda} \sum_{i \in [n]} \langle \lambda_i, C_i \lambda_i \rangle \quad s.t. \quad (7.11a)$$

$$\langle (\lambda_*)_{v(s)}, \mathbf{e} \rangle (\lambda_*)_{v(s)} = (\lambda_*)_{v(s)} \quad \forall s \in [q] \quad (7.11b)$$

$$\langle \lambda_*, \mathbf{e} \rangle = k \quad (7.11c)$$

$$\langle \lambda_i, \mathbf{e} \rangle = 1 \quad \forall i \in [n] \quad (7.11d)$$

$$(\lambda_*)_{v(s)} (\lambda_i)_{v(s)}^\top = (\lambda_i)_{v(s)} (\lambda_i)_{v(s)}^\top \quad \forall s \in [q] \quad \forall i \in [n] \quad (7.11e)$$

$$(\lambda_i)_{v(s)} (\lambda_i)_{v(t)}^\top = \mathbf{0} \quad \forall s \neq t \in [q] \quad \forall i \in [n] \quad (7.11f)$$

$$\lambda_* \geq \lambda_i \geq \mathbf{0} \quad \forall i \in [n] \quad (7.11g)$$

It is not hard to verify that excluding (7.11e), this is a reformulation of (R2). The additional equations (7.11e) are implied by Lemma 7.13 and actually redundant. However, since their degree is low and they directly reduce the number of available monomials, we use them to decrease the size of MM in what follows.

A big advantage of this description is the following property, which will be invaluable for actual computations.

Remark 7.18:

Each polynomial of (7.11) belongs to a polynomial ring of the form $\mathbb{R}[\lambda_*, \lambda_i]$ for some $i \in [n]$. In particular, this sparsity pattern implies that once the values for λ_* are fixed, the problem decomposes into multiple independent problems. This is a crucial property to speed up computations and will be exploited algorithmically in the next section.

Let (R2)[t] denote the t -th stage of MM for the system (7.11). Since the problems are quickly expanding in size with growing t , it will only make sense to consider $t = 1$ in the following.

Simplifying (R2)[1]

While (R2)[1] is still prohibitively large on its own, we can simplify the formulation to end up with a much more tractable formulation that can be solved by current software. To this end, we will first explicitly write down an SDP-representation of (R2)[1] where we use the notation

$$\mathbf{M}_1(\lambda) = \begin{pmatrix} 1 & \lambda_1^\top & \cdots & \lambda_n^\top & \lambda_*^\top \\ \lambda_1 & \Lambda_{11} & \cdots & \Lambda_{1n} & \Lambda_{1*} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \lambda_n & \Lambda_{n1} & \cdots & \Lambda_{nn} & \Lambda_{n*} \\ \lambda_* & \Lambda_{*1} & \cdots & \Lambda_{*n} & \Lambda_{**} \end{pmatrix} \geq \mathbf{0}$$

for the moment matrix involved. We proceed by explicitly writing the first stage of MM of (R2) given by (7.11).

$$\begin{aligned} \min_{\lambda, \Lambda} \sum_{i \in [n]} \langle C_i, \Lambda_{ii} \rangle \quad & \text{s.t.} \quad (7.12a) \\ (\Lambda_{**})_{v(s)} \mathbf{e} &= (\boldsymbol{\lambda}_*)_{v(s)} \quad \forall s \in [q] \quad (7.12b) \\ \langle \boldsymbol{\lambda}_*, \mathbf{e} \rangle &= k, \quad \Lambda_{j*} \mathbf{e} = k \boldsymbol{\lambda}_j \quad \forall j \in [n] \cup \{*\} \quad (7.12c) \\ \langle \boldsymbol{\lambda}_i, \mathbf{e} \rangle &= 1, \quad \Lambda_{ji} \mathbf{e} = \boldsymbol{\lambda}_j \quad \forall j \in [n] \cup \{*\} \quad \forall i \in [n] \quad (7.12d) \\ (\Lambda_{i*})_{v(s)} &= (\Lambda_{ii})_{v(s)} \quad \forall s \in [q] \quad \forall i \in [n] \quad (7.12e) \\ \langle \Lambda_{ii}, \boldsymbol{\Omega} \rangle &= 0 \quad \forall i \in [n] \quad (7.12f) \\ \boldsymbol{\lambda}_* \geq \boldsymbol{\lambda}_i \geq \mathbf{0}, \quad \Lambda_{**} \geq \Lambda_{*i} \geq \Lambda_{ii} \geq \mathbf{0} \quad & \forall i \in [n] \quad (7.12g) \\ \mathbf{M}_1(\boldsymbol{\lambda}) &\geq \mathbf{0} \quad (7.12h) \end{aligned}$$

This can be verified by checking that each line of (7.12) is the consequence from the corresponding line in (7.11), with the addition of (7.12h).

Unfortunately, (7.12) is not tractable in practice, since the dimension of $\mathbf{M}_1(\boldsymbol{\lambda})$ is given by $\mathcal{O}(nm)$, which is prohibitive for SDP-solvers in general. However, (7.12) is not optimized for efficient computation at this point, and we will be able to make the formulation much more compact using the following lemmas.

Lemma 7.19:

Let $\mathbf{M}_1(\boldsymbol{\lambda}|i)$ denote the submatrix of $\mathbf{M}_1(\boldsymbol{\lambda})$ in (7.12) given by

$$\mathbf{M}_1(\boldsymbol{\lambda}|i) := \begin{pmatrix} 1 & \boldsymbol{\lambda}_i^\top & \boldsymbol{\lambda}_*^\top \\ \boldsymbol{\lambda}_i & \Lambda_{ii} & \Lambda_{i*} \\ \boldsymbol{\lambda}_* & \Lambda_{*i} & \Lambda_{**} \end{pmatrix}.$$

Then we can replace $\mathbf{M}_1(\boldsymbol{\lambda}) \geq \mathbf{0}$ in (7.12) with

$$\mathbf{M}_1(\boldsymbol{\lambda}|i) \geq \mathbf{0} \quad \forall i \in [n].$$

Proof. Follows from the application of [Las15, Section 8.1] to Remark 7.18. The idea is that the underlying variables of the polynomial rings $\mathbb{R}[\boldsymbol{\lambda}_i, \boldsymbol{\lambda}_*]$ satisfy the *running intersection property*, which allows to separate them for computations. \square

While this is already a significant reduction of the SDP constraint, the following lemma gives an alternative proof for the fact that we can discard the constant and linear monomials due to Theorem 2.15 or Corollary 3.3.

Lemma 7.20:

Consider a matrix $\Lambda \geq \mathbf{0}$ and a vector \mathbf{a} . Let $\mathbf{a}^\top \Lambda \mathbf{a} = v$ and define $\boldsymbol{\lambda} := \Lambda \mathbf{a}$. Then

$$v\Lambda \geq \boldsymbol{\lambda} \boldsymbol{\lambda}^\top \quad \text{or equivalently} \quad \begin{pmatrix} v & \boldsymbol{\lambda}^\top \\ \boldsymbol{\lambda} & \Lambda \end{pmatrix} \geq \mathbf{0}. \quad (7.13)$$

Proof. Since $\Lambda \geq \mathbf{0}$, there is \mathbf{M} such that $\Lambda = \mathbf{M}^\top \mathbf{M}$ and consequently,

$$v = \mathbf{a}^\top \Lambda \mathbf{a} = \|\mathbf{M}\mathbf{a}\|_2^2.$$

Then for arbitrary \mathbf{x} we have

$$\begin{aligned} \mathbf{x}^\top (v\Lambda - (\Lambda\mathbf{a})(\Lambda\mathbf{a})^\top) \mathbf{x} &= v \cdot \langle \mathbf{M}\mathbf{x}, \mathbf{M}\mathbf{x} \rangle - \langle \mathbf{M}\mathbf{x}, \mathbf{M}\mathbf{a} \rangle^2 \\ &= \|\mathbf{M}\mathbf{a}\|_2^2 \cdot \|\mathbf{M}\mathbf{x}\|_2^2 - |\langle \mathbf{M}\mathbf{x}, \mathbf{M}\mathbf{a} \rangle|^2 \geq 0 \end{aligned}$$

where the last inequality is the Cauchy-Schwarz inequality. The equivalent second formulation of (7.13) follows from the Schur complement Lemma 2.5. \square

We are now ready to simplify (7.12) in the following form.

$$\min_{\lambda, \Lambda} \sum_{i \in [n]} \langle \mathbf{C}_i, \Lambda_{ii} \rangle \quad s.t. \quad (7.14a)$$

$$(\Lambda_{**})_{v(s)} \mathbf{e} = (\Lambda_{*i})_{v(s)} \mathbf{e} \quad \forall s \in [q] \quad \forall i \in [n] \quad (7.14b)$$

$$\Lambda_{j*} \mathbf{e} = k \Lambda_{ji} \mathbf{e} \quad \forall j \in [n] \cup \{*\} \quad \forall i \in [n] \quad (7.14c)$$

$$(\Lambda_{i*})_{v(s)} = (\Lambda_{ii})_{v(s)} \quad \forall s \in [q] \quad \forall i \in [n] \quad (7.14d)$$

$$\langle \Lambda_{ii}, \Omega \rangle = 0 \quad \forall i \in [n] \quad (7.14e)$$

$$\langle \Lambda_{ii}, \mathbf{J} \rangle = 1 \quad \forall i \in [n] \quad (7.14f)$$

$$\Lambda_{**} \geq \Lambda_{*i} \geq \Lambda_{ii} \geq \mathbf{0} \quad \forall i \in [n] \quad (7.14g)$$

$$\begin{pmatrix} \Lambda_{ii} & \Lambda_{i*} \\ \Lambda_{*i} & \Lambda_{**} \end{pmatrix} \geq \mathbf{0} \quad \forall i \in [n] \quad (7.14h)$$

Theorem 7.21:

Problem (7.14) is equivalent to (7.12).

Proof. Problem (7.14) arises from (7.12) by first applying Lemma 7.19 to break down the moment matrix and then applying Theorem 2.15 to reformulate the problem in terms of monomials of degree 2. \square

Now reformulation (7.14) uses n psd. matrices of dimension $2m = 2q(d+1)$, which is still very limiting. However, we don't really exploit $\langle \Lambda_{ii}, \Omega \rangle = 0$, which results in a huge sparsity pattern in the block-diagonal Λ_{ii} .

Based on this observation, a good approach is to drop all variables in Λ_{i*} and Λ_{**} from (7.14) that do not belong to the blockdiagonal structure induced by Ω . Conceptually, the intention is that we only lose information of entries in Λ_{**} that have a minor, indirect impact on Λ_{ii} . On the computational side, this turns each constraint in (7.14h) into q separate psd. constraints of the form

$$\begin{pmatrix} (\Lambda_{ii})_{v(s)} & (\Lambda_{i*})_{v(s)} \\ (\Lambda_{*i})_{v(s)} & (\Lambda_{**})_{v(s)} \end{pmatrix} \geq \mathbf{0} \quad \forall s \in [q]$$

of size $2(d + 1)$. This is already much smaller, but using the fact that

$$(\Lambda_{i*})_{v(s)} = (\Lambda_{ii})_{v(s)},$$

we can additionally compress this to

$$(\Lambda_{**})_{v(s)} \geq (\Lambda_{ii})_{v(s)} \quad \forall s \in [q]$$

as a consequence of the Schur complement Lemma 2.5.

Formally, we end up with the following smaller relaxation of (7.12).

$$\min_{\lambda, \Lambda} \sum_{i \in [n]} \sum_{s \in [q]} \langle (\mathbf{C}_i)_{v(s)}, (\Lambda_{ii})_{v(s)} \rangle \quad s.t. \quad (7.15a)$$

$$\sum_{s \in [q]} \langle (\Lambda_{**})_{v(s)}, \mathbf{J} \rangle = k \quad (7.15b)$$

$$\sum_{s \in [q]} \langle (\Lambda_{ii})_{v(s)}, \mathbf{J} \rangle = 1 \quad \forall i \in [n] \quad (7.15c)$$

$$(\Lambda_{**})_{v(s)} \geq (\Lambda_{ii})_{v(s)} \geq \mathbf{0} \quad \forall s \in [q] \quad \forall i \in [n] \quad (7.15d)$$

Theorem 7.22:

Problem (7.15) is a relaxation of (7.14).

Proof. The objectives of both problems are equal, since only the blocks $(\Lambda_{ii})_{v(s)}$ can be nonzero due to $\langle \Lambda_{ii}, \Omega \rangle = 0$. With the same reasoning, (7.15c) is just a reformulation of $\langle \Lambda_{ii}, \mathbf{J} \rangle = 1$. Equation (7.15b) then follows from (7.14) as

$$k = k \langle \Lambda_{ii}, \mathbf{J} \rangle = (k \mathbf{e}^\top \Lambda_{ii}) \mathbf{e} = (\mathbf{e}^\top \Lambda_{*i}) \mathbf{e} = \sum_{s \in [q]} \mathbf{e}_{v(s)}^\top (\Lambda_{*i})_{v(s)} = \sum_{s \in [q]} \mathbf{e}^\top (\Lambda_{**})_{v(s)} \mathbf{e}.$$

□

While relaxation (7.15) may be slightly weaker than (7.14), it is much more tractable for computation. In fact, we only need to care about $2qn$ psd. constraints for matrices of size $d + 1$, which is much better than the n psd. constraints for matrices of size $2q(d + 1)$ in (7.14) from a computational point of view. In particular, since the system is only weakly coupled, parallel computing schemes can be efficiently used to solve the relaxation for problems of moderate parameters (n, q, d) .

7.4 Related Approaches

Now that we laid down our approach, we are able to compare it to similar approaches found in literature.

7.4.1 Moment Sequences

In the preceding sections, one main idea was (7.8), where we aggregated the optimal parameters λ^j into their sum λ_* , and continued to treat λ_* as a separate, new variable.

One way to interpret this is as aggregation of monomials of the first degree. Then, a possible way to extend this idea is to aggregate whole moment sequences instead and work with the resulting new object. In particular, letting $\mathbf{y}^j \in \mathcal{N}^*(\Delta_\Omega^d)$ denote the moment sequence of λ^j , we can define

$$\mathbf{y}_* = \sum_{j \in [k]} \mathbf{y}^j \in \mathcal{N}^*(\Delta_\Omega^d), \quad (\mathbf{y}_*)_0 = k,$$

to get the implication

$$\mathbf{y} \in \{\mathbf{y}^j\}_{j \in [k]} \Rightarrow \mathbf{y}, \mathbf{y}^* - \mathbf{y} \in \mathcal{N}^*(\Delta_\Omega^d), \quad y_0 = 1$$

as a weaker alternative to (7.9). Then, applying MM on the sets $\mathcal{N}^*(\Delta_\Omega^d)$ *individually* leads to the hierarchy

$$\min_{\mathbf{y}} \sum_{i \in [n]} L_{\mathbf{y}_i}(\mathbf{W}_i) \tag{7.16a}$$

$$\text{s.t.} \quad \begin{aligned} \mathbf{y}_i &\in \mathcal{N}_t^*(\Delta_\Omega^m), & (\mathbf{y}_i)_0 &= 1, & \forall i \in [n], \\ \mathbf{y}^* - \mathbf{y}_i &\in \mathcal{N}_t^*(\Delta_\Omega^m), & \mathbf{y}_0^* - (\mathbf{y}_i)_0 &= k - 1 & \forall i \in [n], \end{aligned} \tag{7.16b}$$

where it can be shown that for $t = 1$, this coincides with (7.15) after properly reformulating.

Remark 7.23:

(7.16) coincides with problem (21) given in [SRS15] in the respective setting.

7.4.2 Mixed Linear Regression

Conceptually, our approach has the following properties:

- (i) global variables are replaced with local estimates for each data point,
- (ii) each local estimate is optimized against its data term,
- (iii) constraints enforce global consistency among the estimates.

In [HJ16], a similar approach is considered, where objective and constraints are exchanged to yield the following properties:

- (i) global variables are replaced with local estimates for each data point,

- (ii) global consistency among the estimates is optimized,
- (iii) constraints enforce optimality of each local estimate against its data term.

We will now formally describe their approach, starting with a reformulation of (7.3) given by

$$\min \left\{ \sum_{i \in [n]} \left\| \mathbf{A}_i \left(\sum_{j \in [k]} u_{ij} \cdot \mathbf{x}^j \right) - \mathbf{b}_i \right\|_2^2 \mid \mathbf{U} \in \mathcal{U}_{n,k}, \mathbf{X} \in \mathbb{R}^{d \times k} \right\}.$$

Assuming that the linear equation systems $\mathbf{A}_i \mathbf{x} = \mathbf{b}_i$ are solvable, we can set

$$\mathbf{x}_i(\mathbf{U}) := \sum_{j \in [k]} u_{ij} \cdot \mathbf{x}^j$$

and, following our approach in Subsection 7.2.1, relax it to a local estimate \mathbf{x}_i to yield

$$\min \left\{ \sum_{i,j \in [n]} \|\mathbf{x}_i - \mathbf{x}_j\|_2 \mid \mathbf{A}_i \mathbf{x}_i = \mathbf{b}_i \ \forall i \in [n] \right\},$$

where the objective describes global coherence of the estimates. As with our approach, solving this relaxations amounts to solving a conic program, where the resulting estimates need to be clustered afterwards. In particular, this optimization problem can be formulated using the Lorentz cone \mathcal{L}_m .

In [HJ16], only matrices $\mathbf{A}_i = \mathbf{a}_i^\top$ with one column are treated, which is known in the literature as *mixed linear regression*. In this setting, the assumption of the solvability of the underlying linear equations $\langle \mathbf{x}, \mathbf{a}_i \rangle = b_i$ is always satisfied, and under certain separability assumptions on the underlying data, they show that recovery of the optimal solution is possible.

7.5 Rounding

Given a solution $(\Lambda_{**}, \Lambda_{ii})$ of (7.15), we still need to determine a proper partition of $[n]$ through a rounding procedure. Ideally, such an algorithm would aim at a high approximation *guarantee*, but both achieving and proving this can become arbitrarily hard, or may even be impossible, as shown in Theorem 6.14 for the colouring problem.

In the following, we will motivate our rounding procedure for (7.15), leaving the question about provable approximation guarantees open. Before we start with a structural lemma, we will introduce some additional definitions. Let $(\Lambda_{**}, \Lambda_{ii})$ be a solution of (7.15) and define the vectors

$$\boldsymbol{\lambda}_* := \Lambda_{**} \mathbf{e}, \quad \boldsymbol{\lambda}_i := \Lambda_{ii} \mathbf{e} \quad \forall i \in [n].$$

Furthermore, we define the set $\Lambda := \{\boldsymbol{\lambda}_i \mid i \in [n]\}$.

Lemma 7.24:

The inequality $\langle \lambda_i, C_i \lambda_i \rangle \leq \langle C_i, \Lambda_{ii} \rangle$ holds for all $i \in [n]$. Furthermore, strict inequality implies that $(\lambda_* \lambda_*^\top, \lambda_i \lambda_i^\top)$ is infeasible for (7.15).

Proof. By Lemma 7.20 and the Schur complement Lemma 2.5, $\Lambda_{ii} \geq \lambda_i \lambda_i^\top$, and the inequality follows from $C_i \geq \mathbf{0}$ due to self-duality of S_+^n . By our optimality assumption, the set $(\lambda_* \lambda_*^\top, \lambda_i \lambda_i^\top)$ cannot be feasible and have strictly smaller objective value. \square

The set $(\lambda_* \lambda_*^\top, \lambda_i \lambda_i^\top)$ already satisfies several properties of (7.12). In particular, the only constraints that can be violated are

$$\langle (\lambda_*)_{v(s)}, \mathbf{e} \rangle (\lambda_*)_{v(s)} = (\lambda_*)_{v(s)} \quad \forall s \in [q], \forall i \in [n], \quad (7.17a)$$

$$(\lambda_*)_{v(s)} (\lambda_i)_{v(s)}^\top = (\lambda_i)_{v(s)} (\lambda_*)_{v(s)}^\top \quad \forall s \in [q], \forall i \in [n], \quad (7.17b)$$

$$(\lambda_i)_{v(s)} (\lambda_i)_{v(t)}^\top = \mathbf{0} \quad \forall s \neq t \in [q], \forall i \in [n], \quad (7.17c)$$

which all enforce the combinatorial structure of the problem through restricting the support of the vectors.

Due to this fact and the nonlinear nature of the problem, our goal is now to use the set Λ in order to construct a good partition \mathcal{T} for the original problem (7.2) instead of deriving a feasible solution for (R2). In particular, we use the vectors λ_i over their matrices Λ_{ii} , since rounding to their native space Δ^m is much more tractable than rounding to the space of doubly stochastic matrices of rank 1.

To motivate how we construct the partitions, we need to introduce another concept. Given Λ , we define the family $\{\lambda_T \mid T \subseteq [n]\}$ by defining

$$\lambda_T := \bigvee_{i \in T} \lambda_i$$

as the meet of $\{\lambda_i \mid i \in T\}$ in the lattice (\mathbb{R}_+^m, \leq) . In other words, λ_T is maximal with regards to \leq such that $\lambda_T \leq \lambda_i$ for all $i \in T$, and as a byproduct, we also have $\lambda_T \geq \mathbf{0}$. The idea behind this concept is the following observation.

Lemma 7.25:

Let $\mathcal{T} \in \mathcal{P}_k^n$, $\lambda_{T_j} \neq \mathbf{0}$ for all $j \in [k]$ and $C_i \geq \mathbf{0}$ for all $i \in [n]$. Then

$$\sum_{j \in [k]} \sum_{i \in T_j} \langle \lambda_{T_j}, C_i \lambda_{T_j} \rangle \leq \sum_{i \in [n]} \langle C_i, \Lambda_{ii} \rangle \leq \max \{ \langle \lambda_{T_j}, \mathbf{e} \rangle^{-2} \mid j \in [k] \} \cdot \sum_{j \in [k]} \sum_{i \in T_j} \langle \lambda_{T_j}, C_i \lambda_{T_j} \rangle.$$

Proof. The first inequality follows from the assumption $C_i \in S_+^m \cap \mathbb{R}_+^{m \times m}$, since

$$(\Lambda_{ii} - \lambda_i \lambda_i^\top) + (\lambda_i \lambda_i^\top - \lambda_{T_j} \lambda_{T_j}^\top) \in S_+^m + \mathbb{R}_+^{m \times m} = (S_+^m \cap \mathbb{R}_+^{m \times m})^*.$$

For the second inequality, first note that $\frac{\lambda_T}{\langle \lambda_T, \mathbf{e} \rangle} \in \Delta^m$ holds for all $T \subseteq [n]$. We now argue that the set $\left\{ \frac{\lambda_{T_j}}{\langle \lambda_{T_j}, \mathbf{e} \rangle} \mid j \in [k] \right\}$ together with $\mathbf{U}(\mathcal{T})$ can be considered as a feasible point for (R1), showing

$$\sum_{i \in [n]} \langle \mathbf{C}_i, \Lambda_{ii} \rangle \leq \sum_{j \in [k]} \sum_{i \in T_j} \langle \lambda_{T_j}, \mathbf{e} \rangle^{-2} \cdot \langle \lambda_{T_j}, \mathbf{C}_i \lambda_{T_j} \rangle$$

since (7.15) is a relaxation of (R1).

To see this, first note that Definition 7.5 states that the constraint $\lambda^j \in \Delta_{\Omega}^m$ in (R1) is redundant. While this normally means the constraint can be ignored completely without changing the optimal value of (R1), we still need to make sure that the constraint $\langle \lambda^j, \mathbf{e} \rangle = 1$ remains in place, since it was used in the construction of the objective matrices \mathbf{C}_i for homogenization. Since this is the case for $\frac{\lambda_{T_j}}{\langle \lambda_{T_j}, \mathbf{e} \rangle}$, we can set $\lambda^j = \frac{\lambda_{T_j}}{\langle \lambda_{T_j}, \mathbf{e} \rangle}$ and use the assignment matrix $\mathbf{U}(\mathcal{T})$ as feasible point of a slight variant of (R1) with the same optimal value. Finally, the second inequality follows by bounding the terms $\langle \lambda_{T_j}, \mathbf{e} \rangle^{-2}$ from above. \square

The preceding lemma shows that in the case of nonnegative objective \mathbf{C}_i , we can actually give an upperbound for the approximation *quality* in the form of

$$\rho(\mathcal{T}) := \max \{ \langle \lambda_{T_j}, \mathbf{e} \rangle^{-2} \mid j \in [k] \}.$$

Maybe unsurprisingly, there is an inverse relation between the governing parameter $\langle \lambda_T, \mathbf{e} \rangle = \|\lambda_T\|_1$ and the diameter of the underlying set. In particular, we define the diameter $\text{diam}(C)$ of a compact set $C \subseteq \mathbb{R}^d$ as

$$\text{diam}(C) := \max \{ \|\mathbf{x} - \mathbf{y}\|_2 \mid \mathbf{x}, \mathbf{y} \in C \}.$$

However, in order to show this, we first need some additional geometric insights. For this, we first recall that the d -dimensional unit ball with radius r and center $\boldsymbol{\mu}$ is given as

$$B_r^d(\boldsymbol{\mu}) := \{ \mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x} - \boldsymbol{\mu}\|_2 \leq r \}.$$

We can now formulate the following result.

Theorem 7.26 (Jung's Theorem, [DGK63, Thm. 2.6]):

Let $C \subseteq \mathbb{R}^d$ be compact. Then there is $\boldsymbol{\mu} \in \mathbb{R}^d$ and

$$\frac{1}{2} \text{diam}(C) \leq r \leq \sqrt{\frac{d}{2(d+1)}} \text{diam}(C)$$

such that $C \subseteq B_r^d(\boldsymbol{\mu})$.

In what follows, let $H = \{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{e} \rangle = 1\}$ be the unique hyperplane satisfying the inclusion $\Delta^d \subseteq H \subseteq \mathbb{R}^d$. Furthermore, we parametrize the $(d-1)$ -dimensional balls in $H \cong \mathbb{R}^{d-1}$ through their centers by setting

$$B_r^H(\boldsymbol{\mu}) := B_r^d(\boldsymbol{\mu}) \cap H \subseteq H.$$

for all $\boldsymbol{\mu} \in H$.

Lemma 7.27:

The inclusion $B_r^H(\boldsymbol{\mu}) \subseteq \Delta^d$ implies $r \leq (d(d-1))^{-\frac{1}{2}}$.

Proof. Let γ be maximal such that $B_\gamma^H(\boldsymbol{\mu}) \subseteq \Delta^d$. Then by symmetry, we can assume that the center is given by $\boldsymbol{\mu} = \frac{1}{d}\mathbf{e} \in \mathbb{R}^d$, and the d facet-defining inequalities of Δ^d are given by $x_i \geq 0$. It follows then from symmetry that

$$\begin{aligned} \gamma^2 &= \min \{ \|\boldsymbol{\mu} - \mathbf{x}\|_2^2 \mid \mathbf{x} \in \Delta^d, x_i = 0 \} \\ &= \frac{1}{d^2} + \min \{ \|\boldsymbol{\mu}' - \mathbf{x}\|_2^2 \mid \mathbf{x} \in \Delta^{d-1} \} \\ &= \frac{1}{d^2} + \text{dist}(\boldsymbol{\mu}', \Delta^{d-1})^2, \end{aligned}$$

where $\boldsymbol{\mu}' = \frac{1}{d}\mathbf{e} \in \mathbb{R}^{d-1}$. Since \mathbf{e} is orthogonal to H , the orthogonal projection of $\boldsymbol{\mu}'$ onto Δ^{d-1} is $\pi_{\Delta^{d-1}}(\boldsymbol{\mu}') = \frac{1}{d-1}\mathbf{e}$, and the latter minimum is

$$\left\| \left(\frac{1}{d} - \frac{1}{d-1} \right) \mathbf{e} \right\|_2^2 = \frac{1}{(d-1)^2 d^2} \cdot (d-1) = \frac{1}{(d-1)d^2}.$$

As a consequence, $\gamma = \sqrt{\frac{1}{d^2} + \frac{1}{(d-1)d^2}} = (d(d-1))^{-\frac{1}{2}}$. □

Theorem 7.28 (Lower Bound Construction):

Let $\Lambda \subseteq \Delta^d \subseteq \mathbb{R}_+^d$ and denote by $\boldsymbol{\lambda}_\vee$ the meet of Λ in (\mathbb{R}_+^d, \leq) . Then there is a lower bound $\mathbf{0} \leq \boldsymbol{\mu}_+ \leq \boldsymbol{\lambda}_\vee$ such that

$$1 - \frac{d-1}{\sqrt{2}} \cdot \text{diam}(\Lambda) \leq \|\boldsymbol{\mu}_+\|_1 \leq \|\boldsymbol{\lambda}_\vee\|_1.$$

Proof. Applying Theorem 7.26 on $\Lambda \subseteq H \cong \mathbb{R}^{d-1}$ yields a translation $\boldsymbol{\mu}' \in H$ such that

$$\Lambda \subseteq B_r^H(\boldsymbol{\mu}') \subseteq H$$

where

$$r \leq \sqrt{\frac{d-1}{2d}} \cdot \text{diam}(\Lambda).$$

Using translation and scaling operations, we have the equivalence

$$\begin{aligned} B_\gamma^H\left(\frac{1}{d}\mathbf{e}\right) \subseteq \Delta^d &\Leftrightarrow \left(\frac{1}{d}\mathbf{e} - \boldsymbol{\mu}'\right) + B_\gamma^H(\boldsymbol{\mu}') \subseteq \Delta^d &\Leftrightarrow B_\gamma^H(\boldsymbol{\mu}') \subseteq \left(\boldsymbol{\mu}' - \frac{1}{d}\mathbf{e}\right) + \Delta^d \\ &\Leftrightarrow \frac{\gamma}{r} \cdot B_r^H(\boldsymbol{\mu}') \subseteq \left(\boldsymbol{\mu}' - \frac{1}{d}\mathbf{e}\right) + \Delta^d &\Leftrightarrow B_r^H(\boldsymbol{\mu}') \subseteq \frac{r}{\gamma} \left(\boldsymbol{\mu}' - \frac{1}{d}\mathbf{e}\right) + \frac{r}{\gamma} \cdot \Delta^d, \end{aligned}$$

and it follows by Lemma 7.27 that all statements are true for $\gamma = (d(d-1))^{-\frac{1}{2}}$. Denoting $\boldsymbol{\mu} := \frac{r}{\gamma}(\boldsymbol{\mu}' - \frac{1}{d}\mathbf{e})$ and $\tau := \frac{r}{\gamma}$, we thus have the inclusion

$$B_r^H(\boldsymbol{\mu}') \subseteq \boldsymbol{\mu} + \tau \cdot \Delta^d =: D \subseteq H. \quad (7.18)$$

Now let $\boldsymbol{\mu} = \boldsymbol{\mu}_+ + \boldsymbol{\mu}_-$ be the Moreau decomposition of $\boldsymbol{\mu}$ into non-negative and non-positive vectors. Since $\Lambda \subseteq D_+ := D \cap \mathbb{R}_+^d$, we have $\mathbf{0} \leq \boldsymbol{\mu}_+ \leq \boldsymbol{\lambda}$ for all $\boldsymbol{\lambda} \in \Lambda$ and in particular $\mathbf{0} \leq \boldsymbol{\mu}_+ \leq \boldsymbol{\lambda}_V$. As a direct consequence,

$$\|\boldsymbol{\lambda}_V\|_1 = \langle \boldsymbol{\lambda}_V, \mathbf{e} \rangle \geq \langle \boldsymbol{\mu}_+, \mathbf{e} \rangle \geq \langle \boldsymbol{\mu}, \mathbf{e} \rangle = 1 - \tau,$$

where the second equality follows from (7.18). Finally, we have

$$\tau = \frac{r}{\gamma} \leq \frac{d-1}{\sqrt{2}} \cdot \text{diam}(\Lambda).$$

□

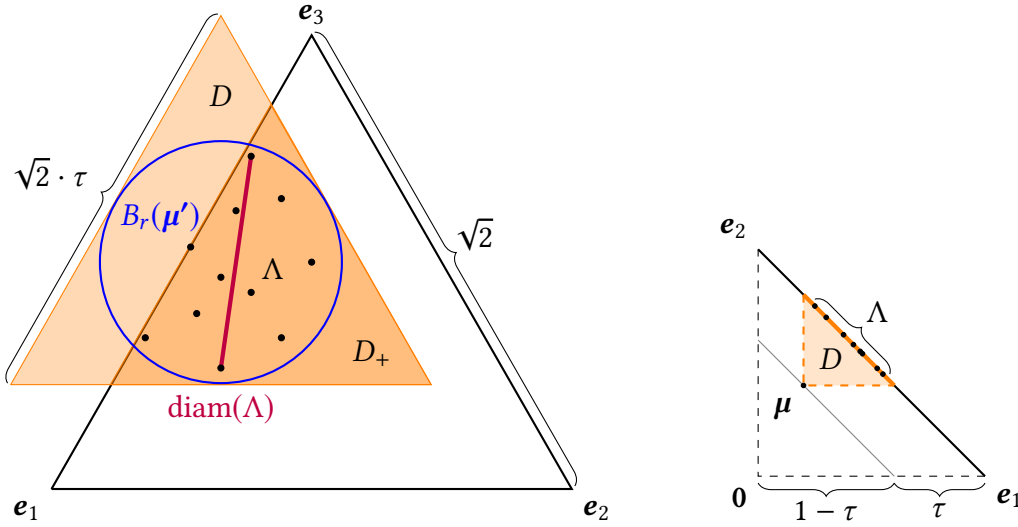


Figure 7.2: Illustration of Theorem 7.28. Left: H for $d = 3$, $\boldsymbol{\mu} \neq \boldsymbol{\mu}_+$. Right: $d = 2$, $\boldsymbol{\mu} = \boldsymbol{\mu}_+$.

Remark 7.29:

It should be noted that the factor $\frac{d-1}{\sqrt{2}}$ in the result of Theorem 7.28 is overly pessimistic and depends heavily on the shape of $\text{conv}(\Lambda)$. In particular, the worst case is achieved when Λ contains the midpoint of every facet of Δ^d , yielding the translation of a rescaled version of $-\Delta^d$, which can be considered as the dual polytope of Δ^d .

However, the constraints (7.17a) actively push the solution away from this configuration, which still has an effect on the relaxations as is implicit in the relations between λ_* and λ_i . Since this influence is hard to quantify, we will just stress that the preceding result describes the worst-case, and actual results may be much better in practice.

Now that we have an interest in using the diameter in place of the quality of the solution, our goal is to find partitions of Λ that minimize the largest diameter of their parts. We will do this indirectly by minimizing a close upperbound on the diameter, which can be derived from the triangle inequality as follows.

Let $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$. Then for any $l \in [n]$, we have

$$\begin{aligned} \text{diam}(X) &= \max \{ \|\mathbf{x}_i - \mathbf{x}_j\|_2 \mid i, j \in [n] \} \\ &\leq \max \{ \|\mathbf{x}_i - \mathbf{x}_l\|_2 + \|\mathbf{x}_l - \mathbf{x}_j\|_2 \mid i, j \in [n] \} \\ &\leq 2 \cdot \max \{ \|\mathbf{x}_i - \mathbf{x}_l\|_2 \mid i \in [n] \} \leq 2 \cdot \text{diam}(X). \end{aligned}$$

In particular, $2 \cdot \max \{ \|\mathbf{x}_i - \mathbf{x}_l\|_2 \mid i \in [n] \}$ is a 2-approximation for every $l \in [n]$, and choosing l to minimize the term will usually give an even better approximation.

Optimizing this 2-approximation of the diameter immediately brings us to the *k-center clustering problem*.

Definition 7.30 (*k-center Clustering*):

Given a set $X = \{\mathbf{x}_i \mid i \in [n]\} \subseteq \mathbb{R}^d$, the *k-center clustering problem* is defined as

$$C_\infty(X, k) := \min \left\{ \max_{i \in [n]} \text{dist}(\mathbf{x}_i, Y) \mid Y \subseteq X, |Y| = k \right\}. \quad (7.19)$$

Fortunately, this is a well-studied problem with efficient and *deterministic* heuristics despite it being NP-hard.

Algorithm 7.1: Farthest Point Clustering (FPN)

Data: Data $X = \{\mathbf{x}_i \mid i \in [n]\} \subseteq \mathbb{R}^d, k \in [n]$

Result: Centers $Y \subseteq X, |Y| = k$

```

1  $b \leftarrow \infty, Y \leftarrow \emptyset;$ 
2 for  $i \in [n]$  do
3    $Y_i \leftarrow \emptyset, \mathbf{y}_1 \leftarrow \mathbf{x}_i;$ 
4   for  $j \in [k]$  do
5      $Y_i \leftarrow Y_i \cup \{\mathbf{y}_j\};$ 
6      $\mathbf{y}_{j+1} \leftarrow \text{argmax} \{ \text{dist}(\mathbf{x}, Y_i) \mid \mathbf{x} \in X \};$ 
7   if  $\text{dist}(\mathbf{y}_{k+1}, Y_i) < b$  then
8      $b \leftarrow \text{dist}(\mathbf{y}_{k+1}, Y_i);$ 
9      $Y \leftarrow Y_i;$ 
10 return  $Y;$ 

```

Algorithm 7.1 greedily builds the set of cluster centers Y by iteratively choosing those points which are farthest away from all prior centers. As initialization, every point is chosen as the first cluster center once and the best overall result is kept as the output of the algorithm. Given its output $Y = \{y_1, \dots, y_k\}$, we can construct a partition $\mathcal{T}(Y) = \{T_1, \dots, T_k\} \in \mathcal{P}_k^n$ by setting

$$T_j = \{i \in [n] \mid \text{dist}(\mathbf{x}_i, Y) = \|\mathbf{x}_i - y_j\|_2\}.$$

As explained by the following theorem, Algorithm 7.1 is optimal for this problem.

Theorem 7.31 (Approximating k -Center clustering [HS85]):

For $d > 2$, achieving an approximation ratio for (7.19) better than 2 is NP-hard, and a 2-approximation is given by Algorithm 7.1.

Our rounding procedure can now be summarized as follows.

Algorithm 7.2: Deterministic Rounding

Data: Solution $(\Lambda_{**}, \Lambda_{ii})$ of (7.15).

Result: Partition \mathcal{T} for (7.2).

- 1 set $\Lambda = \{\Lambda_i \mathbf{e} \mid i \in [n]\}$;
 - 2 compute $Y = \text{FPN}(\Lambda, k)$;
 - 3 **return** $\mathcal{T}(Y)$;
-

7.6 Modifications

When we introduced simplicial covers in Section 7.2, we focused on a basic setting for the sake of exposition. In this section, we return to some underlying assumptions and see how they might be generalized in order to extend the framework suggested in this chapter. In particular, we will focus on the structure of simplicial covers, including the choices of \mathbf{V} and Ω , as well as the objective function.

For most of these generalizations, the constructions in Sections 7.2–7.3 can be carried through analogously, and we will only comment on the results that break down.

Removing redundant vertices

In Section 7.2, there was no assumption on the uniqueness of the columns in \mathbf{V} , and, as pointed out by Remark 7.9, if simplices shared vertices, we used one copy of each vertex for each simplex. The underlying idea was to ensure the blockdiagonal structure of Ω , which in turn is the foundation of Lemma 7.13 and all of its consequences throughout Section 7.2.

However, we can arrive at a slight variation of (R2) even without this assumption. To see this, we merely need to change the definition of the index sets $v(s)$ and adjust Ω

and λ_* accordingly. Given a matrix $\mathbf{V} = (\mathbf{V}_1, \dots, \mathbf{V}_q)$ with multiple identical columns, let \mathbf{V}' collect all the unique columns of \mathbf{V} . Then let $v(s)$ denote the unique set of indices in \mathbf{V}' corresponding to the columns of \mathbf{V}_s and let Ω be the binary matrix which is one everywhere, except for the principal submatrices indexed by each $v(s)$.

Remark 7.32:

If \mathcal{P} defines a simplicial complex with corresponding independence system L as in Remark 7.6, then the new matrix Ω is equal to the matrix Ω_L arising from L as defined in Section 2.2. This is analogous to the connection described in Remark 7.10 and can be seen as generalization.

As a consequence of these new index sets, we might have $\langle (\lambda_*)_{v(s)}, \mathbf{e} \rangle > 1$ for some $s \in [q]$, as can be observed in Example 7.33. To resolve this, we have to weaken (R2) by only assuming inequality instead of equality in

$$\langle (\lambda_*)_{v(s)}, \mathbf{e} \rangle (\lambda_*)_{v(s)} \geq (\lambda_*)_{v(s)}.$$

Then, by construction,

$$\min \left\{ \sum_{i \in [n]} \langle \lambda_i, \mathbf{C}_i \lambda_i \rangle \left| \begin{array}{ll} \lambda_* \in k \cdot \Delta^m, & \langle (\lambda_*)_{v(s)}, \mathbf{e} \rangle (\lambda_*)_{v(s)} \geq (\lambda_*)_{v(s)} \quad \forall s \in [q], \\ \lambda_i \in \Delta_{\Omega}^m, & \forall i \in [n], \\ \lambda_i \leq \lambda_* & \forall i \in [n] \end{array} \right. \right\} \quad (7.20)$$

is still a relaxation of the original problem. However, an analogue to Theorem 7.16 fails to hold in this situation, as showcased by the following example.

Example 7.33:

Consider the 1-dimensional Euclidean 2-clustering problem from Subsection 6.3.2 for the set of points $\{b_1, b_2, b_3, b_4\} = \{-0.6, -0.4, 0.4, 0.6\}$. Enumerating all partitions, it can be verified that the optimal solution is achieved with $x_1 = -0.5$ and $x_2 = 0.5$ for the partition $\{\{1, 2\}, \{3, 4\}\}$ and has value 0.04.

Let the simplicial cover $\mathcal{P} = P_1 \cup P_2 \subseteq \mathbb{R}$ be given by $P_1 = [-1, 0]$ and $P_2 = [0, 1]$. Then \mathcal{P} is separating, and by Theorem 7.16, the optimal solution can be computed by using $\mathbf{V}_1 = (-1 \ 0)$, $\mathbf{V}_2 = (0 \ 1)$ and $\mathbf{V} = (-1 \ 0 \ 0 \ 1)$ to construct and solve (R2).

Now consider the reduced matrix $V' = (-1 \ 0 \ 1)$ consisting of unique columns. By the previous discussion, the indices are given as $v(1) = \{1, 2\}$ and $v(2) = \{2, 3\}$, and thus

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

Then the vectors

$$(\lambda_*, \lambda_1, \lambda_2, \lambda_3, \lambda_4) = \begin{pmatrix} 0.6 & 0.6 & 0.4 & 0 & 0 \\ 0.8 & 0.4 & 0.6 & 0.6 & 0.4 \\ 0.6 & 0 & 0 & 0.4 & 0.6 \end{pmatrix}$$

are feasible for (7.20) and each λ_i represents its data-point b_i . Then this solution is optimal for (7.20) and has objective value 0, which is strictly lower than the real minimum of 0.04.

The example shows that while this approach reduces m and thus the size of the matrices involved in the MM approach of Section 7.3, the quality of the solutions may drop to a point where they do not yield any actual information about the optimum anymore. Additionally, we lose the computational advantages of the diagonal block structure of Ω , resulting in a weakly coupled system.

Furthermore, the example illustrates the mechanism responsible for the reduced quality of the solution as a “discount” in terms of convex combinations. In particular, raising the upperbound given by λ_* on those vertices occurring in multiple polytopes extends the representable points in all of these polytopes, in contrast to the vertices occurring only once. Consequently, the individual λ_i will have a larger spread, as can be seen by comparing Figures 7.6 and 7.7.

However, the same approach can successfully be used to replace the simplices in our simplicial cover by general polytopes, as can be seen in the next subsection.

Polytope covers

In Section 7.2, we introduced \mathcal{P} as a union of separate simplices whose dimension is fixed to that of the ambient space. We can easily generalize this concepts by allowing any combination of

- a union of general polytopes instead of simplices,
- different dimensions for each of the polytopes.

Both generalizations have in common that in principal, they only impact the number of vertices for each polytope in \mathcal{P} . Consequentially, only the definition of the matrix \mathbf{V} in (7.4) needs to be slightly adjusted, and we can arrive at both (R2) and (7.15) without impacting any of the theoretical results.

For practical computations however, there is a slight difference between allowing lower dimensional simplices and using general polytopes instead of simplices. While the former does not impact the relaxations constructed by MM, the latter suffers from the lack of unique representation, as is implicit in the lack of uniqueness in Theorem 2.8 and showcased in the following example.

Example 7.34:

Consider the unit square $C^2 = [0, 1]^2$ and a point $\mathbf{x} \in C^2$. Then $\mathbf{V} = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix}$ describes the vertices of C^2 and any vector $\boldsymbol{\lambda} \in \Delta^4$ that satisfies $\mathbf{V}\boldsymbol{\lambda} = \mathbf{x}$ represents \mathbf{x} . However,

if $\mathbf{x} \in \text{int}(C^2)$, this representation is not unique, since

$$C^2 = \text{conv} \left(\begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right) \cup \text{conv} \left(\begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix} \right) = \text{conv} \left(\begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix} \right) \cup \text{conv} \left(\begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix} \right)$$

can be decomposed in two different ways, which is visualized by Figure 7.3.

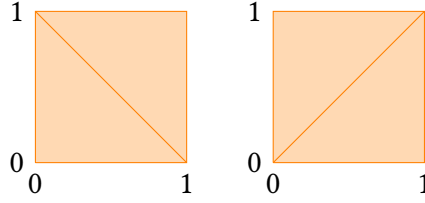


Figure 7.3: Both decompositions of C^2 as union of 2 triangles.

By Caratheodory's theorem, both decompositions define different representations and every convex combination of these representations results in yet another representation.

Most of the time, this ambiguity will weaken the relaxation (7.15) considerably, since there are more degrees of freedom available. Luckily, we can use the ideas of the preceding section to remove the ambiguity by adjusting Ω . To see this, let \mathcal{P} be the union of polytopes P_s for $s \in [q]$ and \mathbf{V}_s their corresponding vertices. Then we can subdivide each polytope P_s into a simplicial complex to rewrite

$$P_s = \bigcup_{i \in q_s} P_{s,i}$$

where each $P_{s,i}$ is a simplex. Then each $P_{s,i}$ defines a unique index-set $v(s, i)$ of the columns of \mathbf{V}_s , and we can define Ω as the binary matrix which is one everywhere, except for the principal submatrices indexed by $v(s, i)$ for both $s \in [q]$ and $i \in [q_s]$. This approach guarantees a unique representation for almost every point in the interior of the polytope P_s . Computationally, this results in a single psd. block of size $|\mathbf{V}_s|$, compared to the unique representation achieved by using the q_s blocks of size $d + 1$ for each simplex $P_{s,i}$.

Example 7.35:

Consider $\mathcal{P} = \bigcup_{s \in [4]} P_{z_s} \subseteq \mathbb{R}^2$ where $P_{z_s} = z_s + C_2$ for $C_2 = [0, 1]^2$ and

$$(z_1 \quad z_2 \quad z_3 \quad z_4) = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix}.$$

Then

$$\begin{aligned} \mathbf{V} &= (\mathbf{V}_1 \mid \mathbf{V}_2 \mid \mathbf{V}_3 \mid \mathbf{V}_4) \\ &= \left(\begin{array}{cc|cc|cc|cc} 0 & 1 & 0 & 1 & 1 & 2 & 1 & 2 & 0 & 1 & 0 & 1 & 1 & 2 & 1 & 2 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 1 & 2 & 2 & 1 & 1 & 2 & 2 \end{array} \right) \in \mathbb{R}^{2 \times 16} \end{aligned}$$

and since each polytope is a translation of C_2 , we can choose the same simplicial decomposition for each polytope. We thus divide each square into a lower-left and upper-right triangle, as depicted by Figure 7.4. Consequentially, we set

$$\Omega = \mathbf{J}_{16} - \mathbf{I}_4 \otimes \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{pmatrix}$$

to get a unique representation involving 4 psd. blocks of size 4 each. This is a computational trade-off, since using the simplicial decomposition directly would have resulted in 8 sdp blocks of size 3.

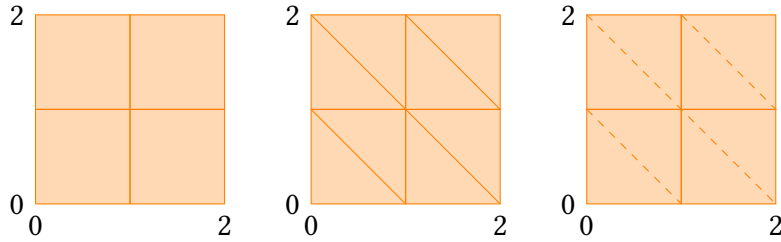


Figure 7.4: Left: Original polytope cover consisting of 4 squares. Middle: Simplicial decomposition using 8 triangles. Right: Polytope cover with implied constraints from simplicial decomposition.

Describing the σ -skeleton

As a special case of the previous subsection, \mathcal{P} might be given as the σ -skeleton of a polytope, given by the following definition.

Definition 7.36:

Let $P \subseteq \mathbb{R}^d$ be a polytope with vertex set \mathbf{V} and $\sigma \in [d]$. Define $G_\sigma(P)$ as the graph with vertex set \mathbf{V} where the edge (v_1, v_2) is contained if and only if the line segment $\text{conv}(\{v_1, v_2\})$ is contained in a face of P with dimension at most σ . Then $\mathbf{A}_{\overline{G_\sigma(P)}}$ denotes the adjacency matrix of the complementary graph and the σ -skeleton $\text{skel}_\sigma(P)$

of P is formally defined as

$$\text{skel}_\sigma(P) := \left\{ \mathbf{x} \in \mathbb{R}^d \mid \mathbf{x} = \mathbf{V}\boldsymbol{\lambda}, \boldsymbol{\lambda} \in \Delta_{\frac{m}{G_\sigma(P)}}^m, \|\boldsymbol{\lambda}\|_0 \leq \sigma + 1 \right\}, \quad (7.21)$$

which is the union of all faces of P of dimension at most σ .

Example 7.37:

The unit square C^2 is given as the convex hull of $\mathbf{V} = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix}$ and $\text{skel}_1(C^2)$ consists of 4 line segments, as shown in Figure 7.5. Choosing $\mathcal{P} = \text{skel}_1(C^2)$ we need 4 2-simplices and consequently $m = 8$ vertices for using the polytope cover approach from Section 7.6. Using the sparsity constraint $\|\boldsymbol{\lambda}\|_0 \leq 2$ in (7.21) allows us to use each vertex only once to end up with $m = 4$ instead.

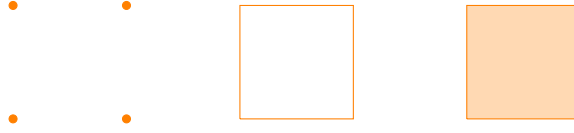


Figure 7.5: The σ -skeletons of C_2 . As σ ranges from 1 to 3, the σ -skeleton describes the union of vertices, edges and the square itself.

Remark 7.38:

The sparse set $\{\boldsymbol{\lambda} \in \Delta^m \mid \|\boldsymbol{\lambda}\|_0 \leq \sigma\}$ can be described by adding the equations

$$\boldsymbol{\lambda}^{e_S} = 0 \quad \forall S \subseteq [m]: |S| \geq \sigma + 1.$$

Using MM, these equations can be incorporated by expanding Ω to sum over higher moments as well. However, this approach becomes impractical very quickly since it requires a stage $t > \sigma$ to work.

It would be interesting to investigate low-degree polynomials as approximations of sparsity constraints.

Simplicial Covers for Semialgebraic Sets

So far, our focus in this section was to modify the construction of the underlying simplicial cover, where we have always assumed that $\text{SE}_B^{\mathcal{A}}$ optimizes a variable $\mathbf{x} \in \mathbb{R}^d$ as in (7.1). However, the powerful results concerning semi-algebraic optimization in Sections 2.6 and 3.1 can be used with our framework to restrict \mathbf{x} to a semialgebraic subset $K \subseteq \mathbb{R}^d$ as well. In particular, let

$$K = \left\{ \mathbf{x} \in \mathbb{R}^d \mid g_i(\mathbf{x}) \geq 0, \quad i \in \mathcal{I} \right\}$$

for multivariate polynomials $g_i \in \mathbb{R}[\mathbf{x}]$. Then each $g_i(\mathbf{x})$ can be easily turned into a polynomial expression $g'(\boldsymbol{\lambda}) = g(\mathbf{V}\boldsymbol{\lambda})$, and we can additionally assume that g' is a homogeneous polynomial of the same degree as g , such that $\mathbf{x} \in K$ where $\mathbf{x} = \mathbf{V}\boldsymbol{\lambda}$ if and only if

$$\boldsymbol{\lambda} \in K' := \{\boldsymbol{\lambda} \in \mathbb{R}^m \mid g'_i(\boldsymbol{\lambda}) \geq 0, \quad i \in \mathcal{I}\}.$$

Consequentially, it suffices to add the membership $\boldsymbol{\lambda}_i \in K'$ for all $i \in [n]$ to problem (R2). Using MM, we then arrive at (7.14) with additional constraints. In general, we only need to cover K by a set of simplices to use our approach. In particular, there is no need for a simplicial cover to *approximate* K with those simplices, as long as K is *covered* by them.

Example 7.39:

Let K be the unit sphere, then $K = \{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{x} \rangle = 1\}$ is described by the homogeneous quadratic equation $\langle \mathbf{x}, \mathbf{x} \rangle = 1$. Substituting $\mathbf{x} = \mathbf{V}\boldsymbol{\lambda}$ yields again a quadric defined by $\boldsymbol{\lambda}^\top (\mathbf{V}^\top \mathbf{V}) \boldsymbol{\lambda} = 1$. Adding this constraint to (7.14) then turns into

$$\langle \Lambda_{ii}, \mathbf{V}^\top \mathbf{V} - \mathbf{J}_m \rangle = 0 \quad \forall i \in [n],$$

where we used homogenization.

Example 7.40:

Let K be the set of vectorized orthogonal 2×2 matrices \mathbf{X} in \mathbb{R}^4 given by the map

$$\mathbf{X} = \begin{pmatrix} x_1 & x_{12} \\ x_{21} & x_2 \end{pmatrix} \mapsto (x_1 \quad x_{12} \quad x_{21} \quad x_2)^\top = \mathbf{x}.$$

Then K can be described as the zeros of four quadratic equations, one for each entry of $\mathbf{X}\mathbf{X}^\top = \mathbf{I}_2$. Denoting them by $\mathbf{x}^\top \mathbf{Q}_l \mathbf{x} = q_l$, substituting $\mathbf{x} = \mathbf{V}\boldsymbol{\lambda}$ yields again an intersection of quadrics defined by $\boldsymbol{\lambda}^\top (\mathbf{V}^\top \mathbf{Q}_l \mathbf{V}) \boldsymbol{\lambda} = q_l$, which turn into the homogeneous equations

$$\langle \Lambda_{ii}, \mathbf{V}^\top \mathbf{Q}_l \mathbf{V} - q_l \mathbf{J}_m \rangle = 0 \quad \forall i \in [n]$$

for (7.14).

However, we note that finding a separating simplicial cover may prove much harder when K has a complex geometry. In particular, while MM will converge towards feasibility in K , we necessarily need to start at a stage determined by the maximum degree of a polynomial describing K to avoid crude approximations.

Clustering Varieties

The task of (7.2) is to partition affine subspaces in such a way that the minimizer of $\text{SE}_B^{\mathcal{A}}$ is as close as possible to all of them. Ideally, this would mean that we get a common intersection of the corresponding affine subspaces, and so we can consider (7.2)

as a generalized common intersection point. Consequentially, we can generalize this concept from subspaces to arbitrary sets and try to adjust our approach accordingly.

An extension to the case of varieties can be done in the following way. We replace $\mathbf{A}_i \mathbf{x} - \mathbf{b}_i$ in (7.1) with $F_i(\mathbf{x})$, where $F_i \in \mathbb{R}[\mathbf{x}]$ is a multivariate polynomial describing the variety $\mathcal{V}_{\mathbb{R}}(\|F_i\|_2^2)$. Following Section 7.2, we may replace \mathbf{x} by $\mathbf{V}\boldsymbol{\lambda}$ and homogenize $\|F_i(\mathbf{V}\boldsymbol{\lambda})\|_2^2$ using $\langle \boldsymbol{\lambda}_j, \mathbf{e} \rangle = 1$ to end up with a variant of (R2), where the objective function has been replaced. The results from Section 7.3 follow according to this replacement, but the stage of MM needs to be twice the maximum degree of any F_i .

Regularization

We want to point out that the reformulation (R2) never explicitly uses the number k except in the equation $\langle \boldsymbol{\lambda}_*, \mathbf{e} \rangle = k$ implicit in $\boldsymbol{\lambda}_* \in k \cdot \Delta^m$. We can thus treat k throughout as a variable instead and include a function of k in the objective function to dynamically search for the number of sets contained in a desired partition. The difficulty then lies in choosing this function in a way to retain a useful solution to the original problem.

7.7 Applications

We will state various applications and visualize the resulting relaxations.

Euclidean Clustering

By choosing $\mathbf{A}_i = \mathbf{I}_n$ in (7.2), we recover the classical problem (6.7) of Euclidean clustering for the points $\{\mathbf{b}_1, \dots, \mathbf{b}_n\} \subseteq \mathbb{R}^d$. As pointed out in Corollary 6.28, it is well known that $\text{conv}(\{\mathbf{b}_1, \dots, \mathbf{b}_n\}) \subseteq \mathcal{P}$ will contain all optimal solutions. In particular, for the simplicial cover assumption it suffices that \mathcal{P} covers a box which includes all $\{\mathbf{b}_1, \dots, \mathbf{b}_n\}$, which can be easily extracted.

We can use Euclidean clustering to get a better intuition of how the relaxation works. Regarding the choice of \mathcal{P} , consider Figure 7.6. Using any simplex containing all the points is the coarsest approximation but yields useless results, since each local estimate $V\boldsymbol{\lambda}_i$ can be chosen as \mathbf{b}_i .

In view of Theorem 7.16, the algorithm will perform best if the simplicial cover is separating, which means that the cluster centers are separated by the polytopes in \mathcal{P} . This suggests that there should be at least k simplices, and that an oversegmentation removes the need of knowing the location of the centers in advance, as can be observed in Figure 7.6 as well. Some of the variants from Section 7.6 that change how \mathbf{V} and $\boldsymbol{\Omega}$ are constructed are visualized in Figure 7.7, but compared to Figure 7.6, they do not offer any noteworthy improvement in quality.

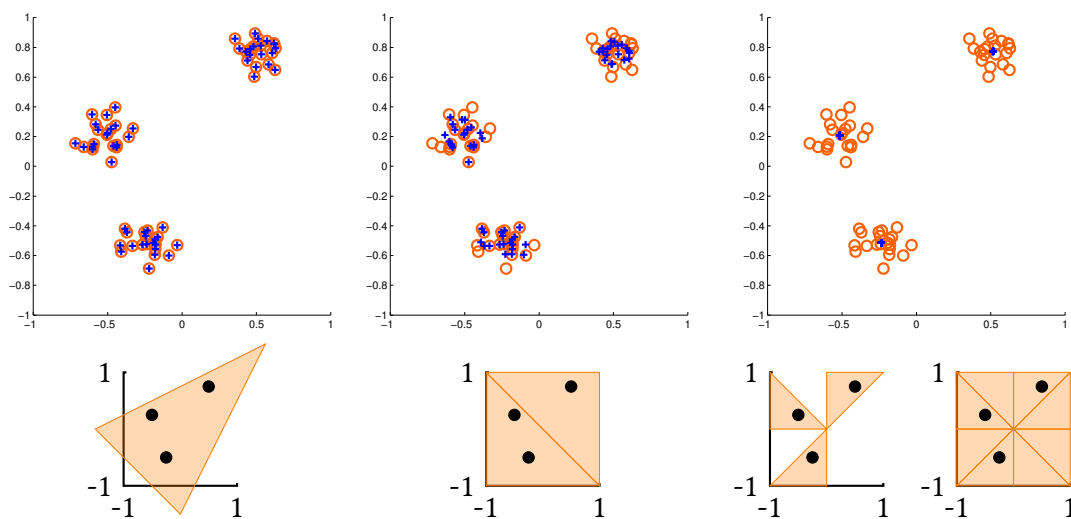


Figure 7.6: *Euclidean clustering* with $d = 2$, $k = 3$, $n = 60$. Top: Circles corresponding to input \mathbf{b}_i and crosses corresponding to points parametrized by the λ_i extracted from (7.15). Bottom: Different choices of \mathcal{P} . From left to right: Minimal cover, non-separating cover, minimal separating cover, oversegmentation. Algorithm 7.2 was able to recover the optimal solution implied by the right plot in each scenario.

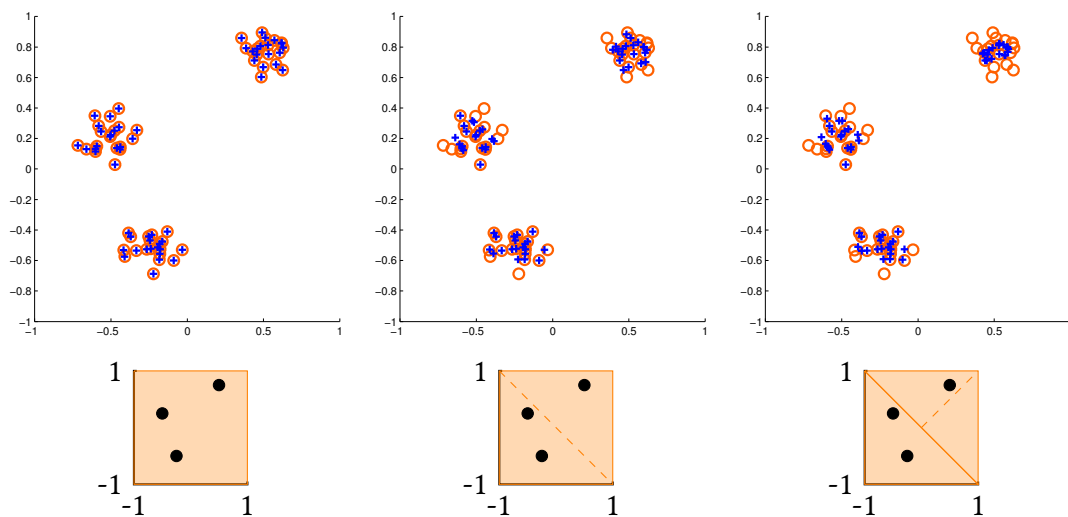


Figure 7.7: *Euclidean Clustering* with $d = 2$, $k = 3$, $n = 60$. Top: Circles corresponding to input \mathbf{b}_i and crosses corresponding to points parametrized by the λ_i extracted from (7.15). Bottom: Variants described in Section (7.6) for describing the feasible set $[-1, 1]^2$. From left to right: \mathcal{P} is the union of 1, 2, 3 polytopes respectively. Vertices at dashed lines are unique rows in \mathbf{V} and used in each bordering polytope.

As pointed out by Remark 7.17, we can easily restrict the feasible set in a way to force the optimal solution into specific regions. For example, by choosing each polytope in \mathcal{P} to be a single vertex, we reduce (7.15) to an LP which aims to choose an optimal collection of locations from a discrete set of points, as can be seen in Figure 7.8. In this case, our experiments always returned the optimal solution.

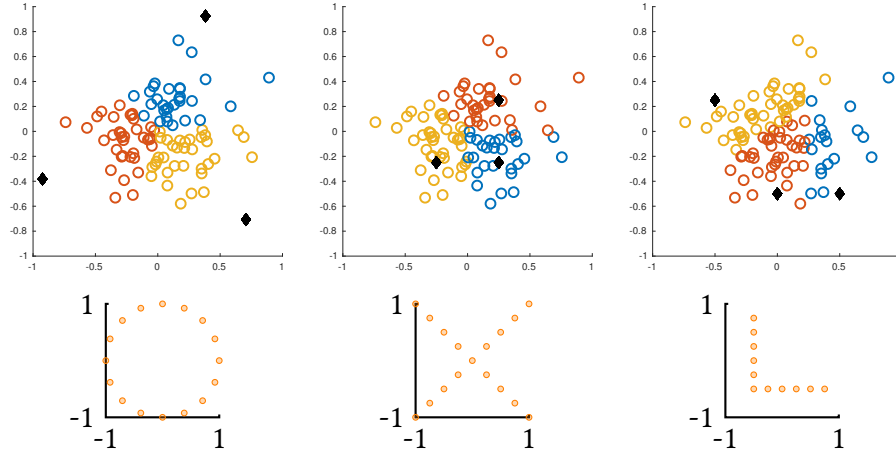


Figure 7.8: Euclidean k -clustering with $d = 2$, $k = 3$, $n = 100$ restricted to discrete \mathcal{P} . Top: Circles corresponding to input \mathbf{b}_i , diamonds corresponding to centers and colours corresponding to clusters. Bottom: Different choices of discrete \mathcal{P} .

Hyperplane Clustering

By choosing $\mathbf{A}_i = \mathbf{a}_i$ as row vectors in \mathbb{R}^d and setting $b_i = 0$, (7.3) becomes the problem of choosing minimal $\langle \mathbf{a}_i, \mathbf{x}_j \rangle^2$ terms. We can interpret this as simultaneously choosing k hyperplanes parameterized by their normal vectors \mathbf{x}_j and assigning the points \mathbf{a}_i to them according to their weighted angle.

We can uniquely parametrize these hyperplanes by choosing an element \mathbf{x} of their complement space which satisfies membership in both $S_{\|\cdot\|}^d = \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\| = 1\}$ for any fixed norm $\|\cdot\|$ and the 'upper halfspace' $H_+^d = \{\mathbf{x} \in \mathbb{R}^d \mid x_1 \geq 0\}$.

Note that the norm will weight each point $\mathbf{x} \in S_{\|\cdot\|}^d \cap H_+^d$ by $\|\mathbf{x}\|$. We will also write S_p^d for the $\|\cdot\|_p$ norms. In particular, even though any polyhedral approximation of $S_2^d \cap H_+^d$ corresponds to a norm and can be used as \mathcal{P} , this will introduce a slight bias. The application of this approach is illustrated by Figure 7.9.

As an application of Example 7.39, we can also work directly with $S_2^2 \cap H_+^2$ by adding the quadratic constraint $\langle \mathbf{x}_j, \mathbf{x}_j \rangle = 1$ and choosing $\mathcal{P} \subseteq H_+^2$. Since S_2^2 is not polyhedral, we need to use 3-dimensional simplices for \mathcal{P} in Figure 7.10 whereas 2-dimensional simplices suffice for the polyhedral approximation shown by Figure 7.9.

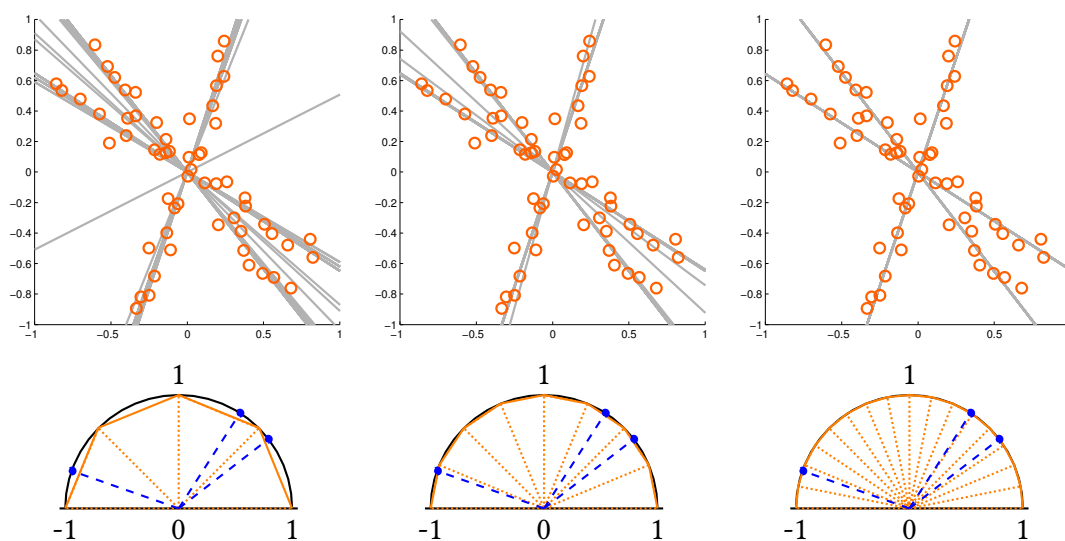


Figure 7.9: *Hyperplane k -clustering* with $d = 2$, $k = 3$, $n = 60$. Top: Circles corresponding to input \mathbf{a}_i and lines parametrized by λ_i extracted from (7.15). Bottom: Approximations of $S_2^2 \cap H_+^2$ by polygonal lines. For better visibility the ends of each line segment are connected to the origin with a dotted line. Dashed lines end in optimal angles.

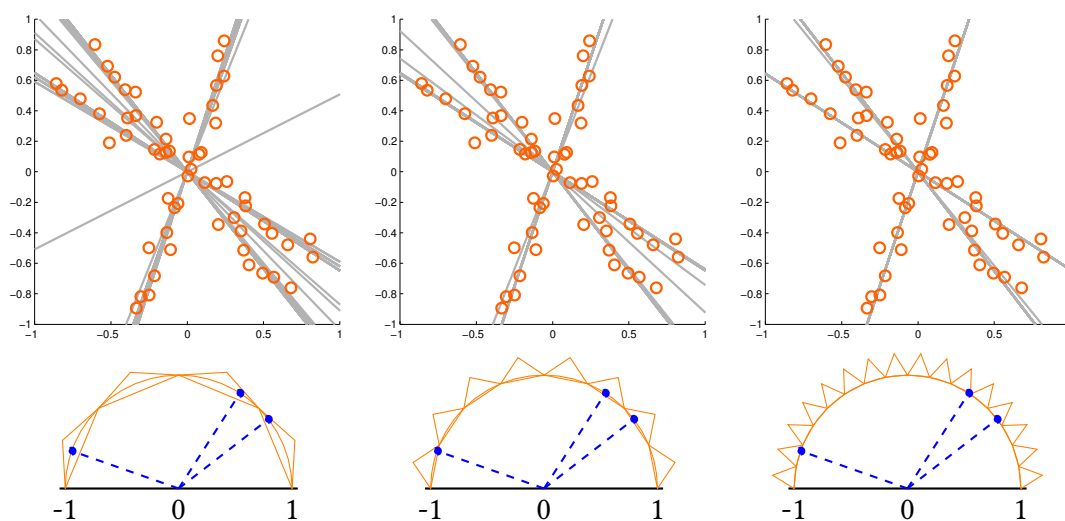


Figure 7.10: *Hyperplane Clustering* with $d = 2$, $k = 3$, $n = 60$. Top: Circles corresponding to input \mathbf{a}_i and lines parametrized by λ_i extracted from (7.15). Bottom: The semicircle $S_2^2 \cap H_+^2$ is covered with triangles in \mathcal{P} . Dashed lines end in optimal angles.

Mixed Linear Regression

We can easily extend the hyperplane clustering to the more general case of *mixed linear regression* that we mentioned already in Section 7.4.2 by changing to homogeneous coordinates. In particular, we can encode the data point $\mathbf{a}_i \in \mathbb{R}^d$ as $(\mathbf{a}_i, 1) \in \mathbb{R}^{d+1}$ and try to find a hyperplane orthogonal to $(\mathbf{x}_j, -z_j) \in \mathbb{R}^d$ to get the minimization of terms like

$$\langle (\mathbf{a}_i, 1), (\mathbf{x}_j, z_j) \rangle^2 = (\langle \mathbf{a}_i, \mathbf{x}_j \rangle - z_j)^2,$$

which approximate membership in the affine hyperplane

$$\mathbf{a}_i \in H_{(\mathbf{x}_j, z_j)} = \{ \mathbf{a} \in \mathbb{R}^d \mid \langle \mathbf{a}, \mathbf{x}_j \rangle = z_j \}.$$

Since the manipulation only amounts to lifting the input data $\{\mathbf{a}_i\}_{i \in [n]}$, this is just an instance of the hyperplane clustering problem in a space with dimension increased by 1. In particular, the problem of Figure 7.11 can be computed as an instance of clustering points from \mathbb{R}^3 into 2-dimensional hyperplanes.

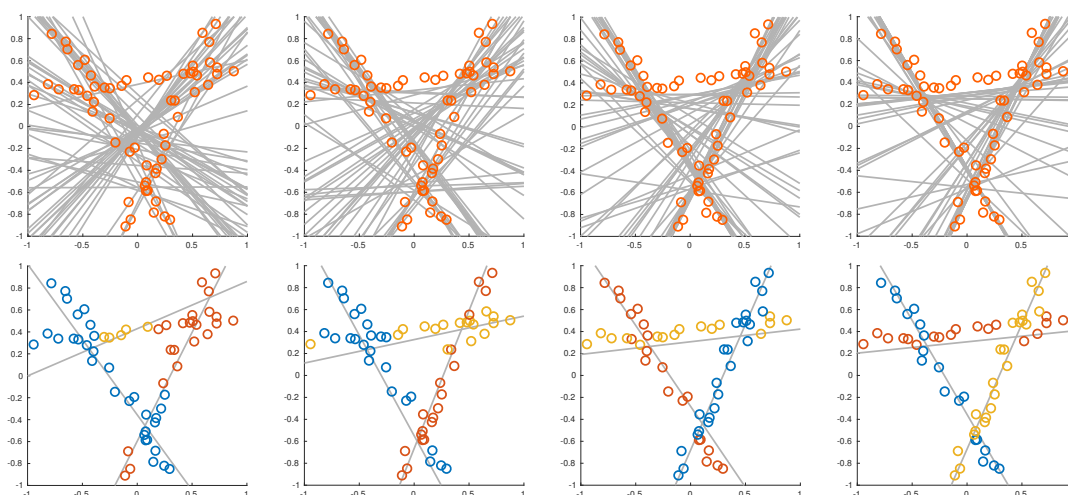


Figure 7.11: *Mixed Linear Regression* with $d = 2$, $k = 3$, $n = 60$ as a special case of *Hyperplane Clustering* with $d = 3$. Top: Circles corresponding to input \mathbf{a}_i and lines parametrized by λ_i extracted from (7.15). Bottom: Gray lines corresponding to rounded solution of (7.15) and coloured data points according to the extracted clustering. From left to right: Discretization of $(S_2^2 \cap H_+^2) \times [-0.3, 0.3]$ into (2×8) , (4×4) , (8×2) and (8×8) line segments, where $S_2^2 \cap H^2$ is approximated like in Figure 7.9.

Chapter 8

Conclusion

In this thesis, different representations for partitions were formulated and used to approximate partition problems. While doing so, special focus was on finding representations that are invariant under permutations of the individual parts, with the intention of evading the computational problems commonly associated with the lack of such invariance. A major tool in doing so was the method of moments, which was used to construct convex relaxations of the underlying polynomial optimization problems that can be solved by semidefinite optimization techniques.

In the first part of the thesis, we set up the machinery surrounding the method of moments in Chapters 2 and 3, while providing an explicit way of solving the underlying optimization problems. The algorithm we presented for this was the CLP-Newton method, which we introduced as a generalization of the LP-Newton method for self-dual cones. Experiments with the CLP-Newton method were provided, which indicated that, while the computation of a Newton step is quite expensive, only a small number of Newton steps were needed to come close to the optimal solution.

The second part of the thesis was concerned with classical partition problems that can be represented using various classes of matrices. For this, we started with assignment matrices and the related theory of orbitopes in Chapter 4, showcasing the problems that can arise in the convex setting when representations of partitions are not invariant under permutations of their parts.

We proceeded by showing how assignment matrices naturally give rise to partition matrices in Chapter 5, which explicitly described a connection that has been implicitly used in literature for a long time. Furthermore, the class of combinatorial projection matrices was introduced in Chapter 6, and its close connection to partitions matrices was established. As a main result, we used combinatorial projection matrices to give a new formulation of the colouring number of graphs and compared the associated relaxations to variants of the Lovász theta number. We were able to show that while the Szegedy number, a slight improvement of the Lovász theta number, is always at least as

good as our relaxation, both relaxations agree on the difficult class of vertex-transitive graphs. In particular, this result implies that relaxations for binary eigenvalues perform worse than relaxations of binary matrix-entries for this problem.

The third and final part of the thesis treated the challenging problem of affine Euclidean k -clustering in Chapter 7, which, to the best of our knowledge, has not been investigated in this generality prior to this thesis. While the classical Euclidean k -clustering allows for a reformulation based on the optimality condition of a subset of the variables, the generalized affine Euclidean k -clustering lacks this property, making it much harder to solve and relax. Assuming a certain separability property on the feasible set associated with the same subset of variables, we proposed a new approach for this problem. Using a simplicial cover of this feasible set, it was shown how to construct a reformulation that can be relaxed to yield bounds on the global optimum. In particular, it was shown how to extend this method to several variations of the affine Euclidean k -clustering problem that contain additional constraints, and the method was shown to work on a number of several different toy examples.

Further Research Directions

At the conclusion of this thesis, several intriguing questions are still open.

On the numerical side, the CLP-Newton method could be improved by using a more efficient way to solve the minimum-norm-point algorithm. While all relaxations in this thesis would benefit from a computational improvement of the underlying CLP solvers, the affine Euclidean k -clustering problem would benefit the most, given that the scaling of the method using simplicial covers makes the method currently prohibitive for practical application. In particular, the runtime of the method might be improved by utilizing recently published SOCP-hierarchies instead of the classical SDP-hierarchies associated with the method of moments. In general, switching to SOCP-hierarchies tends to improve the runtime considerably at only a marginal cost in solution quality.

On the theoretical side, it would be interesting to get a better understanding about the difference between partition matrices and combinatorial projection matrices in terms of relaxation quality. Any insights on this matter might help us to answer the deeper question of whether binary matrix entries or binary eigenvalues are a better modeling paradigm for relaxation based approaches. For the graph colouring problem in particular, it would be interesting to know if there is a class of additional constraints that could be added to combinatorial projection matrices to reduce the gap to the Szegedy number.

In terms of the affine Euclidean k -clustering, it would be helpful to know whether conditions can be given that guarantee a fixed approximation ratio of our method.

While this is unlikely for the general case due to the high expressiveness of the problem, there might be mild conditions under which certain problems have provable guarantees, even in the first stage of the hierarchy we considered. From a practical perspective, the separability assumption is mostly unexplored, and a list of problems with this property would be desirable to get further insights into the applicability of our method.

Bibliography

- [AB09] Sanjeev Arora and Boaz Barak. *Computational Complexity: A Modern Approach*. Cambridge University Press, 2009.
- [ACKS15] Pranjali Awasthi, Moses Charikar, Ravishankar Krishnaswamy, and Ali Kemal Sinop. The hardness of approximation of euclidean k-means. *CoRR*, abs/1502.03316, 2015.
- [ADHP09] Daniel Aloise, Amit Deshpande, Pierre Hansen, and Preyas Popat. Np-hardness of euclidean sum-of-squares clustering. *Machine Learning*, 75(2):245–248, 2009.
- [AL12] Miguel F Anjos and Jean B Lasserre. *Handbook on semidefinite, conic and polynomial optimization, International Series in Operations Research & Management Science, vol. 166*. Springer, New York, 2012.
- [Bac11] Francis R. Bach. Learning with submodular functions: A convex optimization perspective. *CoRR*, abs/1111.6453, 2011.
- [BEF00] Benno Büeler, Andreas Enge, and Komei Fukuda. *Exact Volume Computation for Polytopes: A Practical Study*, pages 131–154. Birkhäuser Basel, Basel, 2000.
- [Ber75] Sven Berg. Some properties and applications of a ratio of stirling numbers of the second kind. *Scandinavian Journal of Statistics*, 2(2):91–94, 1975.
- [BPS66] Richard A Brualdi, Seymour V Parter, and Hans Schneider. The diagonal equivalence of a nonnegative matrix to a stochastic matrix. *Journal of Mathematical Analysis and Applications*, 16(1):31 – 50, 1966.
- [BPT13] Grigoriy Blekherman, Pablo A. Parrilo, and Rekha Thomas, editors. *Semidefinite Optimization and Convex Algebraic Geometry*. SIAM, 2013.
- [BS94] Mihir Bellare and Madhu Sudan. Improved non-approximability results. In *Proceedings of the Twenty-sixth Annual ACM Symposium on Theory of Computing*, STOC '94, pages 184–193. ACM, 1994.

- [CLO15] David A. Cox, John Little, and Donal O’Shea. *Ideals, Varieties, and Algorithms: An Introduction to Computational Algebraic Geometry and Commutative Algebra*. Springer, 4th edition, 2015.
- [DGK63] L. Danzer, B. Grünbaum, and V. Klee. *Helly’s Theorem and Its Relatives*. Proceedings of symposia in pure mathematics: Convexity. American Mathematical Society, 1963.
- [Die05] Reinhard Diestel. *Graph Theory (Graduate Texts in Mathematics)*. Springer, August 2005.
- [DR07] Igor Dukanovic and Franz Rendl. Semidefinite programming relaxations for graph coloring and maximal clique problems. *Math. Program.*, 109(2):345–365, 2007.
- [EZC10] Ernie Esser, Xiaoqun Zhang, and Tony F. Chan. A general framework for a class of first order primal-dual algorithms for convex optimization in imaging science. *SIAM J. Imaging Sciences*, 3(4):1015–1046, 2010.
- [FHYZ08] Satoru Fujishige, Takumi Hayashi, Kei Yamashita, and Uwe Zimmermann. Zonotopes and the LP-Newton method. *Optimization and Engineering*, 10(2):193–205, 2008.
- [FK09] Yuri Faenza and Volker Kaibel. Extended formulations for packing and partitioning orbitopes. *Math. Oper. Res.*, 34(3):686–697, 2009.
- [GL08] Nebojsa Gvozdenovic and Monique Laurent. The operator ψ for the chromatic number of a graph. *SIAM Journal on Optimization*, 19(2):572–591, 2008.
- [Har68] Bernard Harris. Statistical inference in the classical occupancy problem unbiased estimation of the number of classes. *Journal of the American Statistical Association*, 63:837–847, 1968.
- [HJ16] Paul Hand and Babhru Joshi. A convex program for mixed linear regression with a recovery guarantee for well-separated data. *arXiv preprint arXiv:1612.06067*, 2016.
- [HS85] Dorit S. Hochbaum and David B. Shmoys. A Best Possible Heuristic for the k -Center Problem. *Mathematics of Operations Research*, 10(2):180–184, May 1985.
- [Jag13] Martin Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In Sanjoy Dasgupta and David Mcallester, editors, *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*,

- volume 28, pages 427–435. JMLR Workshop and Conference Proceedings, 2013.
- [KP08] Volker Kaibel and Marc E. Pfetsch. Packing and partitioning orbitopes. *Math. Program.*, 114(1, Ser. A):1–36, 2008.
- [Las15] Jean Bernard Lasserre. *An Introduction to Polynomial and Semi-Algebraic Optimization*, volume 52. Cambridge University Press, 2015.
- [Lau03] Monique Laurent. A comparison of the Sherali-Adams, Lovász-Schrijver, and Lasserre relaxations for 0-1 programming. *Mathematics of Operations Research*, 28(3):470–496, 2003.
- [Lov79] László Lovász. On the shannon capacity of a graph. *IEEE Trans. Inf. Theor.*, 25(1):1–7, 1979.
- [Lov87] László Lovász. *An algorithmic theory of numbers, graphs and convexity*, volume 50. SIAM, 1987.
- [NT08] Arkadi S. Nemirovski and Michael J. Todd. Interior-point methods for optimization. *Acta Numerica*, 17:191–234, 2008.
- [OCPB16] Brendan O’Donoghue, Eric Chu, Neal Parikh, and Stephen Boyd. Conic optimization via operator splitting and homogeneous self-dual embedding. *Journal of Optimization Theory and Applications*, 169(3):1042–1068, 2016.
- [PW07] Jiming Peng and Yu Wei. Approximating K -means-type Clustering via Semidefinite Programming. *SIAM J. Optimization*, 18(1):186–205, 2007.
- [PX05] Jiming Peng and Yu Xia. A new theoretical framework for k -means-type clustering. In *Foundations and advances in data mining*, pages 79–96. Springer Berlin Heidelberg, 2005.
- [Ren10] Franz Rendl. Semidefinite relaxations for integer programming. In Michael Jünger, Thomas M. Liebling, Denis Naddef, George L. Nemhauser, William R. Pulleyblank, Gerhard Reinelt, Giovanni Rinaldi, and Laurence A. Wolsey, editors, *50 Years of Integer Programming 1958-2008 - From the Early Years to the State-of-the-Art*, chapter 18, pages 647–686. Springer, 2010.
- [Roc09] R. Tyrrell Rockafellar. *Convex Analysis*. Springer, 2009.
- [Sch79] Alexander Schrijver. A comparison of the Delsarte and Lovász bounds. *Information Theory, IEEE Transactions on*, 25(4):425–429, 1979.

- [Sch03] Alexander Schrijver. *Combinatorial Optimization - Polyhedra and Efficiency*. Springer, 2003.
- [SR16] Francesco Silvestri and Gerhard Reinelt. The LP-Newton method and conic optimization. *ArXiv e-prints*, November 2016.
- [SRS15] Francesco Silvestri, Gerhard Reinelt, and Christoph Schnörr. A convex relaxation approach to the affine subspace clustering problem. In *Pattern Recognition - 37th German Conference, GCPR 2015, Aachen, Germany, October 7-10, 2015, Proceedings*, pages 67–78, 2015.
- [SRS16] Francesco Silvestri, Gerhard Reinelt, and Christoph Schnörr. Symmetry-free SDP relaxations for affine subspace clustering. *ArXiv e-prints*, July 2016.
- [Sta11] Richard P. Stanley. *Enumerative Combinatorics: Volume 1*. Cambridge University Press, New York, NY, USA, 2nd edition, 2011.
- [Stu11] Raman Sanyal ; Frank Sottile ; Bernd Sturmfels. Orbitopes. *Mathematika : a journal of pure and applied mathematics*, 57(2):275–314, 2011.
- [Sze94] Mario Szegedy. A note on the theta number of Lovász and the generalized delserte bound. In *35th Annual Symposium on Foundations of Computer Science, Santa Fe, New Mexico, USA, 20-22 November 1994*, pages 36–39, 1994.
- [Tro15] Nicolas Trotignon. Perfect graphs: a survey. *arXiv preprint arXiv:1301.5149*, 2015.
- [TTT96] Kim-Chuan Toh, Michael J. Todd, and Reha H. Tütüncü. *SDPT3 — a MATLAB software package for semidefinite programming*, 1996.
- [TTT03] Reha H. Tütüncü, Kim-Chuan Toh, and Michael J. Todd. Solving semidefinite-quadratic-linear programs using SDPT3. *Math. Program.*, 95(2):189–217, 2003.
- [Wol] Philip Wolfe. Finding the nearest point in a polytope. *Math. Program.*, 11(1):128–149.
- [XWI05] Rui Xu and Donald Wunsch II. Survey of Clustering Algorithms. *IEEE Trans. Neural Networks*, 16(3):645–678, May 2005.
- [Zie95] Günter M. Ziegler. *Lectures on polytopes*. Graduate Texts in Mathematics. Springer, New York, 1995.

Index

A		E	
adjacency matrix	6	edge	6
assignment matrix	45	Euclidean	
		clustering	45, 82
		affine	87
		distance	10
B		F	
barycenter	48	face	9
basis	13		
Bell number	43	G	
block-inducing constraints	68	Gröbner basis	14
		reduced	14
C		graph	6
characteristic vector	45	bipartite	7
Cholesky decomposition	11	colouring	45
chromatic number	45, 75	complement	7
clustering		homomorphism	6
k-center	108	isomorphism	6
colouring	75	perfect	75
combinatorial projection matrices	68	subgraph	6
cone	10	Grassmannian	67
conic order	10	group action	46
dual	10		
Lorenz-	11	H	
pointed	10	halfspace	9
positive semidefinite	11	hyperplane	9
proper	10	supporting	9
self-dual	10		
convex		I	
function	9	ideal	13
hull	9		
set	9		
counting function	44		

- independence system 7
- K**
- k -partition
 function 43
 problem 43
- \mathcal{K} -box 26
- \mathcal{K} -zonotope 27
- L**
- lattice 8
- Lovász theta number 75
- LP-Newton method 25
 conic 25
- M**
- Max- k -Cut 43
- method of moments 16, 20
- minimum-cover problem 44
- Minkowski sum 5
- mixed linear regression 103
- moment matrix 16
 combinatorial 24
 localizing 17
 reduced 23
- monomial 12
 order 13
- N**
- node 6
- NP-hard 9
- O**
- orbit 46
- order
 graded lexicographic 13
 lexicographic 13
 monomial 13
 partial 8
- strict 8
- total 8
- well 8
- orthogonal projection 10
- P**
- partition 42
 function 44
 regularized 45
 separable 43
- matrix 55
- orbitope 50
- part 42
- problem 44
- polytope 11
- poset 8
- projection matrix 66
- R**
- real variety 13
- regularizer 45
- Reynolds operator 48, 61
- Riesz functional 13
- S**
- Schur complement 11
- semialgebraic 13
- separated 94
- Sherman-Morrison formula 6
- simplex 11
- simplicial cover 90
- squared error function 82
 affine 86
- stable set 44, 74
 number 74
- Stirling numbers of the second kind ... 43
- subdifferential 10
- subgradient 10
- suitable for CLPN 32

T

total degree	12
trace	5
product	6

U

unit cube	12
-----------------	----

V

Veronese map	13
vertex	12
vertex-transitive	6